# Journal of Research

## of the National Bureau of Standards

### Special Issue: Chemometrics Conference Proceedings

---

---

---

# Topical Issue: Chemometrics

This issue of the NBS *Journal of Research* is devoted entirely to one topic: Chemometrics. A conference by that title held earlier this year at NBS brought together experts in analytical chemistry and applied mathematics, disciplines which are the constituents of this new field. This conference was probably the first one in the United States by that title.

The roots of the interdisciplinary effort go back to the late Dr. William (Jack) Youden and we dedicate this issue to him. A brief description of Youden's career serves as the introduction to the collection of conference papers which we present in this volume of the *Journal.* The authors of this biographical sketch, Drs. Ku and DeVoe, worked very closely with Youden while he was at NBS.

With the publication of the papers presented at this conference we hope to stimulate further work in the field of chemometrics. Special recognition goes to the organizers of the conference who also served as invited editors of this special issue of the NBS *Journal of Research:* Drs. Clifford H. Spiegelman of the Center for Applied Mathematics, Robert L. Watters of the Center for Analytical Chemistry, and Jerome Sacks from the University of Illinois.

Hans J. Oser
Chief Editor

# JACK YOUDEN

We are pleased that the Conference proceedings are dedicated to Jack Youden.

Youden, an analytical chemist turned mathematical statistician, was a simplifying and, indeed, enthralling teacher of statistical principles who could hold the attention of sophisticated statisticians as well as scientists and engineers.

Born at the turn of the century, Youden worked for many years at the Boyce Thompson Institute for Plant Research in Yonkers, NY before joining the staff at NBS in 1948. He remained with the Bureau for 17 years and after his retirement continued this association as a guest worker until his death in 1971.

The transformation of Youden from chemist to statistician probably began with his reading of R. A. Fisher's *Statistical Methods for Research Workers*. He studied under Fisher in 1937–38 at the Galton Laboratory of University College, London, thanks to a Rockefeller Fellowship for the discovery of a new class of incomplete block designs, "Youden Squares," which found immediate application in biological and medical research.

At NBS he introduced "Youden Plots" for interlaboratory tests and "Youden's Ruggedness Test" as a check on test methods, but he distinguished himself principally for his ability to reduce a complicated idea to its essentials and to express that idea in a simple, straightforward manner, so that it became understandable to scientists of all disciplines. He worked hard at stripping away needless detail and jargon. And he generated interest in his lectures, making his subject so intriguing that an hour passed as minutes. He was a very rare breed—an outstanding statistician who understood experimental systems in chemistry, physics, and engineering.

His first book, *Statistical Methods for Chemists*, was published in 1951 and was followed by *Statistical Techniques for Collaborative Tests* which consisted of two

Youden lecture series, "Accuracy of Analytical Procedures" and "The Collaborative Test." The latter book was published by the Association of Official Analytical Chemists (AOAC) whose William Horwitz states: "These lectures probably have had a greater influence on improving the quality and interpretation of collaborative studies conducted by AOAC members than any other event in the (85-year) existence of the Association."

A number of Youden's most important papers were collected in the *Journal of Quality Technology* (Vol. 4, No. 1, January 1972) and in Volume 1 of NBS Special Publication 300, *Precision Measurement and Calibration: Statistical Concepts and Procedures*.

Jack Youden was a chemist and a communicator. The Chemical Division of the American Society for Quality Control in 1969 established a Jack Youden prize to be awarded yearly for the best expository paper in its journal, *Technometrics*. But it was Youden the statistician who furthered collaboration and helped to maximize the information content of experimentation, which is what the Chemometrics Conference was about. So it is appropriate that these conference proceedings be dedicated to the memory of Dr. Youden.

## H. H. Ku

**Statistical Engineering Division**
**National Bureau of Standards**

## J. R. DeVoe

**Inorganic Analytical Research Division**
**National Bureau of Standards**

# The Organizers' Goals

The wide range of disciplines represented by the participants and attendees of the Chemometrics Research Conference held at the Gaithersburg Holiday Inn on May 20–22, 1985, exemplifies the depth and diversity of the chemometrics community. The Conference was sponsored by several important professional societies whose members are involved in chemometric activity. These include the Analytical Division of the American Chemical Society, the Section on Physical and Engineering Sciences of the American Statistical Association, the Institute of Mathematical Statistics, and the Society for Applied Spectroscopy. Generous funding for the Conference was provided by the National Bureau of Standards, the Office of Naval Research, and the National Science Foundation.

As organizers, we had two main goals in mind when deciding on the form and substance of the Conference. The first was to provide a forum for reporting on some of the most recent and important research activities in diverse areas relating to chemometrics. Nineteen invited speakers covered topics including experimental design and optimization, kinetic rate constants, Kalman filtering, chromatography, data analysis, artificial intelligence, stochastic processes, regression and factor analysis. Despite the full schedule of papers, each of the five sessions was attended by nearly all of the 134 Conference registrants. Most of the papers are published in this special issue of the National Bureau of Standards Journal of Research to provide the registrants and others the opportunity for careful study of the presentations. In these respects, our first goal was easier to achieve than the second.

Our second and more important goal can only be achieved gradually. This was to increase the willingness of chemists, statisticians, and probabilists to meet as colleagues and to solve problems as a team. This will necessarily involve the exercise of communication skills as well as the combining of scientific skills. We believe that even the best separate efforts of chemists and mathematicians fall far short of the achievements that are possible by joint efforts in chemometric research teams. We underscored this aspect of teamwork between the disciplines of chemistry and mathematics by having each invited paper in one discipline discussed by an invited discussant of the other.

We invited the speakers to take any approach they desired in the exposition of their subject. Hence, the papers range from strictly technical to philosophical in tone. Discussants also had the option of either commenting on the specifics of a given paper, or exploring the relevance of the

subject to their respective disciplines. This free format encouraged open discussion and exchange of ideas at the Conference, and we hope that the same stimulus will be provided by these Proceedings.

Finally, we look forward to future chemometrics conferences organized by researchers with perspectives other than our own. For no matter how broad the coverage of chemometrics topics in a given program, important areas are omitted. We believe that the widest scope of chemometrics research activities can be presented in meetings organized by committees of different backgrounds and insights.

**Clifford H. Spiegelman**
National Bureau of Standards
**Robert L. Watters, Jr.**
National Bureau of Standards
**Jerome Sacks**
University of Illinois, Urbana

# Agenda for Chemometricians

It is most appropriate that the proceedings of this conference are going to be dedicated to the memory of Jack Youden. He was interested in many of the topics that are being considered at this conference, for example, interlaboratory comparisons, calibration, analytical methods, and measurement errors—both systematic and random. He was indeed a pioneering chemometrician, before the name existed. He was also interested in explaining to chemists, chemical engineers, and others how they could benefit by using statistical methods.

I'm sure Youden would have been pleased with this conference, which provides a forum for chemists, statisticians, and others interested in chemometrics to discuss research of mutual interest. He also might have observed that chemometrics as a field has reached a level of maturity that warrants consideration of questions related to spreading the word to others, to non-chemometricians, so that they could take advantage of the techniques that are now available. In other words, perhaps chemometrics as discipline has reached a sufficiently advanced stage of research and development that questions of production should now be addressed. What are our most useful products? Who are out customers? Which products would they find most valuable? What are the obstacles that prevent these customers from using these products now? How can these obstacles be overcome? What are the most important things that can be done in the next three years to reach new customers? What should the agenda be for chemometricians in the next few years?

There are two ways to learn. One is to listen, as in a lecture. The other is to engage in a dialogue, as in a conversation. The first way is passive. The second is active. Let's try the second way to learn from one another how we might answer these questions.

[Participants at this point wrote out answers to these questions, discussed them, and voted on them. The top vote-getters for the most important things that can be done in the next three years to reach new customers were the following, listed in order of decreasing number votes:

1. Organize joint conferences with chemists.
2. Write textbooks on chemometrics.
3. Conduct workshops and teach short courses.
4. Write user-friendly software.
5. Teach chemometrics to graduate students.

6. Write tutorial, expository, and review articles.
7. Undertake joint research projects with chemists.
8. Publicize success stories.
9. Teach chemometrics to undergraduate students.
10. Communicate with management.
11. Hire professionals to help with a public relations effort.
12. Teach chemometrics to high school students.]

I recommend that we take action on the basis of this list. Let me now make a few observations in closing. I would like to suggest a different starting point for statistics courses. Let us represent the relationship between an observed response y and variables $x_1, x_2, \ldots$ as

$$y = f(x_1, x_2, x_3, x_4, x_5, \ldots, x_{126}, \ldots) \ .$$

Many, many, many variables affect y. It is the fluctuation of these variables that gives us different answers when we repeat an experiment two or more times under "identical conditions." We are often interested in creating a mathematical equation (model) that involves a subset of the variables. For purposes of illustration, suppose this subset is $(x_1, x_2)$. We can then write

$$y = f(x_1, x_2) + g(x_1, x_2, x_3, x_4, x_5, \ldots, x_{126}, \ldots) \ .$$

Note that the g function includes $x_1$ and $x_2$ (because of lack of fit of the model) as well as all the other $x$'s. Lack of fit occurs, for example, because the model f may be taken to be linear in $x_1$ and $x_2$ but the actual relationship may be nonlinear in $x_1$ and $x_2$. The function g is most often called experimental error, and it is almost as often endowed by writers with an abundance of desirable and well-known properties. They call it a random variable. A sequence of these experimental errors, they frequently say, can be assumed to be independent, identically distributed according to a Normal distribution with a zero mean and constant variance. I believe that statisticians too readily make this assumption and others like it. *Sometimes* such an assumption makes sense, sometimes not. We should be more careful on this point.

An adequate model is a function that will turn data into white noise, as George Box has said. An analogy that I find useful involves a process for separating gold particles from a slurry. If the process is fully efficient, the waste stream will contain no gold. It is therefore prudent to check the waste stream to see if it contains any gold. Likewise in creating and fitting models, it makes sense to examine residuals to see if they contain any information. The data contain information (that's the gold we want to get), and a good model will extract all the information in those data. Hence the residuals will be manifestations of white noise, an informationless sequence of values.

Chemists and chemical engineers could benefit from knowing more about variance components, statistical graphics, and quality control techniques (including Shewhart and cumulative sum charts). But, above all, I think they would find statistical experimental designs to be the most useful thing of all that chemometricians have to offer. Such designs provide a practical means for increasing research efficiency, which might be defined as the amount of information one obtains per dollar spent.

The damage done by poor experimental design is irreparable. A poor design results in data that contain little information. Consequently, no matter how thorough, how clever, or how sophisticated the subsequent analysis is, little information can be extracted. A good design, for the same expenditure of time, money, and other resources, results in data rich in information. A fruitful analysis is then possible. (Note that analysis is defined as trying to extract all the useful information in the data.)

Two-level factorial and fractional factorial designs can be extremely useful for chemists, chemical engineers, and others who do similar work. One of the best ways for a student to learn about such designs is to set one up, get the data, analyze them, and interpret the results. For a number of years I have had students in our experimental design course undertake such projects.

The main piece of advice I give them is to work on something they care about, something they are really interested in.

Toward the end of an introductory one-semester undergraduate course in statistics, for example, one student said that he was a pilot and that, ever since he started to fly, he had asked instructors and other pilots what he should do if the engine failed on takeoff. He had been told by several people that he should bank the plane, go into a 180° turn, and land on the runway from which he took off. Unfortunately, many different ways of doing this maneuver had been suggested. He successfully organized and executed a replicated $2^3$ factorial design with three variables: bank angle, flap angle, and speed. He measured the loss in altitude. He started each test at 1000 feet instead of ground level. The experiment was a success. He learned which combination of factors he should use for his plane, and he discovered the minimum altitude for attempting such a maneuver.

Factorial designs can be understood and run with profit by graduate, undergraduate, senior high school, and junior high school students. Maybe younger students can use them, too. Students can study the baking of cakes, the riding of bicycles, the making of chemicals, the growing of plants, and the swinging of pendulums. Dalia Sredni, when she was a seventh grader, for instance, studied the effects of changing oven temperature, baking time, and the amount of baking soda when making a cake. Students should be told about factorial designs early so that they can study systems that depend on many variables and learn how they work. Using such designs they can discover interesting things, have fun, and be surprised. Our students deserve more of these pleasures. I have included a list of 101 experiments that have been done by students at Wisconsin, to indicate the variety of things that is possible.

I would like to end by congratulating the conference organizers for the excellent job they have done. It is clear that they have worked hard to make things enjoyable and rewarding for those of us who have been fortunate enough to participate.

## William G. Hunter

**Professor of Statistics and Industrial Engineering**
**Director of Center for Quality and Productivity Improvement**
**University of Wisconsin—Madison**

**Table 1.** List of some studies done by students in an experimental design course at the University of Wisconsin—Madison.

| variables | responses |
| --- | --- |
| 1. seat height (26, 30 inches), generator (off, on), tire pressure (40, 55 psi) | time to complete fixed course on bicycle and pulse rate at finish |
| 2. brand of popcorn (ordinary, gourmet), size of batch (1/3, 2/3 cup), popcorn to oil ratio (low, high) | yield of popcorn |
| 3. amount of yeast, amount of sugar, liquid (milk, water), rise temperature, rise time | quality of bread, especially the total rise |
| 4. number of pills, amount of cough syrup, use of vaporizer | how well twins, who had colds, slept during the night |
| 5. speed of film, light (normal, diffused), shutter speed | quality of slides made close up with flash attachment on camera |
| 6. hours of illumination, water temperature, specific gravity of water | growth rate of algae in salt water aquarium |
| 7. temperature, amount of sugar, food prior to drink (water, salted popcorn) | taste of Koolaid |
| 8. direction in which radio is facing, antenna angle, antenna slant | strength of radio singal from particular AM station in Chicago |
| 9. blending speed, amount of water, temperature of water, soaking time before blending | blending time for soy beans |

Table 1. continued

| variables | responses |
|---|---|
| 10. charge time, digits fixed, number of calculations performed | operation time for pocket calculator |
| 11. clothes dryer (A, B), temperature setting, load | time until dryer stops |
| 12. pan (aluminum, iron), burner on stove, cover for pan (no, yes) | time to boil water |
| 13. aspirin buffered? (no, yes), dose, water temperature | hours of relief from migraine headache |
| 14. amount of milk powder added to milk, heating temperature, incubation temperature | taste comparison of homemade yogurt and commercial brand |
| 15. pack on back (no, yes), footwear (tennis shoes, boots), run (7, 14 flights of steps) | time required to run up steps and heartbeat at top |
| 16. width to height ratio of sheet of balsa wood, slant angle, dihedral angle, weight added, thickness of wood | length of flight of model airplane |
| 17. level of coffee in cup, devices (nothing, spoon placed across top of cup facing up), speed of walking | how much coffee spilled while walking |
| 18. type of stitch, yarn guage, needle size | cost of knitting scarf, dollars per square foot |
| 19. type of drink (beer, rum), number of drinks, rate of drinking, hours after last meal | time to get steel ball through a maze |
| 20. size of order, time of day, sex of server | cost of order of french fries, in cents per ounce |
| 21. brand of gasoline, driving speed, temperature | gas mileage for car |
| 22. stamp (first class, air mail), zip code (used, not used), time of day when letter mailed | number of days required for letter to be delivered to another city |
| 23. side of face (left, right), beard history (shaved once in two years—sideburns, shaved over 600 times in two years—just below sideburns) | length of whiskers 3 days after shaving |
| 24. eyes used (both, right), location of observer, distance | number of times (out of 15) that correct gender of passerby was determined by experimenter with poor eyesight wearing no glasses |
| 25. distance to target, guns (A, B), powders (C, D) | number of shot that penetrated a one foot diameter circle on the target |
| 26. oven temperature, length of heating, amount of water | height of cake |
| 27. strength of developer, temperature, degree of agitation | density of photographic film |
| 28. brand of rubber band, size, temperature | length of rubber band before it broke |
| 29. viscosity of oil, type of pick-up shoes, number of teeth in gear | speed of H.O. scale slot racers |
| 30. type of tire, brand of gas, driver (A, B) | time for car to cover one-quarter mile |
| 31. temperature, stirring rate, amount of solvent | time to dissolve table salt |
| 32. amounts of cooking wine, oyster sauce, sesame oil | taste of stewed chicken |
| 33. type of surface, object (slide rule, ruler, silver dollar), pushed? (no, yes) | angle necessary to make object slide |
| 34. ambient temperature, choke setting, number of charges | number of kicks necessary to start motorcycle |
| 35. temperature, location in oven, biscuits covered while baking? (no, yes) | time to bake biscuits |
| 36. temperature of water, amount of grease, amount of water conditioner | quantity of suds produced in kitchen blender |
| 37. person putting daughter to bed (mother, father), bed time, place (home, grandparents) | toys child chose to sleep with |
| 38. amount of light in room, type of music played, volume | correct answers on simple arithmetic test, time required to complete test, words remembered (from list of 15) |
| 39. amounts of added Turkish, Latakia, and Perique tobaccos | bite, smoking characteristics, aroma, and taste of tobacco mixture |
| 40. temperature, humidity, rock salt | time to melt ice |
| 41. number of cards dealt at one time, position of picker relative to the dealer | points in games of sheepshead, a card game |
| 42. marijuana (no, yes), tequilla (no, yes), sauna (no, yes) | pleasure experienced in subsequent sexual intercourse |

Table 1. continued

| variables | responses |
|---|---|
| 43. amounts of flour, eggs, milk | taste of pancakes, consensus of group of four living together |
| 44. brand of suntan lotion, altitude, skier | time to get sunburned |
| 45. amount of sleep the night before, substantial exercise during the day? (no, yes), eat right before going to bed? (no, yes) | soundness of sleep, average reading from 5 persons |
| 46. brand of tape deck used for playing music, bass level, treble level, synthesizer? (no, yes) | clearness and quality of sound, and absence of noise |
| 47. Type of filter paper, beverage to be filtered, volume of beverage | time to filter |
| 48. type of ski, temperature, type of wax | time to go down ski slope |
| 49. ambient temperature for dough when rising, amount of vegetable oil, number of onions | four quality characteristics of pizza |
| 50. amount of fertilizer, location of seeds (3×3 Latin square) | time for seeds to germinate |
| 51. speed of kitchen blender, batch size of malt, blending time | quality of ground malt for brewing beer |
| 52. soft drink (A, B), container (can, bottle), sugar free? (no, yes) | taste of drink from paper cup |
| 53. child's weight (13, 22 pounds), spring tension (4, 8 cranks), swing orientation (level, tilted) | number of swings and duration of these swings obtained from an automatic infant swing |
| 54. orientation of football, kick (ordinary, soccer style), steps taken before kick, shoe (soft, hard) | distance football was kicked |
| 55. weight of bowling ball, spin, bowling lane (A, B) | bowling pins knocked down |
| 56. distance from basket, type of shot, location on floor | number of shots made (out of 10) with basketball |
| 57. temperature, position of glass when pouring soft drink, amount of sugar added | amount of foam produced when pouring soft drink into glass |
| 58. brand of epoxy glue, ratio of hardener to resin, thickness of application, smoothness of surface, curing time | strength of bond between two strips of aluminum |
| 59. amount of plant hormone, water (direct from tap, stood out for 24 hours), window in which plant was put | root lengths of cuttings from purple passion vine after 21 days |
| 60. amount of detergent (1/4, 1/2 cup), bleach (none, 1 cup), fabric softener (not used, used) | ability to remove oil and grape juice stains |
| 61. skin thickness, water temperature, amount of salt | time to cook chinese meat dumpling |
| 62. appearance (with and without a crutch), location, time | time to get a ride hitchhiking and number of cars that passed before getting a ride |
| 63. frequency of watering plants, use of plant food (no, yes), temperature of water | growth rate of house plants |
| 64. plunger A up (slow, fast), plunger A down (slow, fast), plunger B up (slow, fast) plunger B down (slow, fast) | reproducibility of automatic dilutor, optical density readings made with spectrophotometer |
| 65. temperature of gas chromatograph column, tube type (U, J), voltage | size of unwanted droplet |
| 66. temperature, gas pressure, welding speed | strength of polypropylene weld, manual operation |
| 67. concentration of lysozyme, pH, ionic strength, temperature | rate of chemical reaction |
| 68. anhydrous barium peroxide powder, sulfur, charcoal dust | length of time fuse powder burned and the evenness of burning |
| 69. air velocity, air temperature, rice bed depth | time to dry wild rice |
| 70. concentration of lactose crystal, crystal size, rate of agitation | spreadability of caramel candy |
| 71. positions of coating chamber, distribution plate, and lower chamber | number of particles caught in a fluidized bed collector |
| 72. proportional band, manual reset, regulator pressure | sensitivity of a pneumatic valve control system for a heat exchanger |
| 73. chloride concentration, phase ratio, total amine concentration, amount of preservative added | degree of separation of zinc from copper accomplished by extraction |

**Table I.** continued

| variables | responses |
|---|---|
| 74. temperature, nitrate concentration, amount of preservative added | measured nitrate concentration in sewage, comparison of three different methods |
| 75. solar radiation collector size, ratio of storage capacity to collector size, extent of short-term intermittency of radiation, average daily radiation on three successive days | efficiency of solar space-heating system, a computer simulation |
| 76. pH, dissolved oxygen content of water, temperature | extent of corrosion of iron |
| 77. amount of sulfuric acid, time of shaking milk-acid mixture, time of final tempering | measurement of butterfat content of milk |
| 78. mode (batch, time-sharing), job size, system utilization (low, high) | time to complete job on computer |
| 79. flow rate of carrier gas, polarity of stationary liquid phase, temperature | two different measures of efficiency of operation of gas chromatograph |
| 80. pH of assay buffer, incubation time, concentration of binder | measured cortisol level in human blood plasma |
| 81. aluminum, boron, cooling time | extent of rock candy fracture of cast steel |
| 82. magnification, read out system (micrometer, electronic), stage light | measurement of angle with photogrammetric instrument |
| 83. riser height, mold hardness, carbon equivalent | changes in height, width, and length dimensions of cast metal |
| 84. amperage, contact tube height, travel speed, edge preparation | quality of weld made by submerged arc welding process |
| 85. time, amount of magnesium oxide, amount of alloy | recover of material by steam distillation |
| 86. pH, depth, time | final moisture content of alfalfa protein |
| 87. deodorant, concentration of chemical, incubation time | odor produced by material isolated from decaying manure, after treatment |
| 88. temperature variation, concentration of cupric sulfate concentration of sulfuric acid | limiting currents on totaling disk electrode |
| 89. air flow, diameter of bead, heat shield (no, yes) | measured temperature of a heated plate |
| 90. voltage, warm-up procedure, bulb age | sensitivity of microdensitometer |
| 91. pressure, amount of ferric chloride added, amount of lime added | efficiency of vacuum filtration of sludge |
| 92. longitudinal feed rate, transverse feed rate, depth of cut | longitudinal and thrust forces for surface grinding operation |
| 93. time between preparation of sample and refluxing, reflux time, time between end of reflux and start of titrating | chemical oxygen demand of samples with same amount of waste (acetanilide) |
| 94. speed of rotation, thrust load, method of lubrication | torque of taper roller bearings |
| 95. type of activated carbon, amount of carbon, pH | adsorption characteristics of activated carbon used with municipal waste water |
| 96. amounts of nickel, manganese, carbon | impact strength of steel alloy |
| 97. form (broth, gravy), added broth (no, yes), added fat (no, yes), type of meat (lamb, beef) | percentage of panelists correctly identifying which samples were lamb |
| 98. well (A, B), depth of probe, method of analysis (peak height, planimeter) | methane concentration in completed sanitary landfill |
| 99. paste (A, B), preparation of skin (no, yes), site (sternum, forearm) | electrocardiogram reading |
| 100. lime dosage, time of flocculation, mixing speed | removal of turbidity and hardness from water |
| 101. temperature difference between surface and bottom waters, thickness of surface layer, jet distance to thermocline, velocity of jet, temperature difference between jet and bottom waters | mixing time for an initially thermally stratified tank of water |

# Adaptive Kalman Filtering

## Steven D. Brown

Washington State University, Pullman, WA 99164

and

## Sarah C. Rutan

Virginia Commonwealth University, Richmond, VA 23284

The increased power of small computers makes the use of parameter estimation methods attractive. Such methods have a number of uses in analytical chemistry. When valid models are available, many methods work well, but when models used in the estimation are in error, most methods fail. Methods based on the Kalman filter, a linear recursive estimator, may be modified to perform parameter estimation with erroneous models. Modifications to the filter involve allowing the filter to adapt the measurement model to the experimental data through matching the theoretical and observed covariance of the filter innovations sequence. The adaptive filtering methods that result have a number of applications in analytical chemistry.

Key words: automated covariance estimation; Kalman filter; multicomponent analysis.

## 1. Introduction

The increased computational power available from small computers has prompted a re-evaluation of the methods used in reducing data obtained from a chemical analysis. Many of the responses obtained from chemical analyses are suited to mathematical analysis by methods which estimate the parameters that generate the response; these parameters are generally concentrations. For parameter estimation to be successful, an accurate model of the behavior of the chemical system is necessary. The model used need not be theoretical; empirical models based on experimental results or on a numerical simulation of the chemical system are often satisfactory

as well. When valid models are available, the parameters associated with the model may be obtained with a variety of methods. Some that have seen extensive use in analytical chemistry include analysis of the chemical data using linear least squares [1], nonlinear least squares analysis [2,3], and Kalman filtering [4–6].[1]

The methods mentioned above all work well with accurate models, but are much less satisfactory when used with models containing errors that can arise from many sources. Theoretical models, or models based on simulation, may not describe the physics or chemistry of a system well enough to predict system responses to the accuracy desired. Small changes in the experimental conditions used for data acquisition may perturb experimentally-obtained models, leading to errors when these models are used to analyze subsequent experiments. And, it may be impossible, because of the effects of chemical equilibria, to obtain independent responses for some of the chemical species included in a model for a complex system, leading to "chemical" model errors.

Relatively few methods have been developed to compensate for model errors affecting multicomponent quantitation. Approaches using factor analysis [7] have

---

[1] Bracketed figures indicate literature references.

been developed for situations where the model is unknown, but these approaches are generally limited to very few components [8], and it is difficult to incorporate additional *a priori* information into such methods. An alternative approach has used the Kalman filter. The Kalman filter is a linear, recursive estimator which yields optimal estimates for parameters associated with a valid model [9,10]. Several methods, classified under the term "adaptive filtering," have been developed to permit the filter to produce accurate parameter estimates in the presence of model errors [11–15]. This paper summarizes the development of an adaptive Kalman filter for use in the mathematical analysis of overlapped multicomponent chemical responses.

## 2. Theory

**Kalman Filtering.** The Kalman filter has received some attention for the analysis of multicomponent chemical responses [4,6,16,17]. Because most models relating chemical responses to concentrations are linear, application of the Kalman filter is straightforward. The filter model is comprised of two equations. The system model, which describes the time evolution of the desired parameters, is, in state-space notation

$$X(k) = F(k,k-1)X(k-1) + w(k) \qquad (2.1)$$

where $X$ is a $n \times 1$ column vector of state variables describing the chemical system, where $F$ is an $n \times n$ matrix describing how the states change with time, $w$ is a vector describing noise contributions to the system model, and where $k$ indicates time or some other independent variable which meets the noise requirements given below. For state-invariant systems, $F$ reduces to the identity matrix $I$. Because multicomponent analysis is most often performed under conditions where concentrations are constant over the time frame involved, the case where $X$ is time-invariant is considered here.

The second equation describes the measurement process by relating the measured response $z(k)$, to the filter states. For a single sensor, the measurement model is given by

$$z(k) = H^T(k)X(k) + v(k) \qquad (2.2)$$

where $H^T(k)$ is a $1 \times n$ vector relating the response at point $k$ to the $n$ states, and the scalar $v(k)$ is the noise contribution of the measurement process. For example, in absorption spectrophotometry, $z(k)$ is an absorbance measurement at some wavelength $k$, and $H^T(k)$ is the vector of absorption coefficients at that wavelength for all chemical species included in the model. The mea-

surement model is easily extended to systems with multiple sensors.

The two noise processes in the Kalman filter, $w(k)$ and $v(k)$, are usually assumed to be independent, zero-mean, white noise processes. The matrix $Q(k)$, defined as the covariance of the noise in the system model, is taken as approximately zero for the time invariant system discussed in this paper. The scalar quantity $R(k)$ is the variance of the noise in the measurement process.

The Potter-Schmidt square-root algorithm, one implementation of the Kalman filter [18], is given in table 1. The details of this algorithm have been discussed elsewhere [18,19]. Initial guesses for the filter states and for the covariance matrix $P$ are required to start the filter. Estimates of $X$ and $P$ depend on $k$, and because both are projected ahead of the data (in eqs 2.3 and 2.4) by the filter, the notation $(j \mid k)$ is used to indicate that the estimate is made at point $j$, based on data obtained up through point $k$. The filter output consists of estimates $\hat{X}$, as well as $\hat{P}$. In analytical chemistry, these are often estimates of concentrations and of the error in the concentrations.

**Table 1.** Algorithm equations for the square root Kalman filter.

**State estimate extrapolation**

$$X(k \mid k-1) = F(k \mid k-1) \cdot X(k-1 \mid k-1) \qquad (2.3)$$

**Covariance square root extrapolation**

$$S(k \mid k-1) = F(k,k-1) \cdot S(k-1 \mid k-1) \cdot F^T(k,k-1) \qquad (2.4)$$

where

$$F(k,k-1) = I$$

$$P = S \cdot S^T$$

**Kalman gain:**

$$K(k) = a \cdot S(k \mid k-1) \cdot G(k) \qquad (2.5)$$

where

$$G(k) = S^T(k \mid k-1) \cdot H(k) \qquad (2.6)$$

$$1/a = G^T(k) \cdot G(k) + R(k) \qquad (2.7)$$

$$d = (1 + (a \cdot R(k))^{1/2})^{-1} \qquad (2.8)$$

**State estimate update:**

$$X(k \mid k) = X(k \mid k-1) + K(k)[z(k) - H^T(k) \cdot X(k \mid k-1)] \qquad (2.9)$$

**Covariance square root update:**

$$S(k \mid k) = S(k \mid k-1) - ad \cdot S(k \mid k-1)G(k) \cdot G^T(k) \qquad (2.10)$$

**Adaptive Kalman Filtering.** Errors can occur in both of the models used in the Kalman filter. Errors in the system model arise if the system was taken as time-invariant, but was actually composed of time-dependent states. Errors in the measurement model arise from underestimating the number of components involved in

the state vector (which can be thought of as incorrectly setting values in $H^T(k)$ to zero for one of the possible elements in the state vector), or by use of inaccurate values in $H^T(k)$. Either type of error produces a suboptimal filter, in that the accuracy of the filter's estimates are severely degraded. Many methods for compensating these model errors make use of the filter innovations sequence, $\nu(k)$, defined as

$$\nu(k) = z(k) - H^T(k)\hat{X}(k \mid k-1). \qquad (2.11)$$

The innovations sequence can be used to construct a measure of the optimality of the filter; a necessary and sufficient condition for an optimal filter is that this sequence be a white noise process [10]. An optimal filter is one that minimizes the mean square estimation error $E\{(X-\hat{X})(X-\hat{X})^T\}$. A suboptimal filter may generate results which show large estimation errors, or even a divergence of the errors [11]. The aim of an adaptive filter is to reduce or bound these errors by modifying, or adapting, the models used in the Kalman filter to the real data.

Several methods for controlling error divergence in the filter have been reported [11-15]. Most involve cases where $Q$ is poorly known, the situation which arises when the time-dependence of states is incorrectly modeled. These include methods based on Bayesian estimation and maximum likelihood estimation [14], correlation methods [14], and covariance matching techniques [14,15]. The last method has also been suggested for use when $Q$ is known, but $R$ is unknown, the situation that arises when the number of components in the state is underestimated, or when the measurement model is otherwise incorrect. Because errors in the number of components and in the response factors used in the measurement model are common in multicomponent chemical analysis, covariance matching is used to develop the filter discussed here.

The aim of covariance matching is to insure that the residuals remain consistent with the theoretical covariances. The covariance of the innovations sequence $\nu(k)$ is [14]

$$E[\nu(k) \cdot \nu(k)^T] = H^T(k)P(k \mid k-1)H(k) + R(k). \qquad (2.12)$$

If the actual covariance of $\nu(k)$ is much larger than the covariance obtained from the Kalman filter, either $Q$ or $R$ should be increased to prevent divergence. In either case, this has the effect of increasing $P(k \mid k-1)$, thus bringing the actual covariance of $\nu(k)$ closer to that given in eq 2.12. This also has the effect of decreasing the filter gain matrix, $K$, thereby "closing" the filter to new data which would otherwise be incorrectly interpreted because of errors in the measurement model. In

essence, this amounts to "covering" the errors in the model with noise, then estimating the noise variance. The adaptive estimate of $R$ at the $k$th point, when $Q$ is known, is

$$R(k) = 1/m \left[ \sum_{j=1}^{m} \nu(k-j) \cdot \nu(k-j) \right]$$

$$- H^T(k)S(k \mid k-1)S^T(k \mid k-1)H(k) \qquad (2.13)$$

where $m$ is the width of an empirically chosen rectangular smoothing window for the innovations sequence. The smoothing operation improves the statistical significance of the estimator for $R(k)$, as it now depends on many residuals.

Adaptive estimation of $R$ allows accurate estimates for the states to be obtained, even in the presence of model errors, because only data for which an accurate model is available are used in the filter. A new measurement model can be constructed from the estimated $R(k)$, either by augmenting the $H^T$ vector, or by correcting any one of its existing elements; choice of correction or augmentation is arbitrary. For augmentation, the equations

$$H^*_{n+1}(k) = b(k)[R(k+m)/2)]^{1/2}, \text{ for } b(k) > 0 \qquad (2.14)$$

$$H^*_{n+1}(k) = 0, \text{ for } b(k) < 0 \qquad (2.15)$$

apply, where $H^*_{n+1}(k)$ denotes an element which is incorporated in the $H^T$ vector. The term $(k+m/2)$ arises from the lag induced by averaging $m$ of the squared innovations. The factor $b(k)$ is defined as

$$b(k) = 1, \text{ for } \sum_{j=1}^{m} \nu(k-j+m/2)/m > 0 \qquad (2.16)$$

$$b(k) = -1, \text{ for } \sum_{j=1}^{m} \nu(k-j+m/2)m < 0. \qquad (2.17)$$

Equations 2.16 and 2.17 allow determination of the sign of the model error by evaluating the average of the innovations over the range for which $R$ was calculated. Equation 2.15 reflects the fact that the relation between the chemical response and concentration, given by $H^T$, is generally positive.

For correction of the $i$th component of the vector $H^T$, the expressions

$$H^*_i(k) = H_i(k) + b(k)[R(k+m/2)]^{1/2}, \qquad \text{for}$$

$$H^*_i(k) > 0 \qquad (2.18)$$

$$H^*_i(k) = 0, \text{ for } H^*_i(k) < 0 \qquad (2.19)$$

apply instead of those given in eqs 2.14 and 2.15. In either case, a valid measurement model can be generated from the adaptive estimation of $R$.

Two criteria must be met for this adaptive filter to be useful in the mathematical analysis of multicomponent responses. First, the model to be adaptively corrected must already be correct for some of the values of $k$ where each of the known components of the model has a measureable response. The second requirement is that the adaptive correction must be performed on a single component. For a single sensor, $R$ is a scalar, and it is not possible to distinguish the different portions belonging to the different components. It is often feasible, however, to treat model errors as a single, unmodeled component without affecting the accuracy of some or all of the estimated quantities. Although it has been observed that the adaptive estimation of $R$ by covariance matching is not a sufficient condition for obtaining an improved measurement model [10], application of this approach in the mathematical analysis of multicomponent responses has shown that significant model improvement generally occurs in practice [15,20].

**Automation of the Adaptive Filter.** The adaptive filter requires an initial guess of the states and of their covariances, just as in the ordinary filter. The adaptive estimation of $R$ affects the calculation of $P$, however, and it is found that the diagonal elements of $P$ decrease as $R$ decreases. Since the size of $R$ is directly related to the quality of the measurement model, this relation provides a means by which the quality of the final filter estimates can be judged. Once results are obtained with minimum values for the diagonal elements of the estimated $P$, the resulting corrected measurement model better describes the experimental data available, judging from the deterministic variances of fitting before and after the model correction. Because the innovations are not white in the presence of model error, the filter results are no longer guaranteed to be optimal, but now depend on the initial guess. Thus, the adaptive filter must be run several times, with different initial guesses, $X_o$ and $P_o$, to locate those best estimates. This process is easily automated, however. Simplex optimization [21–23] can be used to minimize the metric based on the diagonal elements of the covariance matrix

$$Y = \sum_{i=1}^{n} \log P_{ii} \qquad (2.20)$$

as a function of the initial guesses input to the adaptive filter. We have previously demonstrated that the minima in the variance surface $Y = f(X_o, P_o)$ correspond well to the minima in an error surface defined by the quantities $(X - \hat{X})$ [24].

## 3. Application in Analytical Chemistry

**Empirical Model Improvement.** Empirical models have been used with the Kalman filter to study the chemical speciation of metal ions. One study [20] reported the adaptive correction of the visible photoacoustic spectrum of $Pr(EDTA)^-$. This spectrum was obtained from data collected on solutions containing both $Pr^{3+}$ and $Pr(EDTA)^-$ species. Direct spectroscopic measurement of $Pr(EDTA)^-$ is not simple. A similar approach was also used to obtain the spectrum of $(UO_2)_3(OH)_5^+$, another ion whose spectrum is difficult to observe in the absence of related chemical species [25]. These studies demonstrate the ability of adaptive filtering to correct for "chemical" errors in the measurement model.

Two other studies used adaptive filtering to model the electrochemical response of an equilibrium mixture of $Cd^{2+}$ and $Cd(NTA)^-$ [20,26]. The adaptively modeled component, attributed to the reduction of $Cd^{2+}$ after dissociation of the $Cd(NTA)^-$ complex, was corrected [26] from an approximate model based on digital simulation [19]. The stability constant for the $Cd(NTA)^-$ species was estimated from the concentrations obtained from the filter. These studies illustrate the correction of "theoretical" errors in the measurement model by adaptive filtering.

The adaptive filter has also been used to correct empirical models for errors which occurred in data acquisition. An example is the correction of models used for the resolution of overlapped electrochemical responses. Resolved peaks are generally needed to obtain estimates of the component concentration. Small changes in experimental conditions, occurring between the time when data are obtained for use in empirical models and the time when the mixtures are measured, change peak positions slightly. The resulting inaccuracy in the model degrades the accuracy of the resolution obtained with the Kalman filter. Adaptive filtering can correct for these type of model errors, resulting in substantially improved concentration estimates from multicomponent electrochemical responses [20].

**Removal of Interferences.** In many multicomponent analyses, substances which interfere with the chemical analysis are often present. Frequently, these species must be chemically separated, because they are not easily removed in the mathematical analysis of the data. Adaptive estimation of these unknown components of the model is an alternative approach. The feasibility of this has been demonstrated [24] in a visible spectrophotometric analysis, where adaptive filtering was used to quantify $UO_2^{2+}$, $Ni^{2+}$, $Co^{2+}$ and picric acid in the presence of the "unknown" contaminant $Cu^{2+}$. The errors in estimating species concentrations were typically less than 5%. An adaptive estimation of $Co^{2+}$ in the presence of "unknown" $Cu^{2+}$, $Ni^{2+}$, $UO_2^{2+}$ and picric acid, where interferent species responses strongly overlap that for the species of interest, gave an estimation

error of 14% with a five-fold excess of interferent species. This estimation's lower accuracy results from the adaptive filter's response when its model restrictions are not met, a situation which occurs here as a consequence of the severe overlap of the analyte and interferent responses. Even though this result is of lower accuracy than many of the others reported, it is still remarkable. Unlike the other fitting, this result does not rely on the use of a complete model. Using peak resolution based on an ordinary filtering approach, with the same incomplete measurement model, an error of 200–300% is likely.

## 4. Conclusion

The automated, adaptive estimation of measurement model covariance permits the application of Kalman filtering in chemical systems where models are poorly known. Although results obtained from the adaptive filter are not guaranteed to be optimal by theory, significant improvement in the accuracy of models and estimated parameters is generally possible in practice.

Restrictions are fairly minor: parts of the model must be known well enough to "open" the filter to the data, and only one component of the model may be adaptively corrected at a time. Adaptive filtering should yield results similar to those obtained from factor analysis using target transformation [27], but the adaptive filter requires only one mixture response, while factor analysis requires several.

## 5. References

[1] Brubaker, T. A.; R. Tracy and C. L. Pomernacki, Linear parameter estimation, Anal. Chem. 50, 1017A–1024A (1978).

[2] Meites, L.; Some new techniques for the analysis and interpretation of chemical data, CRC Crit. Rev. Anal. Chem. 8, 1–53 (1979).

[3] Brubaker, T. A., and K. R. O'Keefe, Nonlinear parameter estimation, Anal. Chem. 51, 1385A–1388A (1979).

[4] Brubaker, T. A.; F. N. Cornett and C. L. Pomernacki, Linear digital filtering for laboratory automation, Proc. IEEE 63, 1475–1486 (1975).

[5] Seelig, P. F., and H. N. Blount, Kalman filter applied to anodic stripping voltammetry: theory, Anal. Chem. 48, 252–258 (1976).

[6] Poulisse, H. N. J., Multicomponent-analysis computations based on Kalman filtering, Anal. Chim. Acta 112, 361–374 (1979).

[7] Lawton, W. H., and E. A. Sylvestre, Self-Modeling curve resolution, Technometrics 13, 617–633 (1971).

[8] Ohta, N. Estimating absorption bands of component dyes by means of principal component analysis, Anal. Chem. 45, 553–557 (1973).

[9] Kalman, R. E., A new approach to linear filtering and prediction problems, Trans. ASME Ser. D, J. Basic Eng. 82, 34–45 (1960).

[10] Gelb, A. (ed.) Applied optimal estimation, MIT Press:Cambridge, MA (1974).

[11] Fitzgerald, R. J., Divergence of the Kalman filter, IEEE Trans. Autom. Control AC-16, 736–747 (1971).

[12] Jazwinski, A. H. Stochastic Processes and Filtering Theory, Academic Press:New York, NY (1970) chapters 7-8.

[13] Mehra, R. K. On the identification of variances and adaptive Kalman filtering, IEEE Trans. Autom. Control AC-15, 175–184 (1970).

[14] Mehra, R. K. Approaches to adaptive filtering. IEEE Trans. Autom. Control AC-17, 693–698 (1972).

[15] Nahi, N. E. Decision-directed adaptive recursive estimators: divergence prevention. IEEE Trans. Autom. Control AC-17, 61–68 (1972).

[16] Brown, T. F., and S. D. Brown, Resolution of overlapped electrochemical peaks with the use of the Kalman filter, Anal. Chem. 53, 1410–1417 (1981).

[17] Rutan, S. C., and S. D. Brown, Pulsed photoacoustic spectroscopy and spectral deconvolution with the Kalman filter for determination of metal complexation parameters, Anal. Chem. 55, 1707–1710 (1983).

[18] Kaminski, P. G.; A. E. Bryson and S. F. Schmidt, Discrete square root filtering: a survey of current techniques, IEEE Trans. Autom. Control AC-16, 727–735 (1971).

[19] Brown, T. F.; D. M. Caster and S. D. Brown, Estimation of electrochemical charge transfer parameters with the Kalman filter, Anal. Chem. 56, 1214–1221 (1984).

[20] Rutan, S. C., and S. D. Brown, Adaptive Kalman filtering used to compensate for model errors in multicomponent analysis, Anal. Chim. Acta 160, 99–119 (1984).

[21] Deming, S. N., and S. L. Morgan, Simplex optimization of variables in analytical chemistry, Anal. Chem. 45, 278A–283A (1973).

[22] Nelder, J. A., and R. Mead, A simplex method for function minimization, Comp. J. 7, 308–313 (1965).

[23] O'Neill, R., Function minimization using a simplex procedure, Appl. Statistics 13, 338–345 (1971).

[24] Rutan, S. C., and S. D. Brown, Simplex optimization of the adaptive Kalman filter, Anal. Chim. Acta 167, 39–50 (1985).

[25] Irving, D., and T. F. Brown, Uranium speciation studies using the Kalman filter, Abstract No. 549; Proceedings, Pittsburgh Conference on Analytical Chemistry and Applied Spectroscopy; New Orleans, LA (1985).

[26] Brown, T. F.; D. M. Caster and S. D. Brown, Speciation of labile and quasi-labile metal complex systems using the Kalman filter, NBS Spec. Pub. 618, 163–170 (1981).

[27] Malinowski, E. R., and D. G. Howery, Factor analysis in chemistry, Wiley:New York, NY (1980).

# The Limitations of Models and Measurements as Revealed Through Chemometric Intercomparison

## L.A. Currie

### National Bureau of Standards, Gaithersburg, MD 20899

Interlaboratory Comparisons using common (reference) materials of known composition are an established means for assessing overall measurement precision and accuracy. Intercomparisons based on common data sets are equally important and informative, when one is dealing with complex chemical patterns or spectra requiring significant numerical modeling and manipulation for component identification and quantification. Two case studies of "Chemometric Intercomparison" using Simulation Test Data (STD) are presented, the one comprising STD vectors as applied to nuclear spectrometry, and the other, STD data matrices as applied to aerosol source apportionment. Generic information gained from these two exercises includes: a) the requisites for a successful STD intercomparison (including the nature and preparation of the simulation test patterns); b) surprising degrees of bias and imprecision associated with the data evaluation process, per se; c) the need for increased attention to implicit assumptions and adequate statements of uncertainty; and d) the importance of STD beyond the Intercomparison—i.e., their value as a chemometric research tool. Open research questions developed from the STD exercises are highlighted, especially the opportunity to explore "Scientific Intuition" which is essential for the solution of the underdetermined, multicollinear inverse problems that characterize modern Analytical Chemistry.

Key words: aerosol source apportionment; chemometric intercomparison; gamma-ray spectra; interlaboratory comparison; inverse problem; linear regression; multivariate data analysis; pattern recognition; reference materials; scientific intuition; scientific judgment; simulation test data.

## Introduction
### Accuracy Assessment

Ideally, the results of chemical analyses performed by a single laboratory using a well-defined Chemical Measurement Process (CMP) should be characterized by reliable measures of accuracy—i.e., imprecision and bias (or bounds for bias). Meaningful statements of uncertainty would then follow directly from these CMP Performance Characteristics [1][1]. Such is almost never the case, however. Once two or more laboratories perform measurements

of the same material, interlaboratory errors become evident. Collaborative tests, using common, homogeneous materials, serve as one of the most powerful means for both exposing and estimating the magnitude of this error component.

A familiar illustration of the outcome of interlaboratory measurement is reproduced in figure 1 [2]. Here, following the spirit of the "Youden Plot" [3], we show the results of pairs of measurements by 10 laboratories of the determination of trace levels of vanadium in two Standard Reference Materials (SRMs). A log transform has been applied to the reported concentrations, in order to expose proportionate errors among laboratories. As is generally the case, intralaboratory precision is comparable among laboratories, and considerably better than the interlaboratory component. Note that the line drawn in the figure is *not* fitted; its location is fixed by the certified values of the SRMs (dashed box),

**About the Author:** L.A. Currie is with the NBS Center for Analytical Chemistry where he leads the atmospheric chemistry group.

---

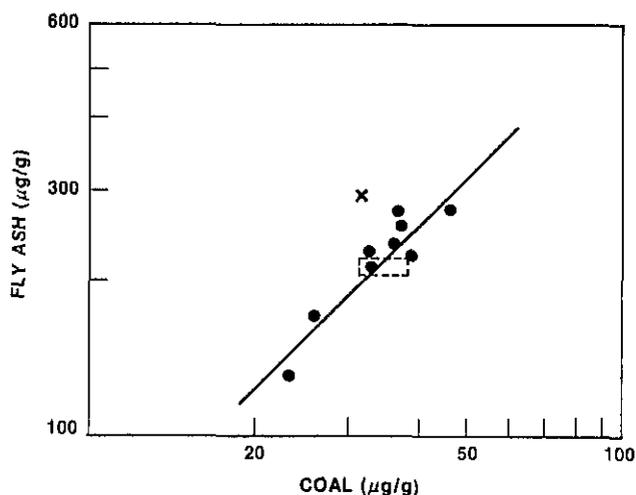[1]Figures in brackets indicate literature references.

**Figure 1**—Interlaboratory results for vanadium (μg/g) in two standard reference materials. The plot shows proportionate interlaboratory errors among nine participants using the same analytical method. The outlier (x) derived from a second method, lacking internal replication. Dashed region indicates the 'truth' [certified values].

and its slope is fixed at 45° (proportionate errors). SRMs or samples having known composition carry a very important attribute in interlaboratory tests, in that one can estimate individual laboratory and CMP bias in addition to interlaboratory variability. Also noteworthy is the "outlier" (marked by the cross) which deviates from the line significantly more than the other members of the interlaboratory set. Investigation of such outliers can sometimes yield important insight into the causes of interlaboratory differences. (In the example at hand the "outlier" resulted from a different analytical method that lacked internal replication.)

## Modern Analytical Chemistry: Importance of Data Evaluation

Enormous advances in analytical methods have brought great improvements in sensitivity, but at the same time, significant complications in data interpretation. In little over a decade, for example, "trace analysis" came to mean the measurement of pg of an analyte rather than μg [4]. Chemical patterns or spectra are central to the interpretation of complex mixtures, as are sophisticated analyte separation techniques such as high resolution gas chromatography. Practical demands on analysts also have accompanied the increase in sensitivity; toxic chemicals, for example, are regulated down to concentrations of $10^{-12}$ g/g. The magnitude of the problem can be appreciated from the fact that the analyst in measuring a substance at a concentation of $10^{-11}$ g/g in drinking water must contend with $\sim 10^5$ compounds which are more concentrated by at least a factor of 1000 [5].

When the Chemical Measurement Process involves significant modeling or numerical operations in the data evalu-

ation or information extraction step, it becomes interesting to consider the data analog of the SRM—the STD or Simulation Test Data set. By providing participants with common, well-characterized sets of data which adequately simulate the observations of real experiments, one can directly assess the imprecision and bias of the data evaluation process, independent of confounding errors or unreliable assumptions connected with the experimental parts of the CMP. This, in turn, makes it possible to estimate the error components associated purely with the experimental steps. Simulaton Data, as opposed to Real Data, are beneficial because "the truth is known"—i.e., the physical model (functional relation) as well as the random error model can be strictly controlled.

One might expect that little may be learned from such "Chemometric Intercomparisons" since numerical operations can be reproduced quite rigorously from laboratory to laboratory; but such is not the case. An illustration involving Real Data comes from the reevaluation (auditing) of several sets of chromatographic data from Love Canal soil and sediment samples for toxic organic compounds. As shown in table 1, compounds identified in common between analytical and auditing labs represented only about 60% of the total identifications, where the discrepancy was due strictly to differences in data evaluation [6].

**Table 1.** Real data—Love Canal soil and sediment samples: Compound identification by GC/MS. (Data Tape Auditing).

| Lab Code | Same Compounds Identified | Total Compounds Identified |
|----------|---------------------------|----------------------------|
| A | 20 | 32 |
| B | 13 | 24 |
| C | 22 | 51 |
| D | 13 | 20 |
| E | 63 | 104 |
| EPA | 14 | 20 |
| | [intersection] | [union] |

A myriad of hidden assumptions, and even algorithm changes exist in many of the pattern recognition and spectrum deconvolution schemes currently in vogue. Since the actual number of degrees of freedom is generally negative— i.e., the chemical model is never really known—numerical solutions often require subtle injections of "scientific intuition" or "scientific judgment." The importance of these issues will be illustrated by two case studies, actual intercomparisons among expert laboratories of the data evaluation phases in Gamma Ray Spectrometry and trace element Aerosol Source Apportionment, respectively. The STD in the first exercise was a data vector (nuclear spectrum); in the second, it was a data matrix (set of samples each having a trace element "spectrum"). The author was a participant in the first intercomparison and instigator of the second.

SRM                           STD

X ——— $\mathbf{Y}$ ┌─────────────┐ —— y —— $\mathbf{V}$ ┌────────────┐ ——— $\hat{x}$, $U_{\hat{x}}(\sigma, \Delta)$
(Sampling)       │ Sample Prep. │          │ Evaluation │        (Reporting)
                 │ Measurement  │          └────────────┘
                 └─────────────┘

Y = Ax + e_y   (1)

Chemical Measurement Process

STD

$\Theta$ ——— ┌──────────┐ ——x—— ┌─────┐ ——$\hat{x}$—$\mathbf{V}$— ┌───────────┐ ——— $\hat{\Theta}$
            │ Environ  │        │ CMP │         │ Environ   │
            │ Process  │        └─────┘         │ Evaluation│
            └──────────┘                        └───────────┘

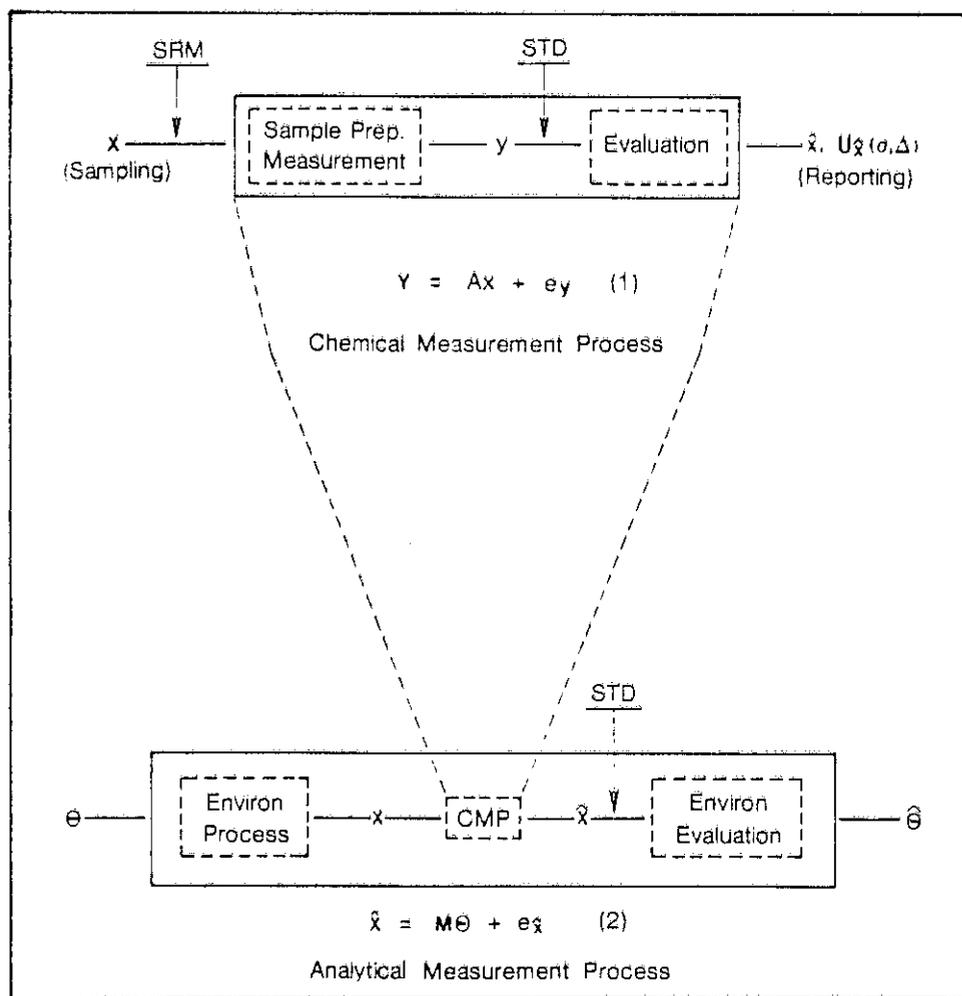$\hat{x}$ = M$\Theta$ + e_$\hat{x}$   (2)

Analytical Measurement Process

Figure 2—The Chemical Measurement Process (CMP) which operates within the laboratory, and the Analytical Measurement Process (AMP) which operates within the larger "Environmental" (or other external) system. (See text for explanation of symbols.)

# Chemometric Intercomparison
## Structure of the Measurement Process

In order to introduce some notation and to put the STD Intercomparison (IC) in perspective, it is useful to consider the structure of the Chemical and Analytical Measurement Processes (CMP, AMP) [4]. These two processes, which symbolize the·environment in which Analytical Chemistry operates, are shown in figure 2. As indicated in the upper portion of the figure, the CMP represents the laboratory process, where a sample of composition x is operated on (chemically) to produce a signal y, which in turn is operated on (mathematically) to generate an estimate of x and an uncertainty interval. The chemometric challenge thus is to obtain a chemically-meaningful and mathematically-consistent solution to the inverse problem as represented by eq 1. Control of the overall measurement process is achieved by injection of an SRM as a surrogate sample; control of the data evaluation process is achieved by injection of an STD as a surrogate signal.

Except in limited laboratory investigations, the real object of Analytical Chemistry is to provide information on an

external attribute (here represented by $\Theta$) through compositional analysis. The lower portion of figure 2 describes this broader context, where an external process (here labeled "environmental") operates an $\Theta$ to produce the sample of composition x. Following this, the imbedded CMP yields the compositional estimate $\hat{x}$. The final step, once again, is the solution of a (generally more difficult) inverse problem eq 2. STD injection in the AMP case means provision of a surrogate sample whose estimated compositional pattern corresponds to eq 2.

A fundamental difference exists between the CMP and the AMP with respect to the chemometric task. That is, in the laboratory, in principle, we can isolate ever-decreasing numbers of analytes (chemical fractions or instrumental signatures), in many cases leaving just a single term (component) in eq 1. For the AMP and the corresponding environmental, geochemical, or biochemical problem, for example, Nature is seldom so cooperative. That is, real samples x are determined by the external process over which we have limited control (beyond the sampling design), so eq 2 nearly always exhibits a multicomponent, multivariate structure. Unique solutions are generally impossible in the absence of

411

scientific knowledge concerning the external ("environmental") system.

The following STD intercomparison consists of univariate data (y) from a simulated CMP. The second example consists of multivariate data ($\hat{x}$) from a simulated AMP. Both ICs took place because of analytical measurement problems having major public import—the first related to accurate monitoring of radioactivity; the second, to accurate apportionment of atmospheric pollutants.

## STD Vector—
## IAEA Intercomparison
## of Gamma Ray Spectrum Analyses

In connection with their Analytical Quality Control Services program, the International Atomic Energy Agency (IAEA) undertook in 1976-77 a broad interlaboratory data evaluation exercise involving computer-simulated high resolution Ge(Li) gamma ray spectra such as might arise in contemporary neutron activation analysis [7]. The purpose of the intercomparison was both to assess the state of the γ-spectrum evaluation art and to provide data sets of known structure to assist in the improvement of that "art" (or science?). To my knowledge this was the first numerical chemometric intercomparison of such scope, using STD vectors. The organization and structure of the IAEA exercise are summarized in table 2.

Several features of the IAEA intercomparison were analogous to those involving chemical measurement intercomparisons with reference materials. First, the STD were well-characterized, with γ-rays of known identity (energy) and amplitude. The "samples" were absolutely homogenous (identical numerical data to all participants), and they simulated observations from actual laboratory samples. SRM

intercomparison organizers strive to also meet such conditions, but of course they can only approach the homogeneity and exact composition knowledge found with STD. The IAEA data sets also had known random error distributions (Poisson), a situation which is actually approached in many nuclear experiments, but which can never be guaranteed. Realism was preserved in the shapes of the γ-peaks, in that they were derived from high precision observations with Ge(Li) spectrometers. The fact that these shapes were not analytic was one of the more discriminating elements of the IC, particularly for the resolution of doublets, where alternative analytic or empirical peak shape functions *had* to be employed [8]. (Each peak was approximately Gaussian near the top but decidely asymmetric near its base.) Referring again to table 2, we can see that important planning took place, over a three-year period, resulting in four categories of data designed to provide initial "calibration," and to test detection, accuracy and precision in quantification, and doublet resolution. The importance of the pilot study cannot be overstated; development of realistic STD of sufficient but not excessive complexity does not come about without careful initial trials and iteration.

Detailed results for the IC may be found in [7]. Some of the highlights follow. Figure 3, for example, shows the spectrum (pattern) offered the participants, in digital and analog form, to address the problem of unknown peak detection. Participants knew only the calibration peak shape (as a function of "energy" or channel number) from the Reference Spectrum #100 (not shown), plus the facts that the unknown peaks were singlets and that random errors were Poisson. The numbers and locations of the trace peaks were to be determined. (The steep rise in the baseline near the center of the spectrum was inserted by the IAEA to simulate a Compton Edge.) The inset shown in figure 3 gives some

---

Table 2. Structure of the IAEA Gamma-Ray STD intercomparison.

**Objectives**

- To permit each participant to assess the accuracy of his data evaluation process.
- To determine the quality of alternative gamma-ray spectrum evaluation methods as applied in representative laboratories.

**Evolution**

- 1973: Proposed at Consultants' Meeting.
- 1975-6: Pilot Study involving a small number of experts.
- 1976-7: Full IC, involving 163 labs in 34 member states.
- Currently: Simulation data offered as continuing part of the IAEA Analytical Quality Control Service.

**Data Sets**

- Reference Spectrum: 20 high-precision peaks spanning 2000 channels.
- Detection Spectrum: 22 subliminal peaks, whose number and locations were unknown to participants; detection criteria (α-, β-errors) were left to individual judgment.
- Precision Spectra: 6 replicate spectra having 20 known plus 2 unknown, large singlet peaks (Poisson statistics).
- Resolution Spectrum: 9 doublets of unknown location and relative amplitude.
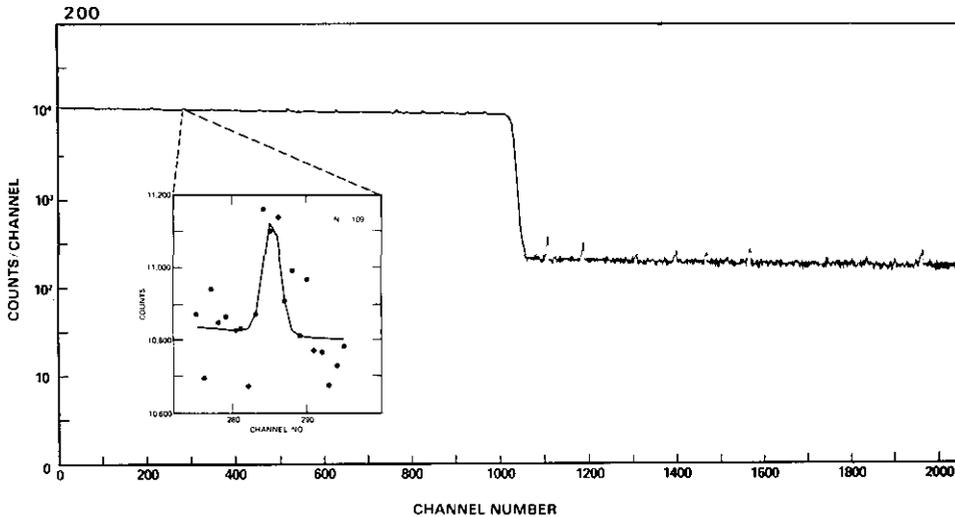
Figure 3—IAEA Gamma-Ray STD; "Detection" Spectrum. Inset shows discrete data for a peak detected by about 50% of the participants. Adapted from [7.]

idea of the discreteness and scatter of the digital data; this "real" peak was detected by about half of the participants.

The results of this intercomparison were somewhat surprising. Though most of the 200-odd participants submitted results, no one correctly identified all 22 subliminal peaks. Some six classes of methods, including one labeled "unclassified," were employed, including: the relative maximum, first and second derivatives, cross correlation, and "visual". It is interesting that the last gave the best result; one "trained eye," using analog data only, identified 19 peaks correctly, with no false positives! Understanding the process (Scientific Intuition) employed by this expert is certainly one of the more intriguing aspects of this work. It seems hardly pure chance, for only 5 out of 212 participants correctly reported this many (19) peaks; yet 2 of the 5 were "visual." (For comparison, the visual technique was employed by about 5% of the participants.) The second derivative and cross correlation techniques were close behind, with up to 18 and 17 peaks correctly identified (w/o false positives), respectively. Performance was quite diverse for all methods, however: correct identifications ranged from 2 to 19 peaks, and false positives ranged from 0 to 23. Apparently, Detection Limits were rarely estimated, for this issue was not even mentioned in [7]. Histograms for the three "best" methods are shown in figure 4. Though all three exhibit considerable dispersion, it is clear that the visual technique gave the best single result as well as the smallest fraction of false positives.[2]

The replications (Spectra 300-5) and resolution (Spectrum 400) exercises also indicated often inadequate and

---

[2]Scientific Intuition, as employed by experts, is alleged to be much more disperse that "rule-based" methods [9]. In the light of figure 4, it is not obvious that this presumption is true, for even the "objective" numerical techniques employed by different laboratories operating on exactly the same data gave broad distributions.
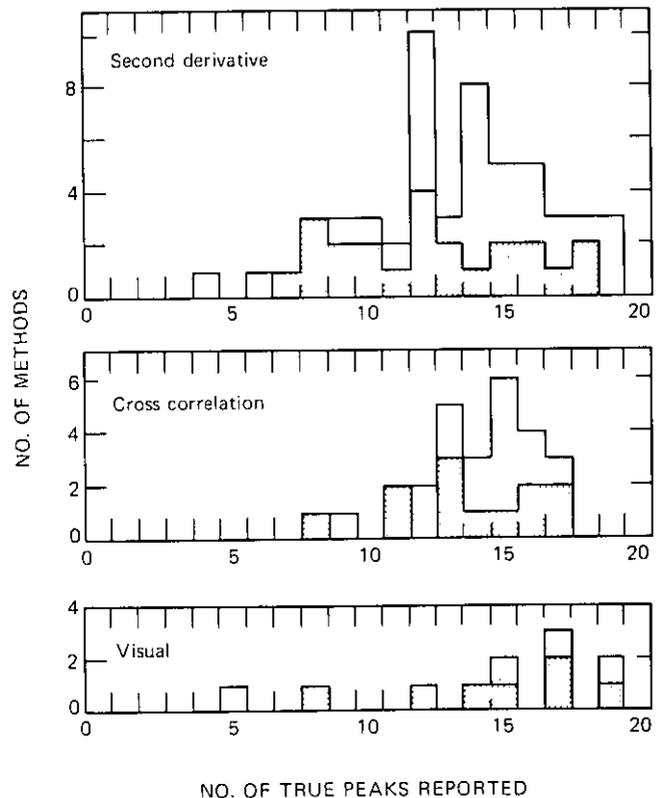
## SCIENTIFIC INTUITION



NO. OF TRUE PEAKS REPORTED

Figure 4—Frequency distribution of results for peak detection according to the type of method used. The upper boundary of each histogram represents the data for all results regardless of the number of spurious peaks reported; the upper boundary of the shaded region is for those results which were accompanied by zero spurious peaks.

413

widely varying performance. The majority of the results submitted contained quite inaccurate or no estimates of uncertainty for large singlet peaks, even though the random error distribution was known; and less than 25% of the participants even submitted results for the most difficult doublet resolution case.

The IAEA Simulation Data Sets have been viewed of sufficient importance that they have become an integral part of the Intercomparison Programme of the Analytical Quality Control Service of that organization. The most recent offering was issued in December 1984 [10] where STD γ-ray spectra are included alongside isotopic and trace element Intercomparison and Certified Reference Materials of importance in many areas of nuclear and environmental analysis.

**STD Matrix—NBS-EPA Intercomparison of Source Apportionment Techniques.** The second case study comprises STD in the form of two-dimensional data matrices, simulating sets of atmospheric aerosol samples each analyzed for up to 20 chemical species [11]. The stimulus for this exercise, which is believed to be the first STD intercomparison involving Data Matrices, was the great potential but great difficulty of identifying multiple pollutant sources via their "chemical fingerprints" as preserved in ambient particles. (A vivid illustration adjoins [11], where one finds discord even in assigning names to pollutant factors deduced from elemental patterns observed in actual measurements of [Houston] aerosol samples [12].) As noted at the beginning of this section, this type of problem is characteristic of the AMP, where superposition of multiple components is intrinsic to the nature of the system, so chemical manipulation cannot simplify the structure of eq (2).

We designed the STD in coordination with (nearly) all of the "Receptor Modeling" (source apportionment) experts in the U.S. with the object of providing a few realistic data matrices covering a range of problems. The overall structure of the study is given in table 3. The data matrix x is given by the superposition of source contributions $(M\theta)_j$, where each source has a characteristic chemical pattern or profile $M_j$ and a temporal (or spatial) intensity pattern $\theta_j$. Three classes of error typify such measurements, as indicated in the table. A significant task involved building the database of source profiles and error terms. Unlike the γ-ray calibration profiles (peak shapes), the aerosol source profiles were not even approximately analytic (fig. 5, top); reliable empirical field data had to be sought and evaluated.

When generating data matrix STDs, one must pay attention to a major new element of complexity which is absent from data vector STDs. That is, unless the simulation data are to be no more than arbitrary superpositions of sources with added random errors, it is essential to generate the source intensity patterns by means of suitable source emission locations (emissions inventory map) plus realistic mixing and transport to the receptor (sampling) site(s). This step is intrinsic to the nature of the data matrix; it underlies its multivariate character, and it highlights information (source map, meteorological patterns,...) external to the data matrix, per se, which may be crucial for a successful data analysis. The source map used for the simplest of our three data sets is shown in the lower portion of figure 5. Stochastic source impacts at the receptor **R** were generated by bringing together the source map, emissions intensities and operating schedules, actual meteorological data (from St. Louis, September 1976), and the 'RAM' atmospheric dis-

Table 3. Structure of the source apportionment simulation test data.

**Generating equation**

$$\hat{x}_{it} = \sum_{j}^{P} [M - e_m + e_H]_{ij} \, \theta_{jt} + e_{it}$$

where:

$p$ = number of active sources ($p \leqslant 13$)

$t$ = sampling period ($1 \leqslant t \leqslant 40$)

$\hat{x}_{it}$ = "observed" concentration of species $i$ for period $t$ ($1 \leqslant i \leqslant N, N \leqslant 20$)

$\theta_{jt}$ = true intensity (at receptor) of source $j (1 \leqslant j \leqslant p)$

$M_{ij}$ = 'observed' source profile matrix (element $i,j$)

$e_i$ = random measurement errors, independent and normally distributed

$e_m$ = systematic source profile errors, independent and normally distributed (systematic because fixed over the 40 sampling periods)

$e_H$ = random source profile variation errors, independent and log-normally distributed

**Data set characteristics**

Set I: $p = 9$ (including one unknown source);* errors=$e_i$, $e_m$; City Plan No. 1 (fig. 5)

Set II: $p = 13$ (all known); errors=$e_i, e_m$; City Plan No. 2

Set III: $p = 13$ (all known); errors=$e_i$, $e_m$, $e_H$; City Plan No. 2

*For Data Set I, participants were told only that $p \leqslant 13$.

mg/g

100

10

$M_{i2}$      1

0.1

C Na        Cl K                Zn        Pb        i (Element)

←1%

μg/m³

10

$\Theta_{2t}$      1

0.1

0.01

10        20        30        40        t (Period)

R   RECEPTOR
0   ROAD (AUTO + SOIL)
1   STEEL A
2   INCINERATOR
3   OIL B
4   SAND BLAST
5   COAL B
▨   SOIL
▦   BASALT
▩   WOOD
EACH AREA SOURCE
IS 2 km SQUARE

**Figure 5**—Source Apportionment STD (Data Set I), Upper portion shows one column (transposed) of the source signature matrix M, and one row of the source intensity matrix θ—both for source-2, IN-CINERATOR. $M_{i2}$ has a discrete pattern (individual chemical elements), the most discriminating elements of which are marked by circles; the dashed line indicates which elements exceed 1% of the (Incinerator) particle mass. $\theta_{2t}$ has a continuous underlying structure (time series) which is sampled at 40 equidistant points; the dashed line indicates samples for which the Incinerator source contributes more than 5% of the average aerosol mass.

The lower portion of the figure displays the aerosol source emission map.

persion model [13]. Illustrations of a source (chemical) signature $[M_{ij}]$ and a source intensity time series $[\theta_{jt}]$ are given in the upper portion of the figure. Note that $M_{ij}$ by its nature (individual elements) is discrete, whereas $\theta_{jt}$ is a sampled continous time series (intensity variations).

The objective of spanning the range of difficulty was achieved. Our initial "pilot" data matrix was so transparent that one of our participant-advisors was able to identify sources by inspection. Caused in part by the narrowness of

the RAM model plumes, this was remedied in the final STD intercomparison data sets, one of which was so difficult (though realistic) that none of the participants submitted results. Results for the data set of intermediate difficulty (Set II) were generated by three laboratories, all using regression techniques. Two of these were identical: "effective variance" weighted least squares (WLS), which took into account errors in the observed chemical concentrations as well as those in the source profiles. The third method was

415

ridge regression, of interest because of the high degree of collinearity among the 13 sources. For the most part, the three sets of results were self consistent, and within a factor of 2 of the truth. For four of the sources, however, widely discrepant results were reported—a surprising outcome, for similar or identical numerical methods were applied to identical numerical data.

An exploratory graphical analysis of the results from the two laboratories using weighted least squares (WLS-1, -2) is shown in figure 6. This approach, which was inspired by the "Youden Diagram" for bi-material intercomparisons, allows us to spot "outliers" from the line of concordance. (As with the line drawn for the Youden plot of figure 1, the 45° line in figure 6 was drawn independent of the data—i.e., it was in no way "fitted".) Exact results would fall at the origin (0, 0), and equivalent data reduction by the two laboratories would produce results lying on the line. Dispersion along the line derives from random errors built into the common data set and systematic errors connected with the common numerical model (WLS). Two of the sources ($Ag$, $O_2$) were well below their detection limits, and will not be further discussed here. The three outliers ($cf$ point 'x' in



Figure 7—Absolute errors for each of the 13 source estimates for Data Set II, by WLS-1 (×), WLS-2 (+) and Ridge Regression (□), plotted as a function of the reported standard errors. (On the average about 2/3 of the points for a correct method should fall below the diagonal.)
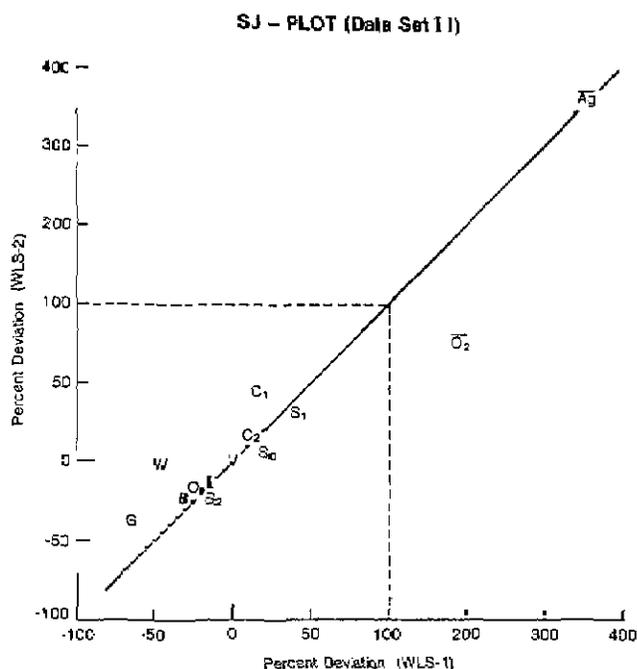
SJ — PLOT (Data Set I I)



Figure 6—"Scientific Judgment" Plot. Correlation (Youden-type) Diagram showing data evaluation results of laboratory-2 vs laboraory-1 operating on the same data set (II) using the same method of numerical analysis (Weighted Least Squares). Deviations from the non-fitted line of concordance imply chemometric "operator error" (SJ).

[Source Codes: Steel-A ($S_1$), Steel-B ($S_2$), Oil-A ($O_1$), Oil-B ($O_2$), Incinerator ($I$), Glass Mfr ($G$), Coal-A ($C_1$), Coal-B ($C_2$), Aggregate ($Ag$), Basalt ($B$), Soil ($So$), Auto ($V$), Wood Smoke ($W$).]

figure 1)—$G$, $W$, and $C_1$—deserve our attention, however. As will be shown in the next section, they may be attributed to what the laboratories involved labeled "Scientific Judgment." Discovering this factor, and understanding its nature, was one of the unexpected but important outcomes of the experiment.

Another major difference and inadequacy among the laboratories relates to the uncertainties (i.e., SEs) reported. All did report standard errors, but an examination of the actual deviations from the truth for the 13 estimates from each lab was illuminating. For one of the three, all 13 (absolute) deviations exceeded the SEs by factors of 2 to 10, whereas for another lab the deviations were all smaller than the SEs by factors of about 1.6 to 30. None of the labs reported bounds for systematic or model error. (See fig. 7.)

Because of the multivariate character of the source apportionment data sets, both simple linear regression and factor analysis (FA) techniques could be applied. The regression methods ("Chemical Mass Balance", i.e., WLS) were the more precise, when model information and source profiles were available. When the model was not fully known and if the number of interfering components was not too great, factor analysis or certain hybrid methods appeared to yield more acceptable results. Such was the case for Data Set I which had 9 sources, one of which was unknown to the participants. In Set II, however, with 13 known sources having a significant degree of multicollinearity, factor analysis was able to discern but four sources.

416

## Subsequent Developments: STD as a Basis for Chemometrics Research

A principal conclusion from the foregoing exercises is that the data evaluation step of multicomponent and multidimensional chemical (nuclear) analysis is *not* free from problems of imprecision and bias. Scientific Intuition (SI) and Scientific Judgment (SJ), often manifest as subtle assumptions, can provide important guidance for pattern recognition, or it can be somewhat misleading. Collaborative STD exercises appear to be an effective means for exposing and perhaps better understanding these "expert" techniques.

The limitations found in the STD intercomparisons— *e.g.*, limitations in accuracy, in model formulation, in uncertainty estimation, in detection, and in utilizing external information (non-negativity, meteorological data,...)—suggest that an important part of experimental inaccuracy, as seen for example in SRM intercomparisons, may lie with the data evaluation process. Attention to this matter offers the possibility of improved overall performance, either through the reduction of needless data evaluation error, or through improved measurement process design to reduce model complexity and multicollinearity.

Both intercomparison exercises exhibited an afterlife. As noted earlier, the Gamma-Ray STD have become a routine part of the IAEA's Analytical Quality Control Service. The source apportionment STD have spontaneously evolved into a research data set for chemical element pattern recognition method development. Initially, new research was undertaken by participants who wished to try improved versions of the their regression or FA methods in light of the IC results and knowledge of the truth. More recently, requests for the data tape have come from others for the testing and development of multivariate pattern recognition methods quite apart from aerosol source apportionment. A sampling of post-intercomparison research with the STD data matrices follows:

| Investigator | Topic | Ref. |
|---|---|---|
| M-D. Cheng, P. Hopke* | Linear Programming | [14] |
| L. Currie* | Detection, Design, Model Error | |
| I. Frank, B. Kowalski | Partial Least Squares | [15] |
| G. Gordon* | Student Instruction, QA | |
| R. Henry | Composite Components (SVD) | [16] |
| P. Lioy | Student Research | |
| D. Lowenthal et al. | Special Error Propagation (Cov) | [17] |
| T. Pace | Sensitivity Analysis | |

*Participants of the original IC, known to be performing advanced studies using the STD. Most of the others listed subsequently requested the data tape specifically for basic investigations of the numerical data evaluation process.

Table 4 includes results from some of these more recent investigations, together with an examination of the "outliers" and Scientific Judgment exposed in the correlation diagram (fig. 6). Table 4A lists results for the best prior results (TTFA) for Data Set I and the new application of Partial Least Squares (PLS). Though PLS results carry no uncertainty estimates, they are clearly closer to the truth. The two deceptions, as well as all of the source profiles, however, were known for carrying out the PLS analysis. Regarding TTFA and other prior analysis, the first deception—the linear combination of AUTO and SOIL to represent the ROAD component, including re-entrained dust— was missed by all participants; the presence of a large, additional component (SANDBLAST) was discovered by all.[3]

**Table 4A.** Source apportionment STD (Data Set I).

| Source | Truth [μg/m$^3$] | PLS[1] | TTFA[2] | Δ/SE[2] |
|---|---|---|---|---|
| Steel (c) | 0.05 | 0.07 | — | ? |
| Oil (c) | 2.0 | 2.5 | 2.1±1.0 | +0.1 |
| Incinerator | 1.3 | 1.4 | 1.9±0.14 | +4.3 |
| Coal-B | 2.4 | 1.9 | 2.2±0.73 | −0.3 |
| "Crustal" (c) | 12.7 | 13.0 | 12.5±0.75 | −0.3 |
| *Road | 7.1 | 7.2 | *4.0±0.09* | *−34.4* |
| Wood | 3.3 | 3.4 | 4.3±0.52 | +1.9 |
| *Sandblast | 4.2 | 4.1 | 4.1±0.20 | −0.5 |
| (Total) | | (+2%) | (−6%) | |

(c) Composite components for multicollinearity reduction.
*Two Deceptions: Road (Soil+Auto); Sandblast (Unk).

[1]I. Frank and B. Kowalski [15].　　[2]P. Hopke [11].

**Table 4B.** Source apportionment STD (Data Set II)—analysis of outliers.

| | Incinerator (I) | Glass (G) | Coal-A (C$_1$) | Wood (W) |
|---|---|---|---|---|
| Truth (μg/m$^3$) [N] | 1.81 [38] | 0.46 [36] | 3.7 [32] | 0.94 [40] |
| WLS-1[1] [N] | 1.5±0.05 [27] | 0.16±0.03 [3] | 4.3±0.3 [11] | 0.5±0.1 [11] |
| WLS-2[2] [N] | 1.5±0.43 [26] | 0.28±0.11 [8] | 5.3±1.4 [22] | 0.91±0.16 [26] |
| WLS-3[3] [N] | 1.51±0.14 [31] | 0.43±0.44 [23] | 4.3±2.2 [29] | 1.33±0.24 [36] |
| LI[4] | 1.60 | 0.48 | 0.44 | 0.45 |

N=Number of actual or assumed non-zero occurrences of the source in question among the 40 aerosol samples.

[1]Heisler, Shah [11].　　[3]Lowenthal, Hanumara, Rahn, Currie [17].
[2]Cooper, DeCesar [11].　　[4]Cheng, Hopke [14].

The outliers $(G, C_1, W)$ from data Set II are examined in table 4B for new and earlier analyses. The INCINERATOR results, which fell directly on the line of concordance (fig. 6), are included for comparison. It is interesting that L1 linear programming gave significantly improved results for sources $I$ and $G$, but much poorer results for the other two sources. A conjecture, which bears investigation, is that in the absence of experimental blunders, least absolute residuals carry too high a penalty in highly collinear problems. We gain some understanding of Scientific Judgment (SJ) from the three WLS results. It is apparent that the exercise of SJ in deleting rare components from the model is operator-dependent and generates important differences in overall results. This is an issue deserving additional research, the outcome of which could be the transformation of Ad-Hoc SJ into scientifically-based SI.[4]

## Conclusion

For modern Analytical Chemistry, where chemometric approaches are mandatory to resolve complex signals from multianalyte mixtures, Simulation Test Data serve two important purposes: 1) The assessment of interlaboratory data-evaluation precision and accuracy; and 2) the exposure of SI and SJ plus the generation of more powerful methods of combining scientific knowledge with advanced techniques of data analysis. The need for this is critical, because in principle nearly all of our multicomponent "inverse" problems are underdetermined—i.e., solutions cannot obtain in the absence of explicit (or hidden) assumptions. We conclude with a summary of recommendations for the preparation of STD data sets, and open research questions which could be fruitfully addressed by chemometricians (table 5).

**Table 5.** STD exercises.

### Observations and Recommendations

- Benefits: Controlled model, errors; Truth is known.
- Prerequisites: defined objectives, 'lab' population, plan, input data/model.
- Pilot study advisable.
- Investigate 'surprises'—unexpected accuracy, discrepancies.
- Blind; anonymous; don't score; discourage premature communication.
- Deceptions: mimic reality; be not too obvious.

### Some Open Questions

- Adequate treatment of uncertainty—esp. bias (bounds).
- Further understanding of 'SJ' and 'SI'—who are the experts?
- Utilization of external information: phenomena, data, "fuzzy" & partial knowledge, constraints,...
- Treatment of the 'total' problem (simultaneous estimation over the entire data matrix), incl. errors in 'y' and 'A'.
- Solution and uncertainties when the linear model doesn't apply (physical-chemical and mathematical strategies; eg, atmospheric transformation).

Special recognition should be given to the designers of the two STD exercises discussed, as well as to the participants and the scientists who are using the STD for continued basic research in chemometrics. R.M. Parr (IAEA - STD) and R.W. Gerlach (Source Apportionment - STD) made essential contributions. Others deserving special recognition are listed as authors of references [7, 11], and [14 - 18].

---

[3]STD deceptions—i.e., realistic complications—are, of course, in order for all but the simplest of exercises. Though organizers should attempt to span the range of difficulty occurring with real data sets, they should not do so in too obvious or regular a manner. Informing the IAEA γ-ray STD participants, for example, that the multiplets had just 2 components was already a considerable simplification, but the regular spacing (1, 3, 10 channels) and relative amplitude (1, 3, 10) of the doublet members was unnatural and could encourage a certain amount of guess work on the part of the participants.

[4]The paradigm proposed here classifies experts' decisions or assumptions as "intuitive" (SI) or "judgmental" (SJ) depending on whether they are based on sound (though possibly subliminal) reasoning or ad hoc judgments, respectively. These asymptotic classes may each yield correct or incorrect results, just as target and contaminating populations may each produce outliers or inliers, though with differing probabilities.

## References

[1] Eisenhart, C., Science 160 (1968) 1201.
[2] Currie, L.A., and J.R. DeVoe, Chapt 3 in J.R. DeVoe, Ed, VALIDATION OF THE MEASUREMENT PROCESS, Amer Chem Soc Sympos Ser #63 (1977).
[3] Youden, W.J., Anal Chem 32 [13] (1960) 23A.
[4] Currie, L.A., Pure & Appl Chem, 54 (1982) 715.
[5] Lamparski, L.L., and T.J. Nestrick, Anal Chem, 52 (1980) 2045.
[6] Kirchhoff, W.H., Edit, NBSIR 82-2511 (1982).
[7] Parr, R.M.; H. Houtermans, and K. Schaerf in COMPUTERS IN ACTIVATION ANALYSIS AND GAMMA-RAY SPECTROSCOPY, US Dept of Energy, CONF-780421 (1979) 544. (See also Zagyvai, P.; R.M. Parr and L.G. Nagy, J. Radioanal. Nucl. Chem. 89 589 (1985).)
[8] Ritter, G.L., and L.A. Currie, ibid, p. 39.
[9] Nalimov, V.V., FACES OF SCIENCE, ISI Press (Philadelphia, 1981).
[10] IAEA Analytical Quality Control Service Program, LAB/243 (1984).
[11] Currie, L.A.; R.W. Gerlach, C.W. Lewis, W.D. Balfour, J.A. Cooper, S.L. Dattner, R.T. DeCesar, G.E. Gordon, S.L. Heisler, P.K. Hopke, J.J. Shah, G.D. Thurston, and H.J. Williamson, Atm Environ 18 (1984) 1517.
[12] Dzubay, T.G.; R.K. Stevens, W.D. Balfour, H.J. Williamson, J.A. Cooper, J.E. Core, R.T. DeCesar, E.R. Crutcher, S.L. Dattner,

B.L. Davis, S.L. Heisler, J.J. Shah, P.K. Hopke, and D.L. Johnson, Atm Environ 18 (1984) 1555.

[13] Novak, J.H., and D.B. Turner, J. Air Polut. Control Assoc., 26 (1976) 570.

[14] Cheng, M-D., and P.K. Hopke, An Intercomparison of Linear Programming Procedures for Aerosol Mass Apportionment, Air Pollution Control Association Paper No. 85-21.8 (1985).

[15] Frank, I.E., and B.R. Kowalski, Statistical Receptor Models Solved by Partial Least Squares, ACS Divn of Environ Chem Sympos Abstracts (Philadelphia, Aug 1984) p. 202.

[16] Henry, R.C., Proc Air Pollut Control Assoc Spec Conf, SP-48 (1982) 141.

[17] Lowenthal, D.H.; R.C. Hanumara, K.A. Rahn, and L.A. Currie, Estimates and Uncertainties in Chemical Mass Balance Apportionments: Quail Roost II Revisited, prepared for submission to Atmospheric Environment (1985).

[18] Stevens, R.K., and T.G. Pace, Atmos. Environ. 18 (1984) 1499.

# DISCUSSION
of the L.A. Currie paper, The Limitations of Models and Measurements as Revealed Through Chemometric Intercomparison

## Leon Jay Gleser
Department of Statistics
Purdue University

The construction and use of Simulation Test Data (STD) to help evaluate alternative chemometric methodologies is a highly welcome contribution to the field. Dr. Currie, and the agencies and colleagues whom he credits, are to be congratulated for an approach which has the potential to promote improvements in the art of quantitative chemical analysis.

What follows is a brief discussion of some previous use of standard data sets in statistical research, along with some warnings about the possible pitfalls connected with the use of such approaches. In particular, the parallel that Dr. Currie draws between the use of standard data sets and interlaboratory comparisons using common reference materials cannot be pushed too far. Many interacting factors lead to bias in modeling and analysis of complex data sets; the contributions of these factors would be confounded in typical interlaboratory comparison designs. One factor, scientific judg-ment, cannot even be identified in standard frequentist reports of statistical data analysis. This suggests that subjective scientific judgments need to be given more explicit mention in reports of statistical analyses, perhaps through the use of the Bayesian approach to inference.

To make standard data sets more closely resemble real-world data, the use of the "bootstrap" is suggested. The "bootstrap" can also help in providing the estimates of statistical precision that Dr. Currie notes were lacking in the two studies conducted to date.

## Standard Data Sets in Statistics

Statisticians have long recognized the usefulness of having common data sets on which new methodologies can be tried out, and their relative merits assessed. For example,

419

new methodologies for classification and discriminant analysis are often applied to the iris data of E. Anderson [1][1], which was featured in a famous paper by R. A. Fisher [2]. (To Anderson's undoubted frustration, these data are usually referred to as "Fisher's iris data.")

Another famous data set is Longley's [3] econometric linear regression data. In these data, the independent variables in the regression are highly interrelated (multicolinear). Longley ran the data through several computer software packages designed to do least squares analysis. In theory, all of these programs solve the same set of linear equations to estimate the regression slopes. However, the solutions obtained by the various algorithms differed, in some cases even by sign! What had happened was that the multicollinearity in the data made the answers obtained highly sensitive to roundoff and truncation of the data, and the algorithms differed by where and by how much roundoffs were done. Longley's paper had the very beneficial consequence that software developers now pay careful attention to numerical analysis in designing statistical algorithms. Further, it stimulated study of the resistance of statistical methodology to data perturbations (robustness).

However, Longley's paper (and particularly his data) may also have had a less salutary effect on software development. Software developers now know that consumers will test out their programs on Longley's data. [See, for example, Lachenbruch's review [4] of STAN, Version II.0 by David Allen.] This may lead them to overcompensate for multicollinearity problems, and consequently overlook or neglect other potential problems or sacrifice desirable features to include subroutines necessary to accurately process multicollinear data.

This last comment points out a real danger in the use of standard data sets, namely that their existence can bias the direction which development of methodology and software takes. The best guard against such bias is the creation of standard data sets of many types.

An *artificial* standard data set (simulated according to a known model for the distribution of errors) can lead to a particularly serious bias. Chemometricians who know that their work will be evaluated by such data sets will tend to use a methodology which is known to be efficient for the given statistical model. Such a methodology, however, may not do well against *real* data, for which the given statistical model is not necessarily a good approximation. Alternatively, chemometricians may object to evaluations on the basis of such data, arguing (with considerable merit) that such data do not reflect their practical experience.

The reason, of course, for using artificial data is that the "truth" or "signal" underlying the "noise" (error) in the data is known. This allows us to separate *bias* (lack of validity) from *precision* (reliability, repeatability). In this respect

there is an obvious parallel, which Dr. Currie correctly points out, with the use of common reference materials in interlaboratory comparisons. The goal of such studies is to eliminate bias (which is usually reflected in interlaboratory variation), and to estimate precision (intralaboratory variation). However, whereas common reference materials are "real" (although they may be ideal examples of materials analyzed in practice), this is not clearly the case with data simulated from specified statistical populations (e.g., Gaussian populations). Real populations may have "heavy tails" and/or other funny features (e.g., several modes) which are not modeled by standard distributions.

One obvious solution is to vary the distributional assumptions which generate the errors in artificial data. This approach is widely used in statistics to study the *robustness* properties of statistical methodologies.

Another possible solution is to use the ideas underlying the "bootstrap" (Efron [5], Diaconis and Efron [6], Freedman and Peters [7,8]) to simulate data which have "real world" error distributions.

## The Bootstrap

In using the "bootstrap," we start by assuming that the observed data $y_i$ are related to unknown parameters $\theta$ and errors $e_i$ by a model

$$y_i = G(\theta, e_i), \qquad i = 1, 2, \ldots, n, \qquad (1)$$

where $G(\cdot, \cdot)$ is known. Given a value for $\theta$ (which may be a vector), we assume that the eq (1) can be inverted to obtain the errors $e_i$. That is,

$$e_i = H_i(\theta, y_1, \ldots, y_n), \qquad i = 1, \ldots, n. \qquad (2)$$

Given "real" data $y_1, y_2, \ldots, y_n$, with $n$ sufficiently large to give us some hope of accurately estimating $\theta$ by

$$\hat{\theta} = \hat{\theta}(y_1, \ldots, y_n),$$

we now construct the residuals (estimated errors)

$$\hat{e}_i = H_i(\hat{\theta}, y_1, \ldots, y_n), \qquad i = 1, 2, \ldots, n. \qquad (3)$$

The resulting finite population $\{\hat{e}_1, \ldots, \hat{e}_n\}$ of residuals is the statistical population from which we can randomly sample new errors $\tilde{e}_i$, $i = 1, 2, \ldots, m$, to create standard data sets

$$\tilde{y}_i = G(\theta^*, \tilde{e}_i), \qquad i = 1, 2, \ldots, m,$$

where $\theta^*$ can be chosen to have any desired value.

---

[1]Figures in brackets indicate literature references.

420

The data sets so simulated are not entirely "real." The model (1) relating observations $y_i$ to errors $e_i$ must still be specified, and need not be correct. However, such models can be specified (and criticized) taking account of chemical and physical theory, without also imposing statistical assumptions about distributions of errors. As such, these models are "prestatistical."

A population of residuals $\hat{e}_i$ that is too small can misrepresent statistical variation. Thus, attempts should be made to constantly enlarge this population with new residuals obtained from real data obtained in contexts described by the model (1). Since the "bootstrap" is a fairly recent statistical development, new insights into problems and advantages connected with the method are constantly being published. Consequently, the input of specialists in "bootstrap" methodology should be sought when applying this method to the generation of standard data sets. In particular, changes in instrumentation, personnel, or experimental design may change the error population over time. Careful attention should be paid to detect such shifts in distribution.

Not all measurement contexts lend themselves to the "bootstrap," since the transformation (2) from observations to errors may not exist, or may not be well defined. (This may be the case, for example, with the Gamma Ray Spectrum Analysis example discussed by Dr. Currie.) However, when the "bootstrap" does apply, it can be used both to create standard data sets, and also to provide nonparametric estimates of precision [5,7,8].

## Estimates of Precision

Although I share Dr. Currie's concern that the laboratories in his two examples either failed to provide estimates of precision, or gave incorrect estimates, I must point out that in Dr. Currie's two examples, it is not clear what measures of precision are appropriate. In all of the analyses, multiple decisions are made. For example, in the Gamma-Ray examples, the locations and amplitudes of several peaks had to be determined *simultaneously*. Although individual standard errors can be given, these do not directly provide measures of simultaneous accuracy [9]. Further in the detection spectrum and precision spectra sets, even the *number* of peaks was unknown. This produces a highly complicated estimation problem for which only large-sample approximations to precision are available. There is some evidence in the literature that such large-sample approximations have considerable bias in moderate samples (see, e.g., [6]).

Similar problems arise in the three data sets for the NBS-EPA Source Apportionment study, particularly in the case of Data Set I, where the number of sources is left unspecified. Both ridge regression and factor analysis are exploratory methodologies, requiring iteration and judgment that are difficult to describe analytically. The only available

measures of precison for such techniques are large-sample approximations which refer to analytical formulas for the estimators not directly related to the way such estimators are actually obtained. For example, I know of no way to specify the precision of estimates of slope obtained by the ridge trace method. Published formulas for the precisions of ridge regression estimators refer to those estimators in which the ridge factor $k$ is a specified function of the data, rather than being obtained by inspection of the ridge trace.

Given the complex natures of the estimation problems that Dr. Currie describes, and the fact that statistical theory has not yet provided reasonable estimates of precision for some of the methodologies used in these problems, it is not surprising that the laboratories either failed to provide measures of precision, or gave estimates that were off the mark. Clearly, there is much theoretical statistical work yet to be done.

In the meantime, it should be mentioned again that the "bootstrap" can provide estimates of precison in cases where the assumptions (1), (2) underlying the bootstrap are applicable.

## Analogy to CMP Interlaboratory Comparisons

As already noted, Dr. Currie makes an analogy between the use of standard data sets in the two examples he discusses, and traditional interlaboratory comparisons using common reference materials. However, this analogy cannot be carried too far, since there are some important differences in context.

In traditional interlaboratory comparisons, differences between laboratories are usually assumed to be due to variations in the calibration (adjustment) of the instruments, or to differences in instrumentation or technique. Consequently, a one factor (additive) components of variance ANOVA model can reasonably be employed to assign variability between inter- and intra-laboratory sources.

In the standard data set context described by Dr. Currie, however, there are at least three factors which can describe variability between laboratories: 1) different models or assumptions used, 2) different statistical methodologies employed, and 3) different numerical algorithms. Further, the "levels" of these factors (particularly factor 1) appropriate to describe a given laboratory's analysis are not always apparent. (Not all assumptions made are clearly stated). For example, "outliers" can be discarded, parameter values can be truncated (e.g., negative estimated amplitudes reported as zero), or several different analyses may be run but only one (the one that the laboratory thinks is "right") reported. Consequently, it will be difficult to separate sources of interlaboratory variation.

Even worse, even if all factor levels can be accurately identified (or set in advance), there is the clear possibility of

421

*interaction* among factors, making interpretation of results difficult. For example, different methodologies work "best" in the context of different models, and different models and methodologies lead to different algorithms and thus to different reasons for numerical instability when such algorithms are applied to data.

How do we then interpret Dr. Currie's examples? First, and most important, we see that in certain far more precisely defined contexts (in terms of model) than would be met in practice there are wide variations in conclusions between laboratories, and also wide variations from the correct answer. Second, the divergences between laboratories cannot be assigned to sampling variation because (and this is the beauty of the standard data sets) the data are fixed. However, the divergence of the centroid of the laboratory conclusions from the truth may be due to sampling variation (the sample did not represent the population), or to poor laboratory conclusion-making processes, or to both. We cannot partition this last variability in terms of possible causes, because no accurate measures of precision over sampling variation are provided by the laboratories, or (in some cases) known.

To better understand the sources of interlaboratory variation, we need to start with pilot studies that control the levels of each of the factors (1-3) listed above. To establish the contribution of algorithms to interlaboratory variability, we need to ask numerical analysts to study the possible numerical errors that can occur in algorithms, describe the situations that produce these errors, and suggest remedies to reduce such errors. (Here, our "pilot sample" design fixes all factors but the "algorithm" factor.) To establish the contribution of methodology (or rather methodology - model interactions) to such variability, chemometricians (particularly statisticians) need to use mathematical analysis and simulaton to identify formulas for the precisions (sampling variability) that can be assigned to the various methodologies in the contexts of various models (Here, we assume a fixed, perfectly accurate algorithm, and vary combinations of method and model.) Finally, we need to study and assess the variability due to choice of model, and also to the other "scientific judgments" made by a laboratory in choosing methodology and algorithms and in announcing measures of precision. It is particularly in this last type of study that standard data sets, both real and artificial, can be most useful.

## Scientific Judgment and Bayesian Inference

The question of how to analyze the biases introduced by "scientific judgment" has a direct relationship to a long-standing controversy between classical (frequentist) statisticians and those statisticians who advocate a Bayesian approach. Scientific decision-making involves subjective judgments about both models and types of permissible conclusions. When such judgments are unstated, we have seen that this can obscure our understanding of how decisions are reached, and thus prevent us from finding sources of "bias" or error.

Bayesian statisticians, who try to mathematically model their subjective judgments in terms of prior probabilities over unknown parameters (and models), are often accused by frequentist statisticians of proposing analyses that lack "scientific objectivity." Clearly the contrary is true. The scientist who claims to base conclusions only on the "objective" evidence provided by observed frequencies is nevertheless often guilty of imposing unstated judgments on such evidence. The Bayesian, at least, tries to bring these judgments into the open, where they can be assessed along with the data. Even if we doubt that probability models can ever serve as adequate models of subjective belief, we can still applaud the Bayesian's efforts to expose the methods by which this belief interacts with the evidence in the data to produce new judgments or belief. Rather than criticize the Bayesian for being "subjective" or "biased", the frequentists need to find ways of making their own decision-making processes available for objective study, so that we can gain the opportunity to learn how to improve scientific judgment.

## References

[1] Anderson, E., The Irises of the Gaspe Peninsula. *Bull. Amer. Iris Soc.* **59**, 2-5 (1935).

[2] Fisher, R.A, The use of multiple measurements in taxonomic problems. *Ann. Eugen.* **7**, 179-188 (1936).

[3] Longley, J.W., An appraisal of least squares programs for the electronic computer from the point of view of use. *Journ. Amer. Statist. Assoc.* **62**, 819-841 (1967).

[4] Lachenbruch, Peter A., Review of "STAN, Version II.0." *The American Statistician* **39**, 146-148 (1985).

[5] Efron, B., *The Jackknife, the Bootstrap and Other Resampling Plans*, Soc. Indust. and Appl. Math.: Philadelphia, PA (1982).

[6] Diaconis, P., and B. Efron, Computer intensive methods in statistics. *Scientific American*, May, 116-130 (1983).

[7] Freedman, D.A., and S.C. Peters, Using the bootstrap to evaluate forecasting equations. Tech. Report, Dept. Statistics, U. Calif. Berkeley (1983).

[8] Freedman, D.A., and S.C. Peters, Bootstrapping a regression equation: Some empirical results. *Journ. Amer. Statist. Assoc.* **79**, 97-106 (1984).

[9] Miller, R.G. Jr., *Simultaneous Statistical Inference*, 2nd edit. McGraw-Hill: NY (1981).

# Statistical Properties of a Procedure for Analyzing Pulse Voltammetric Data

## Thomas P. Lane

Massachusetts Institute of Technology, Cambridge, MA 02139

and

## John J. O'Dea and Janet Osteryoung

State University of New York at Buffalo, Buffalo, NY 14214

O'Dea et al. (1983, *J. Phys. Chem.* **97**, 3911–3918) proposed an empirical procedure for obtaining estimates and confidence intervals for kinetic parameters in a model for pulse voltammetric data. Their goal was to find a procedure that would run in real time, not necessarily one that would have well-defined statistical properties. In this paper we investigate some of the statistical properties of their procedure. We show that their estimation method is equivalent to maximum likelihood estimation, and their confidence intervals, while related to likelihood ratio confidence regions, have a coverage probability that is not fixed and that is potentially quite large. We suggest modifications of their procedure that lead to more traditional confidence intervals. We examine the effect on their procedure of the presence of nuisance paramters. Finally we discuss the possibility of serially correlated errors.

Key words: autocorrelation; confidence intervals; estimation method; kinetic parameters; maximum likelihood estimation; serially correlated errors; statistical properties.

## 1. Introduction

O'Dea et al. (1983) proposed a nonlinear regression procedure for estimating, and obtaining confidence intervals for, kinetic parameters describing the reduction of Zn(II) at a stationary mercury electrode in aqueous

---

**About the Authors:** Thomas P. Lane is with the Statistics Center at MIT while John J. O'Dea and Janet Osteryoung are with the Department of Chemistry at the State University of New York at Buffalo.

---

solutions of $NaNO_3$. In this paper we examine the statistical properties of the procedure and suggest modifications to improve these properties.

In section 2 we describe O'Dea's procedure. In section 3 we show his estimation procedure to be equivalent to maximum likelihood estimation. In section 4 we show his interval estimation procedure produces intervals that are related to higher-dimensional confidence regions obtained by likelihood ratio theory, and we suggest modifications to the procedure that will produce confidence intervals with the desired coverage probability. In section 5 we question the assumption of

independent errors and examine the effect of including another parameter in the model to describe the apparent autoregressive error structure.

The notation used here is similar to that used by O'Dea, but as is customary in literature on regression we use $Y$ as the dependent variable.

## 2. Description of the Procedure

O'Dea models the observed response at time $t_i$ to an arbitrary pulse sequence by

$$Y_i = af_i + c + \epsilon_i ,$$

where $a$ and $c$ are unknown constants that convey no kinetic information, $\{\epsilon_i\}$ is a sequence of errors that are assumed to be independent with mean 0 and unknown variance $\sigma^2$. The function $f_i = f(t_i,\alpha,k,E_\frac{1}{2})$ is the solution of an integral equation. It depends on unknown kinetic parameters $\alpha$, $k$, and $E_\frac{1}{2}$, and it must be obtained by solving the integral equation numerically.

The kinetic parameters are of primary interest, so O'Dea uses a nonlinear optimization procedure to find the values of these parameters that maximize the correlation $R$ between $Y$ and $f(t,\alpha,k,E_\frac{1}{2})$. These values are taken as the estimates. Estimates of $a$ and $c$ can then be obtained by simple linear regression of $Y$ on $f(t,\alpha,k,E_\frac{1}{2})$, and $\sigma$ can be estimated by the standard deviation of the residuals from this regression.

With no error the correlation $R$ calculated above would be equal to unity. O'Dea measures the deviation from unity by $\bar{R} = 1 - R$ and defines $\bar{R}_{min}$ as the optimum value of $\bar{R}$. To measure the uncertainty in his estimate of $\alpha$ he fixes $k$ and $E_\frac{1}{2}$ at their optimal values and finds the two values of $\alpha$ that give $R = 3\bar{R}_{min}$. He calls the interval between these values a "confidence interval" for $\alpha$, but he assigns no confidence level to the interval. He computes similar intervals for $k$ and $E_\frac{1}{2}$.

## 3. Maximum Likelihood Estimation

A more traditional approach to a problem of this sort would be to write the likelihood or log likelihood function for the problem and maximize it as a function of the unknown parameters. For normally distributed errors this is equivalent to choosing parameter values that minimze the sum of squared residuals. In this section we show that the above estimation procedure is also equivalent to maximum likelihood estimation.

If $n$ is the number of observations, the log likelihood $L$ is given by

$$L(a,c,\alpha,k,E_\frac{1}{2},\sigma) = -\frac{n}{2}\log(2\pi\sigma^2) - \frac{1}{2}\sum_i \left[\frac{Y_i - c - af_i}{\sigma}\right]^2 .$$

The maximum is easily found by noting that for any value of $\sigma$, the expression is maximized by choosing $a$, $c$, $\alpha$, $k$, and $E_\frac{1}{2}$ to minimize the sum of squared residuals $SS_R = \Sigma(Y_i - c - af_i)^2$. It is a simple matter to show that if $\bar{Y} + n^{-1}\Sigma Y_i$ and $SS_T = \Sigma(Y_i - \bar{Y})^2$, then $SS_R = (1 - R^2)SS_T$, so $SS_R$ is a minimum when $R$ is a maximum (in absolute value). Therefore O'Dea's estimates are the maximum likelihood estimates.

## 4. Confidence Intervals

Confidence regions for unknown parameters are often found by computing the maximum likelihood estimates and then finding other sets of parameter values for which the likelihood function, or an approximation to the likelihood function, is not much smaller. O'Dea's procedure is related to this approach.

Define $L(\alpha,k,E_\frac{1}{2})$ as the maximum over $a$, $c$, and $\sigma$ of the log likelihood $L(a,c,\alpha,k,E_\frac{1}{2},\sigma)$. Using the relationships between $R$, $SS_R$, and $SS_T$ above and maximizing over $\sigma$ gives

$$L(\alpha,k,E_\frac{1}{2}) = (-n/2)(1 + \log(2\pi) + \log SS_T + \log(1 - R^2) - \log n).$$

O'Dea's procedure involves finding six points—the two endpoints of the confidence intervals for each parameter—with the same correlation $R = 1 - 3\bar{R}_{min}$. The above expression for $L(\alpha,k,E_\frac{1}{2})$ shows that these points have the same log likelihood as well.

Let $\hat{\alpha}, \hat{k}$, and $\hat{E}_\frac{1}{2}$ be the maximum likelihood estimates of $\alpha$, $k$, and $E_\frac{1}{2}$. The quantity $\lambda = \exp(L(\hat{\alpha},\hat{k},\hat{E}_\frac{1}{2}) - L(\alpha,k,E_\frac{1}{2}))$ is called the likelihood ratio. It can be shown that $2\log\lambda$ has an asymptotic chi-square distribution with 3 degrees of freedom if $\alpha$, $k$, and $E_\frac{1}{2}$ are the true parameter values. Then $P[2\log\lambda \leqslant 7.815] = 0.95$, so the parameter values for which $2\log\lambda \leqslant 7.815$ form a 95% confidence region for the true values of the parameters. This region is bounded by the roughly ellipsoidal surface $\lambda(\alpha,k,E_\frac{1}{2}) = \exp(3.908)$, as shown in figure 1. In general, any surface of constant $\lambda$ bounds some confidence region.

The confidence level of the region bounded by the surface containing O'Dea's points can be determined by computing $\lambda$. In his procedure

$$2\log\lambda = -n[\log(1 - (1 - \bar{R}_{min})^2) - \log(1 - (1 - 3\bar{R}_{min})^2)]$$
$$\approx -n\log(2\bar{R}_{min}/6\bar{R}_{min}) = n\log 3,$$

since $(\bar{R}_{min})^2 << \bar{R}_{min}$. Here $n = 81$, so $2\log\lambda \approx 89$. By comparison, the region bounded by the surface for
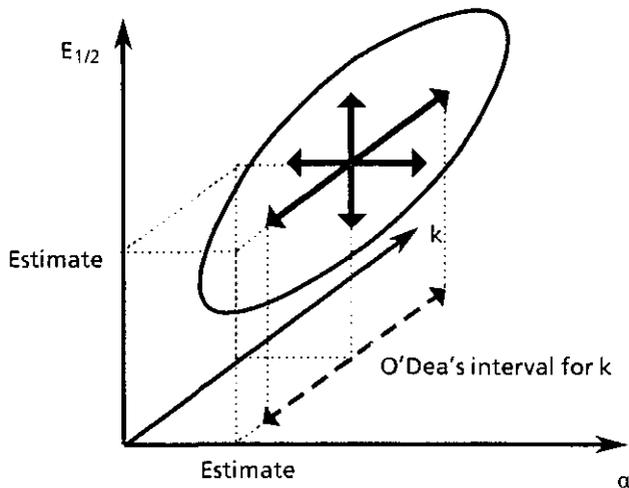
Figure 1–Relationship between O'Dea's intervals and a three-dimensional confidence ellipsoid. The ellipsoid passes through the endpoints of the solid line segments.



Figure 2–Comparison of confidence interval with interval computed by O'Dea's procedure in two dimensions.

which $2\log\lambda = 16.268$ has a confidence level of 99.9%, so a three-dimensional confidence region found using $2\log\lambda = 89$ would be very conservative.

A more customary confidence level is 95%. Since the likelihood ratio can be written as a function of the correlation, it is possible to use a modification of O'Dea's procedure to find points on the boundary of a 95% confidence region. Rather than increasing $\bar{R}$ by a factor of 3, the appropriate factor is the value of $b$ for which 81 $\log b = 7.815$, or $b \approx 1.10$. For example, in a sample data set that does not appear in O'Dea's original paper, $\alpha = .22522$. Increasing $\bar{R}$ by a factor of 3 produces the interval [.22139, .22916], while the factor 1.10 leads to the interval [.22435, .22609].

On the other hand O'Dea's goal was not to find points in a confidence region for all three parameters, but to find separate confidence intervals for each parameter. In order to use the distribution of the likelihood ratio, O'Dea's procedure must be modified so that in computing the endpoints of the confidence interval for one parameter, the likelihood is maximized over the other two parameters. Twice the log of this likelihood ratio has an asymptotic chi-square distribution with one degree of freedom.

In the case of $\alpha$, for example, this is done by comparing $2\log\lambda = 2[L(\hat{\alpha},\hat{k},\hat{E}_{\frac{1}{2}}) - L(\alpha,\tilde{k}(\alpha),\tilde{E}_{\frac{1}{2}}(\alpha))]$ to a chi-square distribution with one degree of freedom, where $\tilde{k}(s)$ and $\tilde{E}_{\frac{1}{2}}(s)$ are the values of $k$ and $E_{\frac{1}{2}}$ that maximize $L(\alpha,k,E_{\frac{1}{2}})$ subject to the restriction $\alpha = s$. The 95% point of this distribution is 3.841, so the values of $\alpha$ for which $2\log \lambda \leqslant 3.841$ form a 95% confidence interval for the true parameter value.

This is best illustrated in two dimensions, as in figure 2. Here approximately elliptical contours of constant $\lambda$
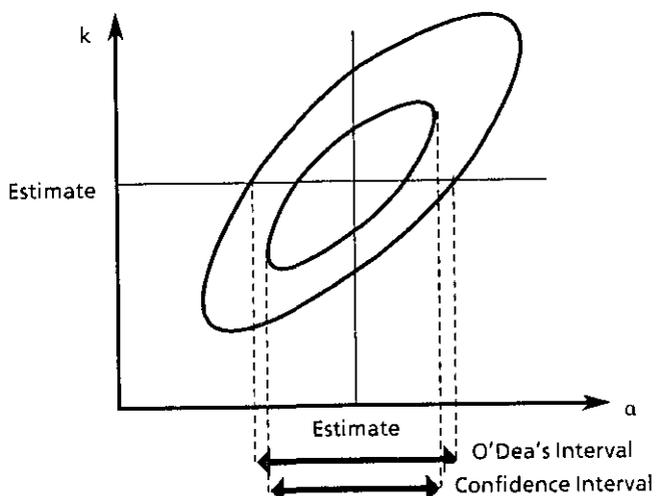
are plotted as a function of $\alpha$ and $k$ for constant $E_{\frac{1}{2}}$. (The complete contours are ellipsoids in three dimensions.) The inner ellipse has $2\log\lambda = 3.841$, while the outer ellipse has $2\log\lambda = 7.815$. The endpoints of the confidence interval for $\alpha$ are the points on the contour that have tangents perpendicular to the $\alpha$ axis. This interval can be compared with the interval found by O'Dea's procedure, which is that portion of the $k = \tilde{k}$ line that is within the outer ellipse.

Which is larger? If the two ellipses have major axes parallel to the coordinate axes, O'Dea's intervals are longer and their coverage probabilities exceed 95%. While this is not desirable, it increases the probability that his interval will contain the true parameters. But if the major axes are not parallel to the coordinate axes and if the lengths of the minor axes are small, O'Dea's intervals are shorter and have a coverage probability less than 95%. Unfortunately it is not possible to determine which is the case by looking only at the points examined in his procedure.

There are two sensible remedies to this problem. The first, the likelihood ratio method, is similar in spirit to O'Dea's original procedure. This method involves finding the confidence interval as described above by finding those values of the first parameter that produce the proper likelihood ratio when the likelihood is maximized over the other two parameters. In this problem, though, it is time consuming to calculate $f_i$ and its derivatives do not have simple analytic expressions, so repeated maximization of the likelihood may be too computationally burdensome.

The other method, an asymptotic normal approximation, is the one we use here. This involves assuming that the maximum likelihood estimates have a multi-

425

variate normal distribution with a mean vector equal to the true parameter values and with covariance matrix equal to minus the inverse of the second derivative of the log likelihood $L$. (This is equivalent to the likelihood ratio method applied to a quadratic approximation to the log likelihood.) Since there is no analytic expression for the second derivative in this problem, we use a numerical approximation.

In the example given above, the estimated covariance matrix is

$$S = \begin{bmatrix} 1.0334 & -.6908 & .0396 \\ -.6908 & 8.6984 & -.7509 \\ .0396 & -.7509 & .1384 \end{bmatrix} \times 10^{-7}.$$

A 95% confidence interval for $\alpha$ is given by $\hat{\alpha} \pm 1.96(S_{11})^{\frac{1}{2}}$, or [.22459, .22585]. This is narrower than the interval obtained above using O'Dea's procedure with a factor of 1.10.

## 5. Residual Autocorrelation

The above derivations are valid if the errors $\{\epsilon_i\}$ are independent normal random variables with a common variance. In practice this assumption must be checked. This is especially true when, as in this case, measurements are taken over time. It is often reasonable to suspect that measurements at neighboring time points may be correlated.

The errors $\{\epsilon_i\}$ are not observed, but they can be estimated by the residuals, or the differences between the observed $Y_i$ and the fitted values $\hat{Y}_i = \hat{c} + \hat{a} f(t_i, \hat{a}, \hat{k}, \hat{E}_{\frac{1}{2}})$. Figure 3 is a plot of $Y$ and $\hat{Y}$ as a function of time. Figure 4 is a plot of the residuals $e_i = Y_i - \hat{Y}_i$ over time. If the residuals were independent we would not expect to find any pattern here, but in fact there is a pronounced tendency for residuals at neighboring time points to have the same sign.

There are three possible causes for this phenomenon. First, it is possible that this is an artifact of the fitting procedure. Even when the errors are independent, fitting an ordinary linear regression produces residuals that have some correlation (for example, they sum to zero). It is possible that minimizing the sum of squared residuals in this more complicated model produces residuals with some autocorrelation. However experiments with the model do not support this hypothesis.

A second possible cause is model inadequacy. The model relates an imposed voltage to an observed current, and the voltage is highly correlated with time. If the equation used here is not the true relationship between the current and the voltage, there may be correlation between the current and the residuals, and this dependence could be masquerading as time dependence. The cure for this difficulty is to propose alternative models that better fit the data.
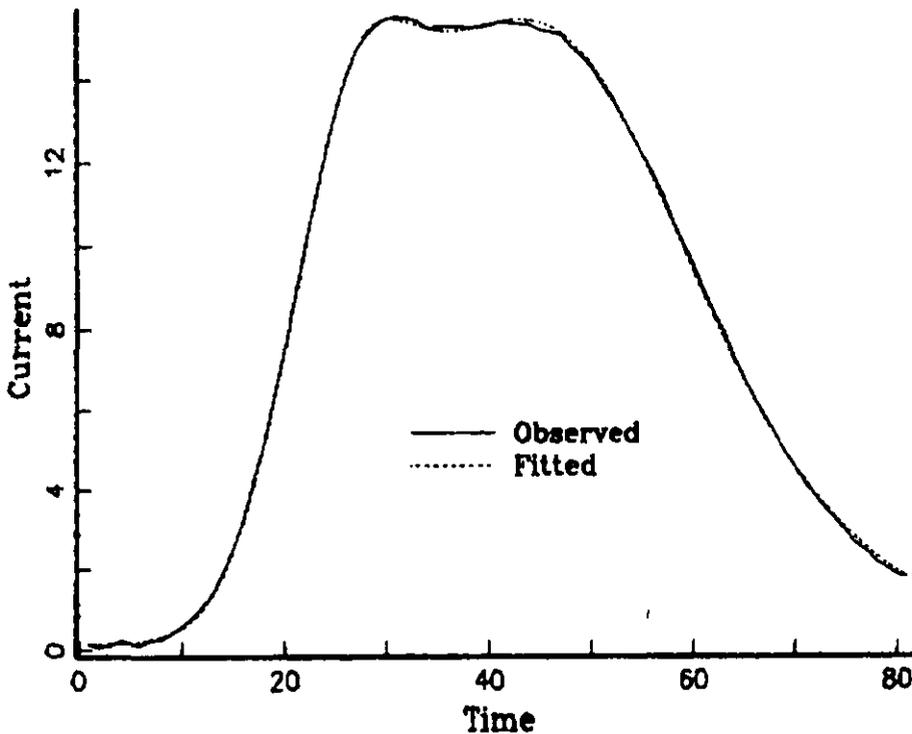


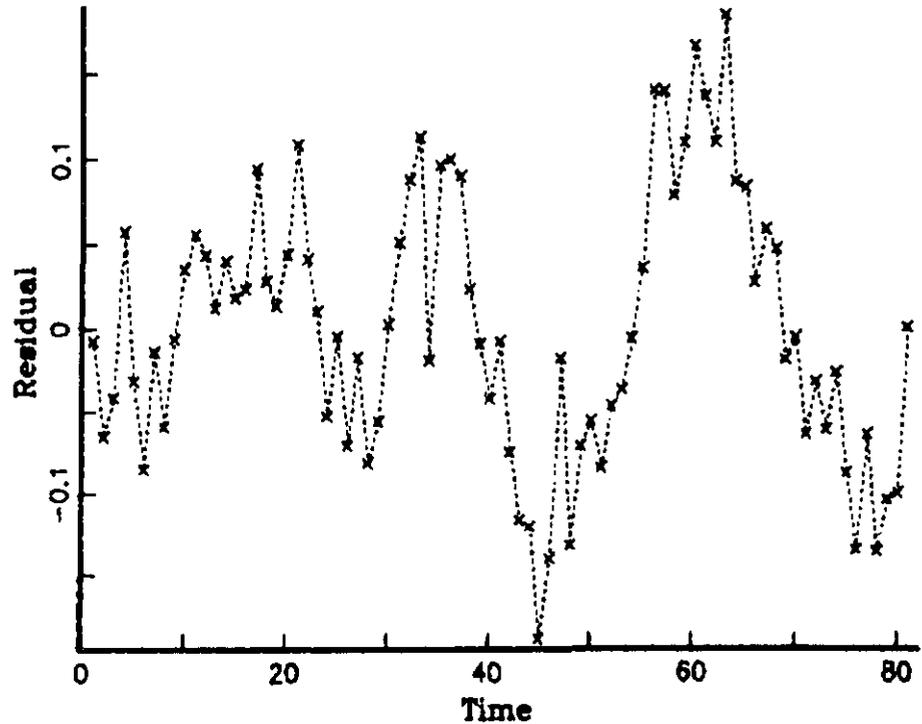Figure 3–Observed and Fitted Current.

Figure 4–Residuals as a Function of Time.

The third possible cause is actual autocorrelation in the errors, and this autocorrelation can be modeled as well. We proceed under the assumption that the true errors have time-dependent correlation.

Two common models for time series are the first order autoregressive model

$$\epsilon_i = \rho\, \epsilon_{i-1} + u_i$$

and the first order moving average model

$$\epsilon_i = u_i + \rho\, u_{i-1},$$

where in both cases $\{u_i\}$ is a sequence of independent normal random variables with mean 0 and common unknown variance, and $\rho$ is an unknown parameter between $-1$ and 1. Other possible models are the higher order models, where terms from earlier time points are used, and mixed models, where $\epsilon_i$ is modeled as a linear combination of $\epsilon_{i-1}, ..., \epsilon_{i-p}$ and $u_i, ..., u_{i-q}$.

Two tools useful for identifying a good model are the autocorrelation function and the partial autocorrelation function. These appear in figures 5 and 6. The sample autocorrelation function is simply the correlation of $e_i$ and $e_{i-k}$ plotted as a function of $k$. For a moving average process of order $q$ the true autocorrelation function is 0 for $k > q$. For autoregressive processes and mixed processes the true autocorrelation function approaches 0 as $k \to \infty$, but it is not identically 0 for all $k$ beyond some finite value. The sample autocorrelation function in fig-

ure 5 seems to be more consistent with that of the autoregressive and mixed models, since there does not seem to be a sharp cutoff.

The partial autocorrelation function is more complicated, but its interpretation is quite simple. It is the "dual" of the autocorrelation function, in that it is 0 for all $k > p$ for an autoregressive process of order $p$, and it approaches 0 as $k \to \infty$ but it does not vanish for moving average and mixed processes. The sample partial autocorrelation function in figure 6 shows a large value at $k = 1$ and smaller values for $k > 1$. It is never exactly 0, but for most $k$ values the sample partial autocorrelation falls inside the boundary that marks the values that are significantly different from 0. The function seems to be consistent with what might be expected from a first order autoregressive process.

The new model

$$Y_i = af_i + c + \epsilon_i \qquad \text{with} \qquad \epsilon_i = \rho\epsilon_{i-1} + u_i$$

is equivalent to

$$Y_i = \rho Y_{i-1} + a(f_i - \rho f_{i-1}) + c(1-\rho) + u_i,$$

which is a nonlinear regression model with independent errors.

There are now four parameters to be estimated, in addition to $a$, $c$, and $\sigma$. But if only the original three parameters are of interest, it is possible to treat $\rho$ as one of the nuisance parameters by using a variant of the Cochrane-Orcutt procedure, as follows:

427

Figure 5–Residual Autocorrelation Function.

········· **Boundary of 0.05 Rejection Region**



Figure 6–Residual Partial Autocorrelation Function.

········· **Boundary of 0.05 Rejection Region**

1)  For any given $\alpha$, $k$, and $E_{\frac{1}{2}}$, compute $\{f_i\}$.
2)  Estimate $a$ and $c$ by linear regression to get $\{e_i\}$.
3)  Estimate $\rho$ by the sample correlation of the $\{e_i\}$.
4)  Regress $Y_i - \rho Y_{i-1}$ on $f_i - \rho f_{i-1}$ to get new estimates of $a$ and $c$, and new residuals $\{e_i\}$.
5)  Repeat steps 3 and 4 until convergence.
6)  Compute the sum of squares $\Sigma u_i^2 = \Sigma(e_i - \rho e_{i-1})^2$.

The computer time needed for these steps is much less than that needed to compute $\{f_i\}$, so the estimation is much faster if the nonlinear optimization program searches only in the three-dimensional space of $(a, k, E_{\frac{1}{2}})$. For each set of trial parameter values the above steps can be performed to minimize the residual sum of squares over the nuisance parameters. The resulting estimate of $\alpha$ is .22473.

The other calculation can also be repeated for this new model. The estimated covariance matrix is

428

$$S = \begin{bmatrix} 4.8344 & -5.6256 & .4766 \\ -5.6256 & 37.5982 & -2.6193 \\ .4766 & -2.6193 & .6513 \end{bmatrix} \times 10^{-7}.$$

This is roughly four times the previous covariance matrix, so the length of the new confidence interval is about twice that of the previous confidence interval. The new interval is [.22336, .22600].

The four intervals around $\alpha$ are compared in figure 7. The procedure used in O'Dea's original paper produces an interval obtained from a three dimensional confidence ellipsoid with a very large confidence level, and it is quite long. The interval is shortened by using a 95% confidence ellipsoid, but it still does not have a 95% coverage probability. The 95% confidence interval is still shorter. Taking into account the apparent autoregressive error structure leads to a confidence interval that is about twice as long as the one in the independence model, but still only a third as long as the interval obtained by O'Dea's procedure.

## References

[1]  O'Dea, J. J.; J. Osteryoung and R. A. Osteryoung, Square wave voltammetry and other pulse techniques for the determination of kinetic parameters. The reduction of zinc (II) at mercury electrodes, *J. Phys. Chem.* **87**, 3911–3918 (1983).

[2]  Morrison, D. F. *Multivariate Statistical Mehods*, McGraw Hill: New York (1976).



Figure 7-Confidence Intervals for Alpha.

# DISCUSSION

of the Lane-O'Dea-Osteryoung paper,
Statistical Properties of a Procedure
for Analyzing Pulse Voltammetric Data

## Janet Osteryoung

Department of Chemistry
State University of New York at Buffalo

To supplement Lane's discussion of a statistical procedure for analyzing pulse voltammetric data, I would like to describe the experiment more fully in the context of the scientific problem being addressed. In a controlled-potential (voltammetric) experiment the current response generally depends on both potential and time. Since the current is the rate of charge transfer, the results of such experiments can be analyzed to yield the values of parameters that characterize the charge transfer process. The current for many charge transfer mechanisms can be calculated, although often the results (a current-potential curve, or voltammogram) can be obtained only numerically. The calculated voltammogram must be inverted to yield the values of the charge transfer parameters. It was the objective of O'Dea, et al. [1][1] to devise a procedure for this inversion that would not depend on either the charge-transfer mechanism or the choice of voltammetric experiment. However the specific problem addressed was that of determining the charge transfer parameters for the reduction of Zn(II) at mercury electrodes in aqueous solutions of $NaNO_3$.

It was assumed at the outset that the mechanism of charge transfer was described by the Butler-Volmer equation

$$i(t) = nFAk\epsilon^{-\alpha}[D_O^{1/2}C_O(O,t) - \epsilon D_R^{1/2}C_R(O,t)]$$

where $i(t)$ is the current at time $t$, $n$ the number of electrons per zinc ion reduced, $F$ the value of the Faraday, $A$ the electrode area, $D_O$ and $D_R$ the diffusion coefficients of the oxidized (O) and reduced (R) forms of zinc. $C_O(O,t)$ and $C_R(O,t)$ are the corresponding concentrations at the electrode surface at time $t$,

$$\epsilon = \exp[(nF/RT)(E(t) - E_{1/2})]$$

where R is the gas constant, T the absolute temperature, E(t) is the imposed potential, which is a function of time, and $k$, $\alpha$, and $E_{1/2}$ are the kinetic parameters as given by

Lane. A mathematical model which describes exactly the current-potential relation is developed by formulating the diffusion problem with the Butler-Volmer relation as a boundary condition and expressing the surface concentrations in terms of convolution integrals of currents to yield an integral equation for the current which can be solved numerically.

The typical procedure used in electrochemical kinetic studies is to measure $E_{1/2}$ independently and to use its value as a known quantity in analysis of the voltammogram. Furthermore, usually $D_O$, $D_R$, $A$, and $n$ are determined in order to compare calculated and experimental currents directly. These additional pieces of information are not necessary, however, and may introduce systematic error into the values of the derived kinetic parameters. Because of the exponential form of the current-potential relation, minor errors in the value of $E_{1/2}$ distort the shape of the response and therefore cause errors in the derived value of $\alpha$. These errors can even suggest a potential-dependence of $\alpha$ which is an artifact. In a differential experiment such as square wave voltammetry the response is generally peak-shaped, and the height of the peak reflects the values of $k$ and $\alpha$. Errors of normalization (e.g., measurement of $A$) therefore also introduce error into the values of the derived parameters.

From the discussion of Lane et al., it is clear that the normalization factor, $a$, is an unnecessary "nuisance" parameter, and thus it is foolish to confound the results of kinetic measurements by employing a method of data analysis which requires that $a$ be known. The question of $E_{1/2}$ is more subtle, for experimental and chemical factors must be considered. In principle, knowing the true value of $E_{1/2}$ simplifies the problem. Potential differences can be measured accurately, but it is difficult to maintain a laboratory reference potential at a known value over time. The data of [1] and the data employed by Lane et al. display confidence intervals for $E_{1/2}$ at the 95% level of $\leqslant 0.001$ V. Working laboratory standards are not maintained with that precision. Chemical factors

---

[1] Number in bracket is literature reference.

must also be considered. The value of $E_{1/2}$ is measured using either a voltammetric experiment with a much longer time scale or an equilibrium experiment. In the latter case the diffusion coefficients must also be known to yield $E_{1/2}$. In either case the change in time scale introduces the possibility of a change in mechanism, which produces a value of $E_{1/2}$ inappropriate for the conditions of the kinetic experiment. Therefore $E_{1/2}$ should be treated as an unknown parameter of the experiment together with $k$ and $\alpha$.

A further objective of this work was to obtain confidence intervals for the expectation values of the kinetic parameters. Typically in experiments of this type uncertainty is estimated by estimating the coefficient of variation of the current response and assigning that coefficient of variation to the derived parameters, because more realistic procedures have not been available. The procedure presented by O'Dea et al. [1] has the merit of computational simplicity and thus provides a well-defined quantity that can be used by the experimenter as a working figure of merit during the course of experiments. The procedure of Lane et al., which provides conventional confidence intervals with known confidence bound, relies on quadratic approximation of the model, which should be adequate for well-behaved response surfaces.

Recent advances in theory and in computational capabilities raise the possibility of fully quantitative theoretical descriptions of at least some classes of electrochemical reactions. Theoretical developments can be guided and tested using accurate data that have been analyzed using appropriate statistical techniques.

## References

[1] O'Dea, John J.; Janet Osteryoung, and R. A. Osteryoung, J. Phys. Chem. 1983, 87, 3911–3918.

# Fitting First Order Kinetic Models Quickly and Easily

## Douglas M. Bates

University of Wisconsin-Madison, Madison, WI 53705

and

## Donald G. Watts

Queen's University at Kingston, ON, Canada

Kinetic models described by systems of linear differential equations can be fitted to data quickly and easily by taking advantage of the special properties of such systems. The estimation situation can be greatly improved when multiresponse data are available, since one can then automatically determine starting values and better discriminate between rival models.

Key words: compartment model; determinant criterion; multiresponse estimation.

## 1. Introduction

In this article we summarize the work of a series of papers [1–3][1] in which we deal with fitting first order kinetic models to uniresponse and multiresponse data.

[1] Numbers in brackets indicate literature references.

We consider systems in which the expected responses at $K$ points in the system, $\eta(t) = (\eta_1(t), \eta_2(t), \ldots \eta_K(t))'$, are described by the system of linear differential equations

$$\partial\eta/\partial t = \dot{\eta}(t) = A\,\eta(t) + \iota(t) \qquad (1.1)$$

where $A$ is a $K \times K$ system *transfer matrix* depending on *rate constants* $\theta_1, \theta_2, \ldots$, and $\iota(t)$ is a vector input function to the system. We assume further that there are $K$ initial conditions $\eta_0 = (\eta_{01}, \eta_{02}, \ldots, \eta_{0K})'$, some possibly unknown, and that $\tau = t - \theta_0$, where $\theta_0$ is a (possibly unknown) time delay. All the unknown parameters are gathered into a $P \times 1$ parameter vector $\theta$.

**Example: Oil Shale.1**

As an example of a chemical system described by a set of linear differential equations, we cite the pyrolysis of oil shale in which the model, fitted by Ziegel and Gorman [4] has the system diagram



In this system, $\eta_1$ denotes kerogen, $\eta_2$ bitumen, and $\eta_3$ oil. The model implies that kerogen decomposes to bitumen with rate constant $\theta_1$, and to oil with rate constant $\theta_4$, and bitumen produces oil with rate constant $\theta_2$, and unmeasured by-products with rate constant $\theta_3$.

## 2. Expectation Functions and Derivatives With Respect to the Parameters for First Order Kinetic Systems

Several methods are used to estimate the parameters in first order kinetic models. The most obvious method is to solve the system of differential equations corresponding to the particular compartment model and use the resulting expectation function in a standard nonlinear estimation program. A second approach is to fit a general sum of exponentials model by "peeling" [5]. A third approach is to use a standard nonlinear estimation program, using numerical integration to solve the equations. A superior approach, proposed by Jennrich and Bright [6], is to obtain the general solution to the system of equations by calculating values for the model function $\eta(t)$ and its derivatives directly, given values of $\theta$ and $t$ and $\iota(t)$.

### 2.1 The General Solution

The solution to a linear system of differential equations can be expressed in terms of convolutions using the matrix exponential [7]. The solution is

$$\eta(t) = e^{At}\eta_0 + e^{At}*\iota(t) \qquad (2.1)$$

where the '*' denotes convolution,

$$e^{At}*\iota(t) = \int_0^t e^{A(t-\xi)}\iota(\xi)\,d\xi. \qquad (2.2)$$

The integral of the vector function is evaluated componentwise. Computational methods for evaluating the convolution integral are given in Moler and Van Loan [8], when $A$ is diagonalizable, and in Bavely and Stewart [9], when $A$ is nondiagonalizable.

### 2.2 Derivatives of the Expectation Function

To use a Gauss-Newton procedure to estimate the parameters we need derivatives with respect to the parameters. As shown by Jennrich and Bright [6], a great advantage to compartment models is that the derivatives can be evaluated in the same fashion as the model function itself. Instead of differentiating $\eta(t)$ directly as in Jennrich and Bright, however, we differentiate eq (1.1) and solve the resulting linear system of differential equations. This idea was discussed in another context in Smith [10] and was used by Kalbfleisch et al. [11].

To simplify notation, we use a subscript $p$ to denote differentiation with respect to the parameter $\theta_p$ so

$$\eta_p(t) = \frac{\partial \eta(t)}{\partial \theta_p}$$

for $p = 1, 2, \ldots P$. The derivative of (1.1) with respect to $\theta_p$ is then

$$\dot{\eta}_p(t) = A\eta_p(t) + A_p\eta(t) + \iota_p(t)$$

for which the solution is

$$\eta_p(t) = e^{At}\eta_p(0) + e^{At}*[A_p\eta(t) + \iota_p(t)] \qquad (2.3)$$

$$= e^{At}\eta_p(0) + e^{At}*\iota_p(t) + e^{At}*A_pe^{At}\eta_0 + e^{At}*A_pe^{At}*\iota(t).$$

Dead time, $\theta_0$, can be incorporated by modifying (1.1) to

$$\dot{\eta}(\tau) = A\eta(\tau) + \iota$$
$$\eta(0) = \eta_0$$

where

$$\tau = \begin{array}{ll} t - \theta_0 & t > \theta_0 \\ 0 & t \le \theta_0 . \end{array}$$

If $\theta_0$ is known, we simply replace $t$ by $\tau$ in eqs 1.1, 2.1, and 2.3, but if $\theta_0$ is unknown and depends on a parameter, the expression for the derivatives is extended to

$$\eta_p(\tau) = e^{A\tau}\eta_p(0) + e^{A\tau}*[A_p\eta(\tau) + \iota_p]$$

$$+ \tau_p[A\eta(\tau) + \iota].$$

It is easy to evaluate $A_p$, $\eta_p(0) = \partial\eta o/\partial\theta_p$, and $\tau_p$ since they are constants and, in fact, usually two of the three are zero. Also, for any $t < \theta o$, $\tau_p$ is zero for all $p = 1, \ldots, P$. Note that the method can be extended to higher order derivatives.

### Specifying the Model

A simple unambiguous computer notation can be used to specify first order kinetic models in a *parameter table* consisting of three columns, the first column giving the parameter number. For a rate constant, the second column entry gives its source and the third column entry its sink, a sink with compartment number 0 denoting elimination. For initial conditions, $\eta o$, the second column entry is the number of the component in $\eta o$ and the third column entry is $-1$. For a step input, $\iota$, the second column entry is the number of the component in $\iota$ and the third column entry is $-2$. Dead time is coded as 0 in column 2 and 0 in column 3.

### Example: Oil Shale.2

The oil shale parameter table is presented in annotated form. It is to be noted that a single parameter may represent more than one rate constant.

Table 1. Oil shale parameters.

| Parameter number (type) | Column 2 | | Column 3 | |
|---|---|---|---|---|
| 1 (rate constant) | 1 | (source) | 2 | (sink) |
| 2 (rate constant) | 2 | (source) | 3 | (sink) |
| 3 (rate constant) | 2 | (source) | 0 | (elimination) |
| 4 (rate constant) | 1 | (source) | 3 | (sink) |
| 5 (dead time) | 0 | | 0 | |

## 3. Multiresponse Estimation

In the multiresponse situation when the errors have unknown variances and covariances but are assumed to be temporally uncorrelated, the appropriate criterion derived via a likelihood or Bayesian approach is to minimize the $M \times M$ determinant [12].

$$|V(\theta)| = |Z'Z| \qquad (3.1)$$

In eq.(3.1), $Z = Y - H$ is the $N \times M$ matrix of residuals $\{z_{nk}\} = \{y_k(t_n) - \eta_k(t_n)\}$, $k = 1, \,2\,, \ldots, M$, $n = 1, 2, \ldots, N$. The *expected responses* $\eta$ are assumed to depend on $P$ parameters $\theta$: except where explicitly required, we suppress the dependence on $\theta$. For first order kinetic systems with all responses measured, $M = K$.

To determine the minimum of $|Z'Z|$, it is advantageous to calculate the gradient and Hessian and exploit Gauss-Newton optimization techniques. Efficient numerical procedures for computing the gradient and approximate Hessian are given in [2], and an algorithm which performs the calculations in [1].

Alternative expressions for the components of the gradient $\gamma$ and the Hessian $\Gamma$ are

$$\gamma_p = \partial |V| / \partial\theta_p = 2 |V| tr[V^{-1}Z'Z_p],$$
$$p = 1,2,\ldots,P, \qquad (3.2)$$

and

$$\Gamma_{pq} = \partial^2 |V| / \partial\theta_p \partial\theta_q = \qquad (3.3)$$

$$= \gamma_p\gamma_q + 2 |V| tr[V^{-1}Z'Z_qV^{-1}Z'Z_p]$$
$$+ 2 |V| tr[V^{-1}Z'Z_qV^{-1}Z'_p Z]$$

$$+ 2 |V| tr[V^{-1}Z'_pZ_q] + 2 |V| tr[V^{-1}Z'Z_{pq}],$$
$$p,q = 1,2,\ldots,P.$$

The second derivative terms $Z_{pq}$ in eq (3.3) are ignored to produce an approximate Hessian.

## 4. Practical Aspects

### Linear Constraints

Sometimes the data matrix $Y$ involves dependencies as a result of imputation of responses or mass-balance calculations. If these dependencies also occur in the expected responses, then important modifications to the multiresponse estimation procedure must be made so as to avoid convergence to spurious optima [13,14]. It is therefore necessary to examine the residual matrix $Z(\theta)$ for singularities, which can be done by arranging the rounding units in the columns of $Y$ to be approximately equal and taking a singular value decomposition of $Z$ [15]. As explained there, singular values on the order of the rounding unit indicate singularity and should prompt the analyst to search for constraints in the data. Such examinations should be done at the beginning of the analysis using the initial parameter values and at the end of the analysis using the converged values. To aid convergence, logarithms of parameters are used during estimation.

Linear constraints can be dealt with easily by combining the linear constraint vectors into a matrix, performing a $QR$ decomposition of that matrix, and letting the rotation matrix $W$ be the columns of $Q$ which are orthogonal to the constraint vectors. We then simply minimize $|(ZW)'(ZW)|$, where $ZW = YW - HW$. Clearly, the gradient and Hessian of this determinant

are obtained from eqs 3.2 and 3.3 by replacing $Z$ by $ZW$ and $Z_p$ by $Z_pW$.

## Constraints on the Number of Parameters, Responses and Observations

The determinant criterion implies two constraints on the number of observations [2,3]. First, $N$ must be at least equal to $M$ since otherwise the determinant is identically zero. Second, $N$ must exceed $P$ otherwise the criterion can be made zero by fitting any one response perfectly, which can generate up to $M$ distinct minima. Thus the residual matrix has effectively $N-P$ degree of freedom. It may seem that there should be more degrees of freedom since there are $NM$ separate observations, but the criterion can be locally controlled by any one response so the effective number of observations is $N$ rather than $NM$.

## Starting Values

An important part of fitting nonlinear models is determining good starting values. For uniresponse data, an effective method is to use peeling in which we plot the logarithm of the response versus time and fit a straight line to the segment at large $t$ values. The slope of the line gives an estimate of the smallest eigenvalue of the $A$ matrix. Using the fitted line to generate residuals and plotting the logarithm of the residuals versus $t$ should again reveal a straight line portion at large $t$ values, so the process is repeated, thereby obtaining estimates for the eigenvalues. As mentioned in section 2, this process is often used for parameter estimation, but we do not recommend it.

In the case of multiresponse data for first order kinetics, the problem is easily solved using linear least squares by exploiting the linear relation between the rates and the responses! As noted in [3], if we could measure the rates $\dot{\gamma}$ and the responses $\gamma$ at a particular time $\tau$, then using $\dot{\gamma}(\tau)=A\gamma(\tau)$ produces a linear relation between the "dependent" variable $y=\dot{\gamma}$ and the "independent" variables $x_p = A_p\gamma$ in the form $y = X\theta$. We can thus solve for $\theta$ by using linear least squares. A simple procedure for obtaining starting values, then, is to use approximate rates from finite differences of the responses at successive time points and $x_p$ values from the corresponding averages. Alternatively, one could smooth the data for each response by fitting splines so as to obtain better rate and response values, and then use these in a linear least squares routine.

### Example: α-pinene.1

Data on the thermal isomerization of α-pinene at $189.5°$ were reported by Fuguitt and Hawkins [16], and

a thorough multi-response analysis was presented in [13]. The fitted model was described by the system

$$\dot{\gamma} = \begin{bmatrix} \dot{\gamma}_1 \\ \dot{\gamma}_2 \\ \dot{\gamma}_3 \\ \dot{\gamma}_4 \\ \dot{\gamma}_5 \end{bmatrix} = \begin{bmatrix} -\theta_1\gamma_1-\theta_2\gamma_1 \\ \theta_1\gamma_1 \\ \theta_2\gamma_1-\theta_3\gamma_3-\theta_4\gamma_3+\theta_5\gamma_5 = A\gamma \\ \theta_3\gamma_3 \\ \theta_4\gamma_3-\theta_5\gamma_5 \end{bmatrix}$$

which can also be written $\dot{\gamma}=X\theta$, where

$$X = \begin{bmatrix} -\gamma_1 & -\gamma_1 & 0 & 0 & 0 \\ \gamma_1 & 0 & 0 & 0 & 0 \\ 0_1 & \gamma_3 & -\gamma_3 & -\gamma_5 & \gamma \\ 0 & 0 & \gamma_3 & 0 & 0 \\ 0 & 0 & 0 & \gamma & -\gamma_5 \end{bmatrix}$$

Substituting estimated rates for each time and joining them into a single vector y, and calculating $X$ matrices for each time and joining them into a single $X$ matrix, allows us to use linear regression to estimate starting values for $\theta$. The starting values obtained are listed in column 2, table 2. Note their closeness to the converged values, column 4.

Table 2.   Estimation for α-pinene at 189.5°.

| Parameter | Start | Box et al. | Bates and Watts Full Model | Bates and Watts Reduced Model |
|---|---|---|---|---|
| 1 | 5.84 | 5.95 | 5.94 | 5.95 |
| 2 | 2.65 | 2.85 | 2.84 | 2.82 |
| 3 | 1.63 | 0.50 | 0.45 | – |
| 4 | 27.77 | 31.5 | 31.21 | 30.75 |
| 5 | 4.61 | 5.89 | 5.79 | 5.72 |
| Determinant | 600 | | 28.4 | 29.0 |

Parameter estimates are listed in table 2 in minutes$^{-1}\times 10^{-5}$; that is a table value of 5.84 is actually a rate constant of $5.84\times 10^{-5}$ minutes$^{-1}$.

When some of the responses are not measured, it is still possible to use approximate rates provided other information, such as a mass-balance, is substituted.

## 5.   Two Examples

### 5.1   Oil Shale

The model and data presented by Ziegel and Gorman [4] were fitted using the procedures described here. In this problem the concentration of $\eta_1$ was not measured, which introduces some complexity in determining start-

436

ing values. Second, the elimination compartment, corresponding to coal and gas, was not measured. Third, there was a time delay caused by the shale having to reach the reaction temperature. The observed responses are $y_2$ and $y_3$, measured in percent of the initial kerogen $\eta_1$, so $K=3$, $M=2$, with $N=14$.

To determine starting values in this case, we plotted the data and obtained a rough estimate of $\theta_0=6$ min. Estimating the initial slopes of $\eta_2$ and $\eta_3$ from the graph gave values of 0.029 min$^{-1}$ for $\theta_1$ and 0.013 min$^{-1}$ for $\theta_4$. This delay estimate and the rate constants were used to obtain starting estimates of $\theta_2$ and $\theta_3$ as shown in column 2, table 3, following the procedure of section 4. The final results, together with those of Ziegel and Gorman, are shown in columns 3 and 4 of table 3.

Table 3. Parameter estimates for oil shale data.

| Parameter | Start | Bates-Watts | Ziegel-Gorman |
|---|---|---|---|
| 1 | 0.029 | 0.0172 | 0.0173 (k$_1$) |
| 2 | 0.012 | 0.0090 | 0.0092 (k$_2$f$_2$) |
| 3 | 0.028 | 0.0205 | 0.0201 (k$_2$(t−f$_2$)) |
| 4 | 0.013 | 0.0104 | 0.0104 (k$_3$) |
| $\theta_0$ | 6.0 | 7.8 | 7.7 |
| Determinant | | | |
| | 67980 | 429 | |

### 5.2  α-pinene

In their analysis of the α-pinene data, Box et al. [13] noted that response 4 was imputed using $y_4=0.03(100-y_1)$, and that the data set was subject to a mass-balance constraint, $y_1+y_2+y_3+y_4+y_5=100$. To avoid convergence to spurious optimal parameter values, they recommended that these data dependencies be taken into account by using observation vectors consisting of linear combinations of $y_1$, $y_2$, $y_3$, $y_4$ and $y_5$ which are orthogonal to the space defined by the vectors $(0.03, 0, 0, 1, 0)'$ and $(1, 1, 1, 1, 1)'$. We therefore treat $y_4$ as an unmeasured component for estimation purposes, and use linear combinations of the responses which are orthogonal to the vectors $a_1'=(0, 0, 0, 1, 0)$ and $a_2'=(1, 1, 1, 1, 1)$. The rotation matrix $M$ and the modified responses can therefore be determined by performing a $QR$ decomposition on the matrix $(a_1, a_2)$ and using the last 3 columns of $Q$ coupled with all the responses. In this case, $K=5$, $M=3$, with $N=8$.

Approximate 95% confidence limits for $\ln\theta_3$ were very wide, suggesting that $\theta_3$ was badly estimated and could be zero. We therefore fitted a reduced model in which there was no path from $\eta_3$ to $\eta_4$, see column 5. The change in the determinant is 0.6 on 1 degree of freedom, which, when compared with the scaling factor $s^2=28.4/3=9.46$ on 3 degrees of freedom, is clearly

small, verifying that the reduced model is adequate for this data set.

To further substantiate the adequacy of the reduced model, we fitted both models to a second set of data taken 204.5° [16]. The results of this fitting procedure are presented in table 4. The reduced model appears to be adequate for both data sets.

Table 4. Estimation for α-pinene at 204.5°.

| Parameter | Start | Full Model | Reduced Model |
|---|---|---|---|
| 1 | 23.0 | 22.6 | 22.6 |
| 2 | 13.2 | 12.6 | 12.6 |
| 3 | 8.3 | 0.02 | – |
| 4 | 76.9 | 72.8 | 72.8 |
| 5 | 16.0 | 15.6 | 15.6 |
| Determinant | | | |
| | 116 | 0.55 | 0.55 |

## 6.  Conclusions

Several advantages of the direct multiresponse estimation approach for systems of differential equations are apparent. First, the model can be specified directly from the network diagram. Second, there is no need to obtain the analytic solution to the differential equations describing the reactions. Third, there is no need to code the model functions in a nonlinear estimation routine. Fourth, the bothersome and error-prone step of obtaining and coding derivatives of the expected responses with respect to the parameters is eliminated. Fifth, excellent starting values can be determined automatically.

## References

[1] Bates, D. M., and D. G. Watts, A multi-response Gauss-Newton algorithm, Commun. Statist.—Simul. Comput. B(13), 705-715 (1984).

[2] Bates, D. M., and D. G. Watts, A generalized Gauss-Newton procedure for multi-response parameter estimation, SIAM Journal on Scientific and Statistical Computing, to appear.

[3] Bates, D. M., and D. G. Watts, Multiresponse estimation with special application to systems of linear differential equations, Technometrics 27, 190-201 (1985).

[4] Ziegel, E. R., and J. W. Gorman, Kinetic modelling with multi-response data, Technometrics 22, 139-151 (1980).

[5] Anderson, D. H., Compartmental modeling and tracer kinetics, Springer-Verlag (1983).

[6] Jennrich, R. I., and P. R. Bright, Fitting systems of linear differential equations using computer generated exact derivatives, Technometrics 18, 385-392 (1976).

[7] Noble, B., Applied linear algebra, Prentice Hall (1969).

[8] Moler, C., and C. Van Loan, Nineteen dubious ways to compute the exponential of a matrix, SIAM Review 20 (1978).

[9]  Bavely, C. A., and G. W. Stewart, An algorithm for computing reducing subspaces by block diagonalization, SIAM Journal of Numerical Analysis 16, 359-367 (1979).

[10]  Smith, W. R., Parameter estimation in nonlinear models of biological systems, Fisheries and Marine Services Technical Report 889, Fisheries and Environment Canada (1979).

[11]  Kalbfleisch, J. D.; J. F. Lawless and W. M. Vollmer, Estimation in Markov models from aggregate data, Biometrics 39, 907-919 (1983).

[12]  Box, G. E. P., and N. R. Draper, The Bayesian estimation of common parameters from several responses, Biometrika 52, 355-365 (1965).

[13]  Box, G. E. P.; W. G. Hunter, J. F. MacGregor, and J. Erjavec, Some problems associated with the analysis of multiresponse data, Technometrics 15, 33-51 (1973).

[14]  McLean, D. D.; D. J. Pritchard, D. W. Bacon, and J. Downie, Singularities in multiresponse modelling, Technometrics 21, 291-298 (1979).

[15]  Dongarra, J. J.; J. R. Bunch, C. B. Moler, and G. W. Stewart, Linpack User's Guide, Chap. 11, Philadelphia: S.I.A.M. (1979).

[16]  Fuguitt, R. E., and J. E. Hawkins, Rate of thermal isomerization of α-pinene in the liquid phase, J. Amer. Chem. Soc. 69, 319-312 (1947).

# DISCUSSION

of the Bates-Watts paper,
Multiresponse Estimation With Special
Applications to First Order Kinetics.

## Michael Frenklach

Department of Materials Science and Engineering
Pennsylvania State University [1]

The authors presented an interesting approach to parameter estimation for first order kinetic systems. The method is user oriented and particularly suited for computer implementation as a "canned" program. Indeed, present chemical kinetic codes input reaction mechanism in a natural chemical language, that is, specifying reactions (usually in unformatted READ routines) as they are conventionally written on the paper. This information is automatically converted to a so-called reaction matrix and, based on it, to differential equations describing the kinetics of reaction species. The reaction matrix, which contains all the stoichiometry of the system, can conveniently provide the required input infor-

Another important feature, from the user's point of view, is that the presented method is applicable to multiresponse data. It should be realized that modern problems of interest to chemical kinetics get tougher, as for example, formation of pollutants in hydrocarbon combustion. The experimental answer to the growing complexity of the systems is the employment of multiple diagnostics for simultaneous monitoring of various process variables. However, interpretation of the experimental results cannot be fully realized without reliable and convenient multiresponse methods.

The following are some of my thoughts on the needs in this area:

1) Often, kineticists exhibit a philosophical resistance to a multiparameter approach to experimation for automatic coding of the method of Bates and Watts.

___

[1] Michael Frenklach's contribution to the subject stems from work performed in the Department of Chemical Engineering, Louisiana State University.

438

mentation. A classical way is to "isolate" a given reac- iton of interest; under such conditions the rate coefficient parameters can be determined by a simple well-established straight-line treatment. Determination of more than one rate coefficient in a single set of experiments is considered "not clean experimentation." In principal, however, the isolation is not possible: there are always other reactions occurring simultaneously with the one of interest. The researchers usually engage in an elaborate line of reasoning to assume, sometimes unjustifiably, single-reaction conditions. These kineticists must realize that multiparameter analysis using rigorous multiresponse techniques can provide more accurate and informative answers. Neglecting, for instance, a chemical reaction with the rate contribution of, let us say, 10%, can lead to a much larger than 10% distortion in the estimation of the main parameter. Statisticians, on the other hand, should demonstrate the techniques they develop on examples of current interest and difficulty.

2) Although first order kinetic models constitute an important class, higher order kinetics are of more general interest and there is a great need for development of statistical methods for these nonlinear systems.

3) Most estimation methods, including the one presented by Bates and Watts, concentrate on determining the solution which minimizes the objective function and only approximate confidence limits. What is of interest to many applications is the joint confidence region. It should be noted that in the problems of chemical kinetics these regions are usually not ellipsoidal, for which second order approximation methods are sufficient, but crescent shaped.

4) While estimating parameters, it is most important to check the model adequacy. This point was excellently demonstrated by Box and Draper (1965). These authors warn that "the investigator should not resort immediately to the joint analysis of responses. Rather he should... consider the consistency of the information from various responses." To my knowledge, however, a formal multivariate lack-of-fit test for a general nonlinear case has not been developed.

5) A question on the number of degrees of freedom

was brought up by Bates and Watts. Using fast digital sampling electronics, the number of observations per response can be very large (in our laboratory this number was approximately 1000). Does this number determine the degrees of freedom? If so, then one can easily increase this number by orders of magnitude by using faster electronics. This point should be clarified.

Finally, I would like to point out that in an attempt to resolve some of the issues brought up above, a method for multiresponse parameter estimation applicable to a dynamic model of general order was developed in our laboratory (Miller and Frenklach, 1983; Frenklach, 1984; Frenklach and Miller, 1985). The method is based on approximating the solution of the differential equations describing the kinetics of reactive system instead of the equations themselves. The approximation is developed following the methods of empirical model building (Box et al. 1978) and the concept of computer experiment of Box and Coutie (1956). Once the approximations to all responses are obtained, the parameter estimation, determination of joint confidence region, and lack-of-fit test are easily performed following the approach of Box and Draper (1965).

## References

[1] Box, G. E. P., and G. A. Coutie, Application of digital computers in the exploration of functional relationship, Proc. I.E.E. 103B (Suppl. No. 1) 100–107 (1956).

[2] Box, G. E. P., and N. R. Draper, The Bayesian estimation of common parameters from several responses, Biometrika 52, 355–365 (1965).

[3] Box, G. E. P.; W. C. Hunter and J. S. Hunter, Statistics for Experimenters. An Introduction to Design, Data Analysis, and Model Building, Wiley: New York (1978).

[4] Frenklach, M., Modeling, Combustion Chemistry (Edited by W. C. Gardiner, Jr.), Chap. 7, Springer-Verlag: New York (1984).

[5] Frenklach, M., and D. L. Miller, Statistically rigorous parameter estimation in dynamic modeling using approximate empirical models, AICHE J. 31, 498–500 (1985).

[6] Miller, D., and M. Frenklach, Sensitivity analysis and parameter estimation in dynamic modeling of chemical kinetics, Int. J. Chem. Kinet. 15, 677–696 (1983).

# The Use of Kalman Filtering and Correlation Techniques in Analytical Calibration Procedures

## H. C. Smit

University of Amsterdam, Nieuwe Achtergracht 166, 1018 WV Amsterdam, The Netherlands

Different chemometric methods to improve calibrations are described. A Kalman filter is applied for processing and predicting slowly varying parameters of a linear calibration graph. The results are used for the evaluation of unknown samples, and for deciding whether to calibrate again or to analyze the next unknown sample. Another approach of the calibration problem, particularly in chromatography, is the use of correlation techniques. The noise reduction property of correlation chromatography is used to extend the calibration graph to very low concentrations. Furthermore, an experimental technique to determine a calibration curve and the unknown sample simultaneously under exactly the same conditions is described.

Key words: calibration; chromatography; correlation techniques; Kalman filter.

## 1. Introduction

The computer has added a new dimension to analytical chemistry. Chemometrics, the application of mathematical and statistical techniques, is improving the quality of the analytical results concerning accuracy, precision, time, and costs, and has created new possibilities. An extended number of chemometric procedures are now available and are being increasingly applied in practice.

However, the "off-line" application of chemometric procedures, i.e., the processing of data or signals already obtained with common analytical methods like titration, chromatography, spectroscopy, etc., is dominating. The computer is generally used as an off-line calculation machine. The incorporation of the computer into existing analytical methods or the development of new methods based on the capabilities of the computer is not far developed; the intelligent analyzer is still in its infancy.

In this paper, some examples of chemometric on-line computer applications in an analytical procedure and analytical method are given. The calibration, which is very important in analytical chemistry, is emphasized in both examples—a generally applicable procedure using an optimum recursive parameter estimation technique (vector Kalman filter), and a method developed for one particular analytical technique (chromatography). More details concerning the basic theory are given in [1][1].

## Calibration and Optimum Estimation

An analytical system is usually very complex and includes chemical, optical, electrical, and mechanical parts. All these parts are subject to several influences, like contamination, changes in temperature, humidity, etc., and, of course, aging. These influences result in a decrease of the quality of the analytical data. Two main

---

[1] Figures in brackets indicate literature references.

components can be distinguished: a stochastic component resulting in stationary random fluctuations, and a semi-random component caused by irreversible processes like the mentioned contamination, aging etc. This semi-random component has a non-stationary nature and a regular calibration is required to maintain the quality of the results respectively to reduce the influence of the drift in the calibration parameters. Ignored drift seriously affects the accuracy in the analytical results and various off-line drift correction procedures have been proposed [2-6].

A Kalman filter enables on-line drift compensation, particularly suitable in the case of automated analytical procedures. An optimum recursive estimator like a Kalman filter requires a model of both the system or the signal process, including the system noise and the measurement (observation) noise. Using an appropriate model, the Kalman filter can predict (estimate) future values of the changing parameters and the samples can be evaluated using these predicted parameters. The estimation may also be used to determine when a recalibration is required. The criterion is a given preset precision of the results. The final goal is to analyze samples with a predetermined minimum accuracy.

## System Model

A state space model is used to describe the system. The linear discrete dynamic system

$$x(k) = F(k) \, x(k-1) + w(k-1) \tag{1}$$

$$z(k) = h'(k) x(k) + v(k) \tag{2}$$

where

| | |
|---|---|
| $k$ | :a time or a sequence number |
| $x(k)$ | :$n \times 1$ state vector |
| $F(k)$ | :$n \times n$ transition matrix |
| $h'(k)$ | :$1 \times n$ measurement vector |
| $z(k)$ | :measured signal |
| $w(k-1)$ | :$n \times 1$ system noise vector |
| $v(k)$ | :scalar measurement noise |

is representative for many analytical systems (fig. 1). The model is linear because neither in the transition matrix $F(k)$ nor in the measurement vector $h'(k)$ parameters $x(k)$ are present. A commonly used calibration graph is

$$y = a \cdot c + b \tag{3}$$

where

$a$: sensitivity
$b$: intercept

$c$: concentration
$y$: measurement

Reformulation gives

$$y = (c, 1) \begin{pmatrix} a \\ b \end{pmatrix} \tag{4}$$

Using

$$x = \begin{pmatrix} a \\ b \end{pmatrix} \quad F = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} \tag{5}$$

$h' = (c, 1)$ and $z = y$

and adding system noise and measurement noise results in a dynamic calibration curve

$$\begin{pmatrix} a\,(k) \\ b\,(k) \end{pmatrix} = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} \begin{pmatrix} a(k-1) \\ b(k-1) \end{pmatrix} + \begin{pmatrix} w_1(k-1) \\ w_2(k-1) \end{pmatrix} \tag{6}$$

$$z(k) = (c, 1) \cdot \begin{pmatrix} a(k) \\ b(k) \end{pmatrix} + v(k). \tag{7}$$

A common situation in practice is that the parameters of the calibration curve are slowly varying in time. Stochastic variations can be represented by the system noise, introduced in the model. However, in practice often a deterministic variation of the calibration parameter can be observed; particularly an extension of the model with a linear drift increases the usability.

A single parameter (state) $x$, affected by a linear drift in the sequence $k$ can be written as

$$x(k) = dk + e \tag{8}$$

where

$d$: (constant) drift parameter
$e$: $x(0)$

Evaluation of eq (8) leads to

$$x(k) = d(k-1) + d + e$$
$$= x(k-1) + d. \tag{9}$$

The introduction of system noise and measurement noise gives

$$\begin{pmatrix} x(k) \\ d(k) \end{pmatrix} = \begin{pmatrix} 1 & 1 \\ 0 & 1 \end{pmatrix} \begin{pmatrix} x(k-1) \\ d(k-1) \end{pmatrix} + \begin{pmatrix} w_1(k-1) \\ w_2(k-1) \end{pmatrix} \tag{10}$$

$$z(k) = (1, 0) \begin{pmatrix} x(k) \\ d(k) \end{pmatrix} + v(k).$$
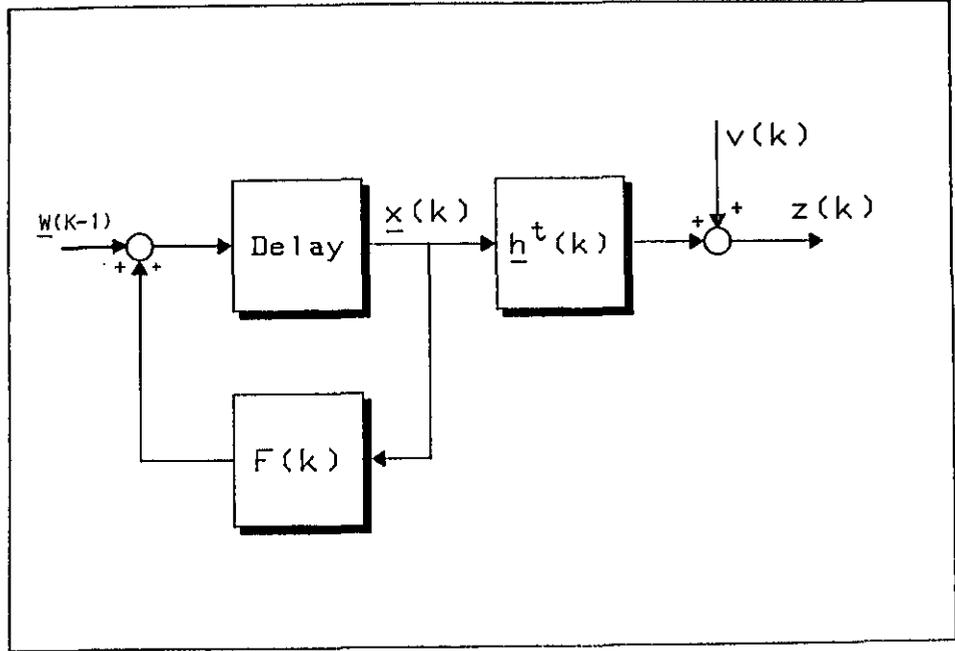
442

Figure 1-Linear dynamic system.

If both the sensitivity $a$ and the intercept $b$ are influenced by random drift, two extra parameters have to be introduced in the original model. The final model is

$$\begin{matrix} a(k) \\ b(k) \\ \alpha(k) \\ \beta(k) \end{matrix} = \begin{matrix} 1\ 0\ 1\ 0 \\ 0\ 1\ 0\ 1 \\ 0\ 0\ 1\ 0 \\ 0\ 0\ 0\ 1 \end{matrix} \begin{matrix} a(k-1) \\ b(k-1) \\ \alpha(k-1) \\ \beta(k-1) \end{matrix} + \begin{matrix} w_1(k-1) \\ w_2(k-1) \\ w_3(k-1) \\ w_4(k-1) \end{matrix}$$ (11)

$$z(k)=(c,1,0,0) \begin{matrix} a(k) \\ b(k) \\ \alpha(k) \\ \beta(k) \end{matrix} +v(k).$$

The observability matrix is

$$M= \begin{matrix} c_1 & c_2 & c_3 & c_4 \\ 1 & 1 & 1 & 1 \\ 0 & c_2 & 2c_3 & 3c_4 \\ 0 & 1 & 2 & 3 \end{matrix}.$$ (12)

The system is observable if there are at least two different concentrations.

## Sample Evaluation and Recalibration

Assuming the transition matrix $F(k)$ in eq (11) is known exactly and the measurement vector $h'(k)$ is known from the calibration, and assuming the statistical properties of the system noise and measurement noise are known ("white" noise with zero mean and normal probability density function (pdf), the usual Kalman filter algorithms, given in eqs (13–17), can be used to esti-

mate the slope, the intercept, and the drift parameters of the calibration curve [7]

$$x(k/k-1)=F(k)\hat{x}(k-1/k-1)$$ (13)

$$P(k/k-1)=F(k)P(k-1/k-1)F'(k)+Q(k-1)$$ (14)

$$x(k/k)=\hat{x}(k/k-1)+k(k)\{z(k)-h'(k)\hat{x}(k/k-1)\}$$ (15)

$$P(k/k)=P(k/k-1)-k(k)h'(k)P(k/k-1)$$ (16)

$$k(k)=P(k/k-1)h(k)\{h'(k)P(k/k-1)h(k)+ R(k)\}^{-1}$$ (17)

where:

$Q(k)$: system noise
$R(k)$: measurement noise
$k(k)$ : Kalman gain factor (correction factor)

Equation (18) gives the innovation, i.e., the difference between the experimental and the estimated measurement

$$v(k)=z(k)-h'(k)\hat{x}(k/k-1)$$ (18)

with

$$E\{v(k)\}=0$$

$$E\{v(k)v'(l)\}=\{h'(k)P(k/k-1)h(k)+R(k)\}\delta(k,l).$$ (19)

The variance of the predicted measurement is also

443

given by eq (19). The prediction of the Kalman filter can be used to evaluate a constituent in an unknown sample. Rewriting the calibration relation eq (3) gives

$$\hat{c}_{un} = (z(k) - \hat{b})/\hat{a}$$

$$\text{var}(\hat{c}_{un}) = (1/\hat{a}^2)\{h'_{un}P(k/k-1)h_{un} + R_{un} \qquad (20)$$

with: $h_{un}' = (\hat{c}_{un}, 1, 0, 0)$ and $\hat{a}, \hat{b}$ from $\hat{x}(k/k-1)$.

The relative imprecision of an unknown concentration $c_{un}$ is given by

$$N_{un} = 2/(c_{un}\hat{a})\{h'_{un} P(k/k-1)h_{un} + R_{un}\}^{\frac{1}{2}} \qquad (21)$$

with: $h'_{un} = (c_{un}, 1, 0, 0)$.

The available calibration standards are used to compute $N_{un}$ in eq (21). The computed maximum imprecision $N_{max}$ is compared with a predefined imprecision $N_{crit}$. If $N_{max} \geqslant N_{crit}$, a recalibration is performed; if $N_{max} < N_{crit}$, the samples are processed.

## Application in Practice

The described state estimation is applied in automated flow injection analysis. The quality of the results is improved by smoothing all the stored estimates of the Kalman filter. An extensive description of the smoothing procedure is given in [8].

Figure 2 shows an automated flow injection system used for the determination of chloride in aqueous samples. Thiocyanate originating from $Hg(SCN)_2$ is substituted by $Cl^-$ in the presence of $Fe^{3+}$. Red coloured $Fe(SCN)_3$ is formed and measured spectrophotometrically at 470 nm. Each sample requires 40 seconds; 90 samples/hour can be processed.

Figure 3 shows the results of the repeated injections of 2, 4, 6, 8, and 10 ppm samples. The chi-square values given in table 1 are obtained with a noise variance $R = 15.10^{-6}$ and system noise covariances $Q_{33} = Q_{44} = 10^{-10}$.

As can be seen from table 1, a first order calibration graph (4 parameter state, $x1 - x4$) is obviously not satisfying and the model has to be extended to a second order drifting calibration graph (6 parameter state, $x1 - x6$). In this case smoothing does not yield significant improvement in the estimation.

Figure 4 shows the result of the measurement of standards and "unknown" samples (6 ppm) at fixed positions in the sequence. The on-line Kalman estimation is depicted in figure 5 and the improvement by the smoother is shown in figure 6. The histograms figure 7 and figure 8 show the evaluated results.

The on-line processing of the uncorrected peak heights permits one to decide to recalibrate or not. Figure 9a and figure 10 give an impression of the results, respectively for the estimation of the state by the Kalman filter and after the smoothing. If a given preset criterian $N_{crit}$ is exceeded by the maximum imprecision $N_{max}$, the system is recalibrated. After 9–12 calibrations the system starts to evaluate the unknown samples and recalibrates regularly.



Figure 2-Flow injection system. S=sample holder, C=column, P=pump, D=spectrophotometer, I=injection valve, W= waste.

Figure 3-Repeated injections.



Figure 4-Standards and unknown samples.

Table 1. Values of $x^2$ for the flow injection peaks of figure 3.

| | Peak height | Peak height baseline corrected | Peak integral/300 | Peak integral/300 corrected |
|---|---|---|---|---|
| Kalman filter (1st order) | 618.8 | 662.7 | 134.9 | 797.2 |
| Smoothed (1st Order) | 578.7 | 620.1 | 122.0 | 744.2 |
| Kalman filter (2nd order) | 69.8 | 67.5 | 61.0 | 395.4 |
| Smoothed (2nd order) | 62.0 | 59.3 | 53.8 | 354.4 |



Figure 5-On-line Kalman estimation.

445

Figure 6–Improvement by smoothing.



Figure 7–Histogram of the Kalman estimation.



Figure 8–Histogram of the smoothed results.



Figure 9–On-line calibration system. State and concentrations used by the Kalman filter.

446

Figure 10-On-line calibration system. State and concentrations used by the smoother.

## Correlation Techniques

The introduction of correlation techniques permits a completely different approach to the calibration problem, particularly in chromatography. In correlation chromatography (CC) the usual impulse-shaped injection is replaced by multiple random injections. An example of an application is given in [9]. The resulting random response of the chromatographic system is cross-correlated with the used input function. The cross-correlation function of two signals, in this case the input signal $x(t)$ and the output signal $y(t)$ of a linear process, is by definition:

$$R_{xy}(t_1,t_2)=E[x(t_1)y(t_2)] \qquad (22)$$

$E[\ \ ]$ denotes the expected value of the expression between the brackets. The output signal $y(t)$ of a linear system can be calculated as a convolution of the input signal $x(t)$ and the impulse response $h(t)$ of the system:

$$y(t)=x(t)*h(t)=\int_0^\infty h(\tau)x(t-\tau)d\tau . \qquad (23)$$

Combining eq (22) and eq (23) gives:

$$R_{xy}(t_1,t_2)=E\left[x(t_1)\int_0^\infty h(\tau)x(t_2-\tau)d\tau\right]. \qquad (24)$$

In this case integration and averaging can be interchanged, hence:

$$R_{xy}(t_1,t_2)=\int_0^\infty h(\tau)E[x(t_1)x(t_2-\tau)]d\tau \qquad (25)$$

However, the autocorrelation function of $x(t)$ is defined as:

$$R_{xx}(t_1,t_2)=E[x(t_1),x(t_2)] \qquad (26)$$

and eqs (25) and (26) can be combined to:

$$R_{xy}(t_1,t_2)=\int_0^\infty h(\tau)R_{xx}(t_1,t_2-\tau)d\tau \ (\tau>0) \qquad (27)$$

If $x(t)$ is stationary, then:

$$R_{xx}(t_1,t_2-\tau)=R(t_2-\tau-t_1)=R(t-\tau) \qquad (28)$$
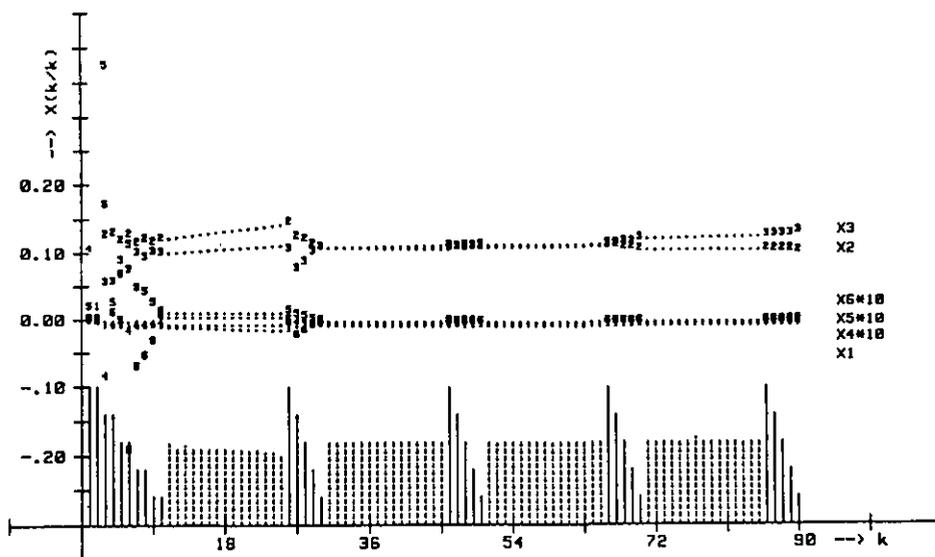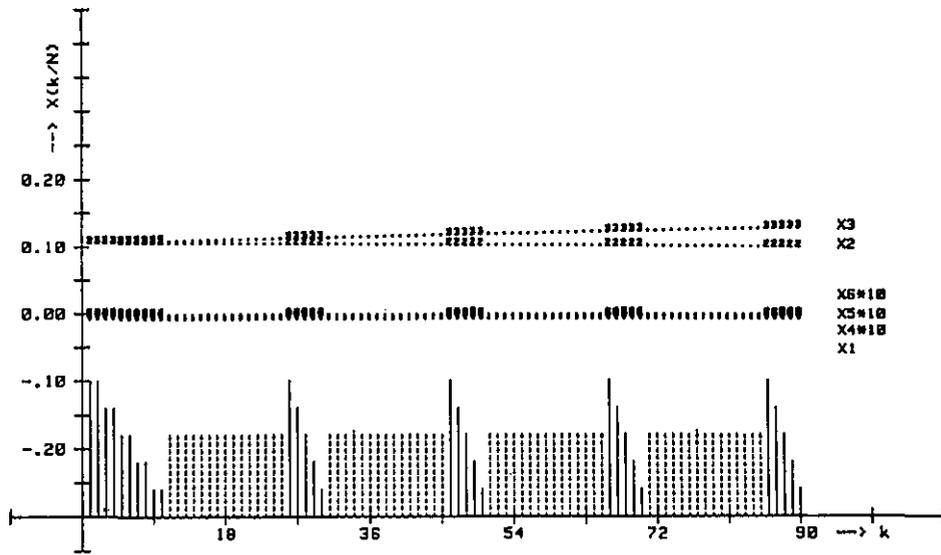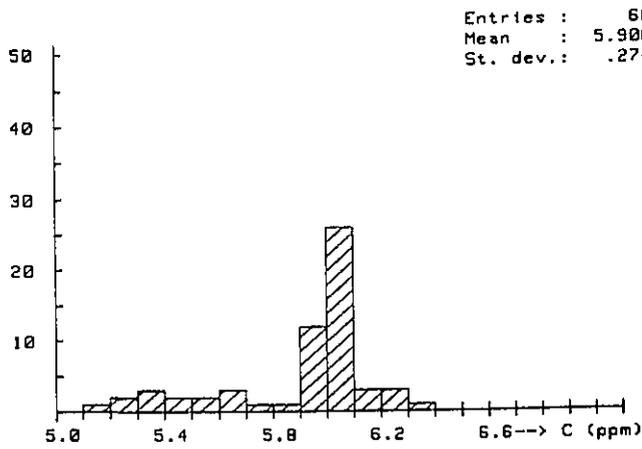
and eq (27) becomes:

$$R_{xy}(t_1,t_2)=R_{xy}(t)=\int_0^\infty h(\tau)R_{xx}(t-\tau)d\tau \qquad (29)$$

Comparing eq (23) and eq (29) shows that the output signal $y(t)$ of a linear system with an input signal equal to the autocorrelation function $R_{xx}(t)$ of a signal $x(t)$ is similar to the cross-correlation function $R_{xy}(t)$ of the input signal $x(t)$ and the output signal $y(t)$ resulting from $x(t)$.

A white noise, that is, white with respect to the bandwidth of the system, has in impulse-shaped autocorrelation function and can be used as input function $x(t)$ to determine the impulse response, in which case $R_{xy}(t)=h(t)$.

On further consideration a chromatographic procedure can be regarded as the determination of an impulse response; a chromatogram shows the response on the impulse-shaped injection of the sample. The prime objective of correlation chromatography is to determine the chromatogram by stochastically injecting the sam-

447

ple into the column and cross-correlating the input and the resulting output. If the chromatographic system is contaminated with noise, this noise is not correlated with the input and its contribution to the overall cross-correlation function converges to zero with increasing correlation time. A considerable improvement of the signal to noise ratio can be achieved in a relatively short time.

## Application in Practice

The most suitable random input function, controlling the input flow of the sample, is the pseudo random binary sequence (PRBS). This function is to be preferred to other random inputs with approximately impulse-shaped autocorrelation functions for the following reasons:

1) It is a binary noise, the only two levels being +1 and −1 or +1 and 0. The levels can be used to control simple on/off valves and correspond with the injection of sample and eluent, respectively;

2) It can be easily generated and reproduced; and,

3) Its special properties offer the possibility of reducing the correlation noise, which is caused by a limited correlation time.

The PRBS is a logical function combining the properties of a true binary random signal with those of a reproducible deterministic signal. After a certain time (a sequence) the pattern is repeated. It is important that the *estimated* autocorrelation function of a PRBS, if computed over an integral number of sequences, is at all times exactly equal to the autocorrelation function. Figure 11 shows a correlation HPLC set-up.



Figure 12-Calibration graph of phenol with fluorimetric detection.

The analytical performance of CC is demonstrated in figure 12. A calibration curve of phenol was measured over five decades of concentration: $0.01-100$ $\mu g$ $l^{-1}$. Conventional HPLC equipment with fluorimetric detection and a newly developed injection device for correlation HPLC was used. The two higher concentrations $(10-100$ $\mu g$ $l^{-1})$ were determined by conventional (reverse phase) HPLC and the two lower concentrations $(0.01-0.1$ $\mu g$ $l^{-1})$ by correlation HPLC with 16 and 3 sequences of correlation time, respectively.

Measurements at the $1 \mu g$ $^{-1}$ level were performed both by conventional and correlation HPLC (1 sequence).

The bars indicated on the calibration graph represent the peak area $\pm 3\sigma_i$ (arbitrary units), when $\sigma_i$ is the standard deviation of the integrated noise [10]. The inner bars at the 1 $\mu g$ level represent the correlation re-
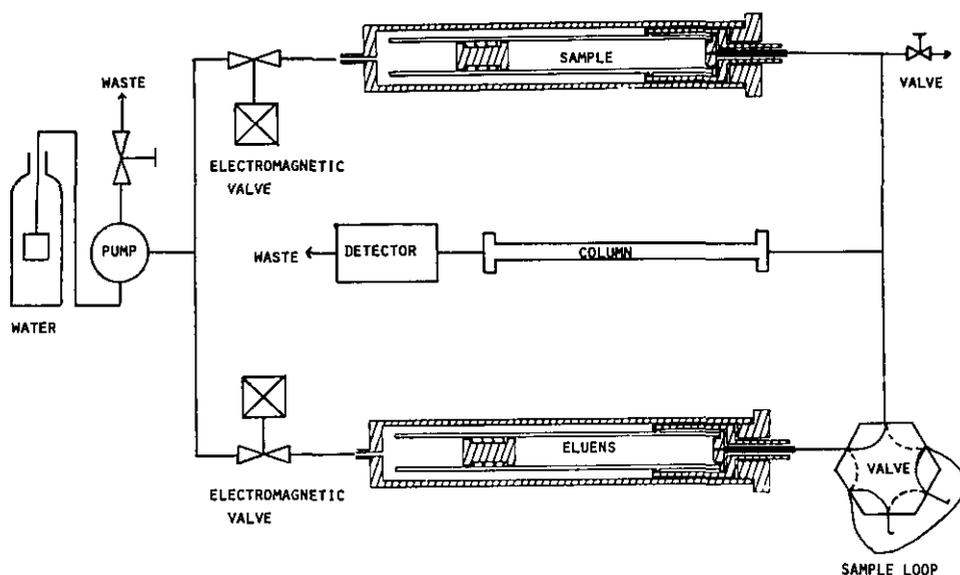


Figure 11-Set-up of a correlation HPLC system. The constant water flow is depending on a PRBS pattern directed either to the sample or to the eluent reservoir.
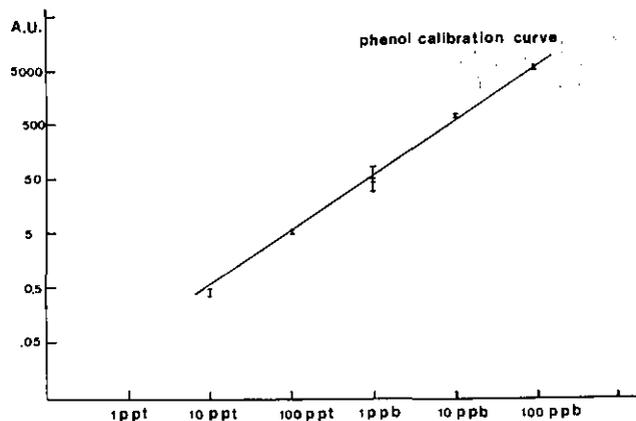
448

sults and the outer bars the single injection results. The detection limit for the single injection experiments, defined as $3\sigma_1$, was about 0.5 $\mu g\ l^{-1}$. The detection limit with the 10 ng $l^{-1}$ concentration was estimated to be 3 ng $l^{-1}$ (3 ppt).

## Simultaneous Correlation Chromatography

On further consideration, the question arises whether it is possible to inject different samples simultaneously, each sample injection controlled by its own unique PRBS. If these pseudo random sequences are mutually uncorrelated, then the correlogram (chromatogram) of each sample can be determined without any influence of the other samples, even if the components in the samples are the same. The problem is to find uncorrelated sequences.

A possible solution is the following. A binary pseudo random noise sequence is generated by a digital shift register with a suitable modulo-2 feedback. Appropriate digital multiplexing yields multiple uncorrelated pseudo random outputs from a single $n$-bit shift register. Each of the $K$ output sequences is identical to the single shift register sequence; they are staggered by $2^n/K$ bits, so they will remain uncorrelated for $1/K(2^n-1)$ output bits. These output sequences can be used for the mentioned simultaneous input patterns of a chromatographic column.

Figure 13 shows the result of a simulation experiment. A separation of two components is simulated with four different samples which were "injected" simultaneously. The "separation" is excellent.

The injection of $n$ samples requires a correlation time of $n$ times the total elution time. Therefore, no time gain can be expected. However, the method can be used for a number of interesting applications. For example, simultaneous CC permits measurement and calibration at



Figure 13–Simulated simultaneous chromatogram.

the same time under the same chromatographic conditions.

Figure 14 shows an experimental set-up of a simultaneous chromatograph (HPLC) with four reservoirs (three samples and eluent) together with four valves, each controlled by a sequence uncorrelated with the others. The flow stability is maintained by dividing a clock period in four parts. In each part either eluent or sample (respectively sample 1, 2, and 3) can be injected, depending on the status (0 or 1) of the sequence concerned.

The results of the analysis of three samples with naphtalene, anthracene, and 1,2 benzanthracene is shown in figure 15. The anthracene concentration in each sample is the same; the concentration of the other components in the different samples differ by a factor of 2.

The experimental injection system is not yet perfect, and because of this, serves as a source of so-called correlation "noise." This noise is not really random, but is



Figure 14–Experimental set-up of a simultaneous HPLC system.

449

Figure 15–Simultaneous chromatogram of three samples, each containing naphtalene, anthracene and 1,2 benzanthracene with different concentrations.



Figure 16–Calibration graph for benzanthracene, simultaneously determined.



Figure 17–Calibration graph for benzanthracene, successively determined.

composed of deterministic signals (ghost peaks). The peaks of benzanthracene are used to construct a calibration curve (fig. 16). Comparison of this curve with the calibration curve determined in the usual way (fig. 17) shows the performance of the method. The advantages are twofold: the random fluctuations are reduced by the multiple injection and averaging property, and both the unknown sample and the calibration sample are measured simultaneously under exactly the same conditions.

The final conclusion is that on-line chemometric techniques, such as Kalman filtering and correlation procedures, create promising new possibilities in analytical chemistry. The given improvements of the calibration procedure are a typical example of the power of these techniques.

---

Major contributions to this paper were made by P. C. Thijssen, J. M. Laeven and C. Mars.

## References

[1] Thijssen, P. C.; S. M. Wolfrum, G. Kateman, and H. C. Smit, Anal. Chim. Acta 256, 87 (1984).

[2] Nisbet, J. A., and E. Simpson, Clin. Chim. Acta 39, 367 (1972).

[3] Nisbet, J. A., and J. A. Owen, Clin. Chim. Acta 92, 367 (1979).

[4] Bennett, A.; D. Gartelman, J. J. Mason, and J. A. Owen, Clin. Chim. Acta 29, 161 (1970).

[5] Svehla, G., and E. L. Dickson, Anal. Chim. Acta 136, 369 (1982).

[6] Mclelland, A. S., and A. Fleck, Ann. Clin. Biochem. 15, 281 (1978).

[7] Gelb, A., (Ed.), Applied Optimal Estimation, MIT Press, Cambridge, MA, USA (1974).

[8] Thijssen, P. C.; G. Kateman and H. C. Smit, Anal. Chim. Acta (submitted for publication, 1985).

[9] Smit, H. C.; T. T. Lub and W. J. Vloon, Anal. Chim. Acta 122, 267 (1980).

[10] Smit, H. C., and H. L. Walg, Chromatographia 8, 311 (1975).

450

# DISCUSSION

of the H.C. Smit paper, The Use of Kalman Filtering and Correlation Techniques in Analytical Calibration Procedures.

## Diane Lambert

Department of Statistics, Carnegie-Mellon University

I would like to thank Professor Smit for his thoughtful paper on one of the most important steps in the measurement process: the calibration which determines how instrument response is translated into concentration units. As Dr. Smit points out, calibration in analytical chemistry is often difficult because the response of analytical instruments changes with time. The response may vary with ambient temperature, line voltage, and contamination that accrues with use, for example. To compensate for such fluctuations, certain U.S. EPA/CLP protocols for GC/MS instrumentation require that the calibration factor be determined every eight hours. At the beginning of a shift, a standard sample with known concentrations of the target chemicals is analyzed and its response factors for the target chemicals are used to quantitate all other samples analyzed in the same shift. In some cases, a standard "check sample" is analyzed at the end of the shift and its response factors are required to be within some percentage of those observed earlier. This calibration method has two major shortcomings. First, the calibration factor is determined by only one sample, and if there are any anomalies in its response factors, they affect all the samples analyzed in the shift. Even small variability in the response factors of standard samples may introduce unacceptable variability in measured concentrations between shifts. Second, the calibration factor is changed every eight hours regardless of how slowly or rapidly instrument response is changing.

In contrast, Dr. Smit proposes that the calibration factor be updated smoothly, based on the behavior of past samples, and that a new standard sample be analyzed only when the estimated imprecision of a measured concentration becomes intolerably large. There are also other perhaps less evident advantages to Dr. Smit's approach. First, the assumptions about the measurement process that justify the updating scheme are all explicit. Drift, measurement noise and system noise are modelled parametrically, so that the adequacy of models can be checked and the updating scheme can be modified if the models are found lacking. For example, in Professor Smit's application concening the determination of chloride in aqueous samples, a quadratic rather than linear model of drift is fit. Second, estimates of model parameters such as background and average drift are convenient for monitoring instrument performance. Third, the procedure automatically provides information about the uncertainty of measured concentrations. It is as important to report how trustworthy reported concentations are as it is to report the measurements themselves.

Prof. Smit has also considered simultaneous injection and measurement of standard and "unknown" samples. There are, however, some questions about his procedure for conventional sequential analysis of samples that I believe have not yet been resolved. For example, how effective is the procedure if samples are analyzed at a rate of one per hour rather than one per minute? What happens when several or hundreds of chemicals, perhaps all at trace levels, are measured for the same sample? In some examples, accuracy was increased by smoothing Kalman filter estimates and in others it was not. What guidance can be given to a laboratory technician? In short, what are the limits of applicability of this calibration method and when should it be authorized?

Prof. Smit has taken an important step towards improving the chemical measurement process. I look forward to his future work on his procedure.

# Intelligent Instrumentation

Alice M. Harper

University of Texas at El Paso, El Paso, TX 79968

and

Shirley A. Liebman

Aberdeen Proving Grounds, MD 21105-50666

Feasibility studies on the application of multivariate statistical and mathematical algorithms to chemical problems have proliferated over the past 15 years. In contrast to this, most commercially available computerized analytical instruments have used in the data systems only those algorithms which acquire, display, or massage raw data. These techniques would fall into the "preprocessing stage" of sophisticated data analysis studies. An exception to this is, of course, are the efforts of instrumental manufacturers in the area of spectral library search. Recent firsthand experiences with several groups designing instruments and analytical procedures for which rudimentary statistical techniques were inadequate have focused efforts on the question of multivariate data systems for instrumentation. That a sophisticated and versatile mathematical data system must also be intelligent (not just a number cruncher) is an overriding consideration in our current development. For example, consider a system set up to perform pattern recognition. Either all users need to understand the interaction of data structures with algorithm type and assumptions or the data system must possess such an understanding. It would seem, in such cases, that the algorithm driver should include an expert systems specifically geared to mimic a chemometrician as well as one to aid interpretation in terms of the chemistry of a result. Three areas of modern analysis will be discussed: 1) developments in the area of preprocessing and pattern recognition systems for pyrolysis gas chromatography and pyrolysis mass spectrometry; 2) methods projected for the cross interpretation of several analysis techniques such as several spectroscopies on single samples; and 3) the advantages of having well defined chemical problems for expert systems/pattern recognition automation.

Key words: data systems, intelligent; instrumentation; multivariate algorithms, statistical and mathematical; pattern recognition; preprocessing; pyrolysis gas chromatography; pyrolysis mass spectrometry.

Modern computer hardware and software technologies have revolutionized the direction of analytical chemistry over the past 15 years. Standard multivariate statistical techniques applied to optimization and control of instrumentation as well as routine decision making are at the forefront of new instrumental methods such as biomedical 3-dimensional scanners and pyrolysis MS and GC/MS as well as more established measurement techniques. Despite these advances, little attention has been paid to the exploitation of intelligent computerized instrumentation in the design phase of chemical research.

Instrumental intelligence is the ability of a scientific instrument to perform a single or several intelligence func-

About the Authors: Alice M. Harper is with the Chemistry Department at the University of Texas at El Paso. Shirley A. Liebman is with the Ballistics Research Laboratories at Aberdeen Proving Grounds.

tions in such a way that operations normally performed by the scientist are completely under automated computer control and decision making. Under this definition, intelligent instruments are quite common. Indeed in recent years, manufacturers of small scientific equipment have used the term "intelligent" in conjunction with single purpose items such as recorders to describe the addition of software and/or programmability to the device. Concurrent with this, larger scientific instruments have been marketed with data systems hosting a wide variety of intelligence functions including control and optimization of instrumental variables, optional modes of experimental design, signal averaging, filtering and integration as well as post analysis data massaging and library search interpretation. Although instruments with the software to perform sophisticated intelligence operations exist, they are not so readily marketed as intelligent instruments. For example, modern pulsed Fourier transforms nuclear magnetic resonance spectrometers (NMR) have microcomputers built into the system, operate over a wide range

of NMR experimental designs, and control instrumental parameters; however, decision making is, for the most part, an operator based function. For this reason, such instruments have limited "intelligence" in comparison to the level of intelligence required to carry an experiment to completion without extensive interaction with the operator. In fact, it could be that more intelligent research instruments might also be less versatile. The limitation is the current state of technology of intelligence programming. The instrument, if used for routine analysis where problem statements can be well defined, could operate with no loss of utility as a totally automated and intelligent instrument.

Figure 1 gives insight into the problems arising in creating intelligence programs for instrumentation. Figure 1a) is an analytical chemist's perception of a totally automated experimental design [1][1]. As can be seen, the experiment remains unspecified without an initial problem statement. The issue of intelligent data systems for single instruments is one of either defining the set of all possible problem statements in an evolutionary design or restricting the analysis to a single well defined problem. Figure 1b) is an example of a first stage multivariate data analysis system proposed for research applications in pyrolysis mass spectrometry [2]. The design is one that leaves the problem statement, data interpretation and decision making entirely in the hands of the scientist. What is really described in this figure is a statistical package residing on a microcomputer which receives data from the mass spectrometer. However, as Isenhour has

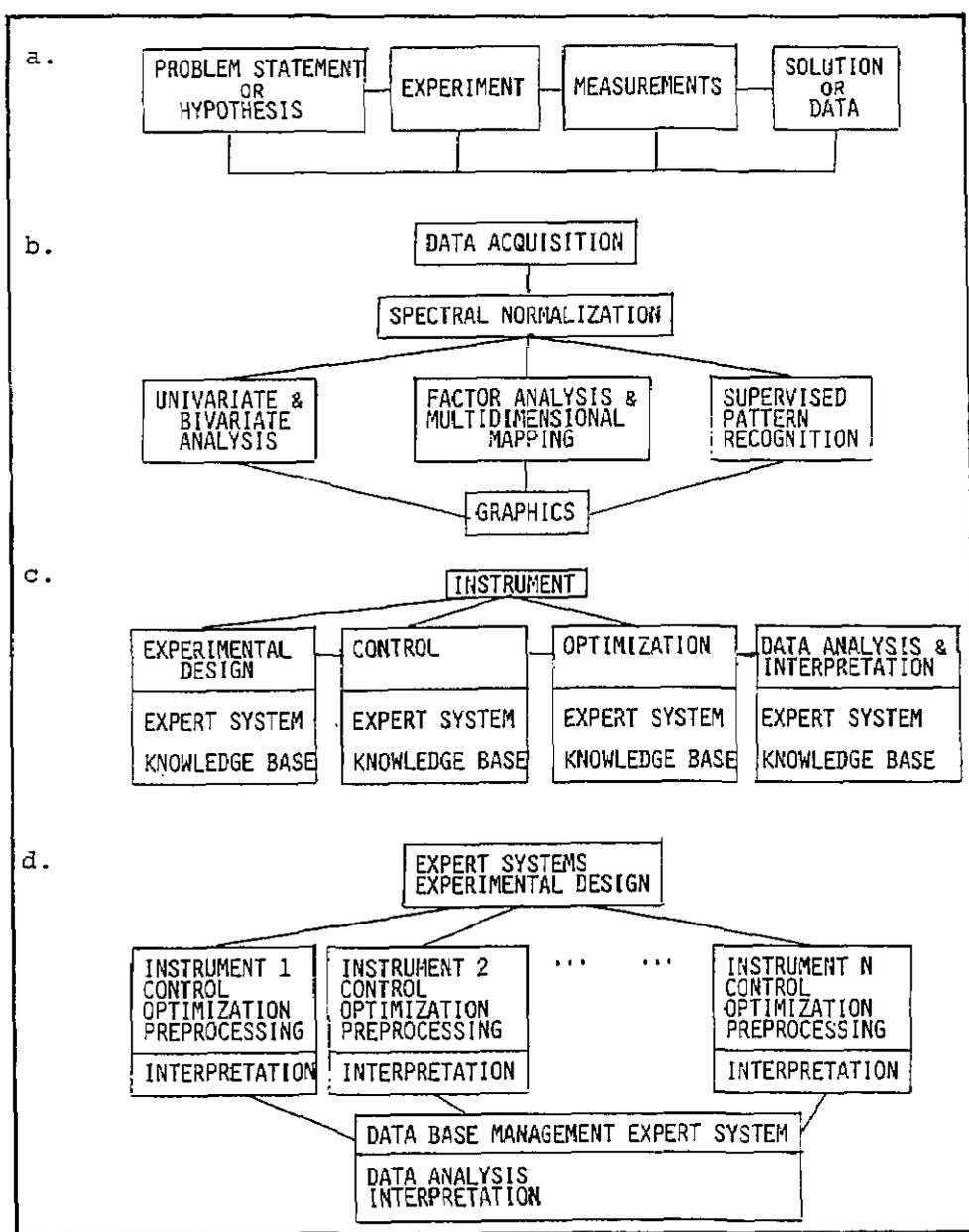---

[1]Figures in brackets indicate literature references.



Figure 1—Four designs for instrument data systems.
a. Totally automated experimental design and decision making.
b. Multivariate analysis for spectral instruments.
c. Expert systems driven single purpose instrument.
d. Laboratory automation using expert systems drivers.

demonstrated, multivariate factoring of spectral libraries offers advantages in the interpretation of complex single spectra [3]. It would seem that the incorporation of a multivariate statistics package in the data system is a key element in intelligence programming for many instruments. Figure 1c) incorporates the expert system approach to intelligence for instrumental systems. Under this design, the instrument can accommodate a series of problem statements and decision networks at various analysis stages and uses an intelligent driver for the multivariate analysis and interpretive stages of the analysis. Figure 1d) places such a system into the research laboratory controlling a variety of instruments and interpreting results based on one or more analyses.

Although the diagrams of figures 1b-1d are not comprehensive designs for total automation, they do provide a hierarchy for linking problem statements and decision making into multivariate research problems. After a brief discussion of components of intelligence designs, an example will be presented of the feasibility of developing expert system data reduction for pyrolysis analysis problems.

**Problem Statements:** Ideally a problem statement is analogous to the standard hypothesis in statistical analysis. Under the hypothesis a knowledge base can be collected and the hypothesis tested. For example, a patient does or does not carry a genetic trait [4]. Unfortunately, chemical research problems are often ill-specified and the problem statement may become a hierarchy of data investigations leading to one or more problem statements. For example, 1) are there differences in the chemical composition of a series of samples? if differences do occur; 2) what are the nature of the differences?; 3) do the chemical differences correlate with observed changes in physical properties?; and 4) can physical properties be predicted from the chemical differences? [5,6].

**Knowledge Base:** In order for a instrument to operate as an intelligent system, it must have the knowledge base necessary to arrive at a solution for each of its problem statements. This knowledge base may contain data, rules ("Mass peak 94 is phenol"), programs, and heuristic knowledge.

Consider the knowledge base required for setting up and operating routine analyses of polymer composition by pyrolysis gas chromatography. Possible problem statement areas include optimization of pyrolysis parameters, chromatographic conditions and interface characteristics, control of instrument and data acquisition parameters, data reduction, and data interpretation. The knowledge base must include all information necessary to each of the proposed problem statements. For optimization of parameters, rules governing the detection of an optimum, algorithms (e.g., "simplex" [7]) for efficiently moving toward an optimum, rules for hierachical movements within the algorithm, rules for the detection of poor optimization surface structure, and representative previous optima data might be employed in the decision making. Such optima

will be determined, in part, by the polymer degradation characteristics and therefore will not, in this case, be independent of the samples used in the analysis. Instrumental control might employ a knowledge base of rules for the automation of events such as the initiation of sample pyrolysis and data collection. Data reduction for this method will require the rules and data necessary for baseline correction, chromatographic normalization, and peak matching. The knowledge base requires a "memory" of previously collected data to aid peak matching protocols, rules for baseline determinations, transformations for baseline correction, normalization rules and algorithms, and rules for acceptance or rejection of the chromatogram under consideration. Data interpretation on the other hand might require a library of previous chromatograms as well as rules and/or algorithms for interpretation of the current event based on a knowledge of past data with verified interpretations.

Development of the knowledge base is the expensive and time consuming operation in the development instrumental intelligence even when the application is highly specific. It must be remembered that each operation that a human might perform automatically from experience must be programmed into the data system. For this reason, among the attributes of the system, there needs to be evolutionary operation. In other words, in addition to long term knowledge, facts about the current data and new conjectures under consideration must also be easily accommodated.

**Expert Systems Driven Multivariate Data Systems:** The response of many modern chemical instruments (e.g., spectrometers and chromatographs) is inherently multivariate. For such instruments, data reduction and interpretation often consume a greater portion of analysis time than data collection, and requires scientific expertise. The time delay between data collection and decision making has become an acute problem for newer hyphenated techniques, such as gas chromatography-mass spectrometry and mass spectrometry-mass spectrometry, which are capable of collecting thousands of mass spectral peaks in a short period of time.

One possible solution to the problems imposed by large bodies of data is to incorporate into the instrument a data reduction system consisting of multivariate analysis methods. The problem encountered in the actual implementation of such a system is that few experts in instrumental analysis have the expertise for carrying multivariate statistical analyses. It has become increasing apparent that instruments employing versatile multivariate based data systems should be capable of operating in a transparent data analysis mode. In order to accomplish this, the expertise of a chemometrician will need to be programmed into an expert systems driver for the instrument data system. Given a problem statement, and a knowledge base of rules from previous experience, the computer could decide from a variety of possibilities how to reduce the data and represent the results in a meaningful form. A relatively simplistic example of this is the problem

statement "Is there a correlation between two independent variables?" Invisible to the user would be the operations for determining the integrity of the variable distributions, a computation of the correlation and a determination of the significance of the correlation coefficient of the relationship. The result might take a simple form such as "there appears to be a significant correlation between the first variable and the logarithm of the second variable. Would you like to see the computed values?"

# Demonstration of Expert System for Data Analysis in Curie-point Pyrolysis Mass Spectrometry

Pyrolysis mass spectrometry has come to the attention of both mass spectrometrists and chemometricians because of its utility in the analysis of polymeric materials and the complexity of the mass spectra produced by natural polymers and biopolymers. The technique involves degradation of the solid material by pyrolysis followed by mass spectrometry of the pyrolysis fragments. It has been demonstrated by the pioneering work of Meuzelaar (for example see [8]) that Curie-point pyrolyses of samples of biomaterials can produce profiles which, when properly normalized [2], are reproducible and diagnostic of the chemical similarities and dissimilarities among groups of samples and are quantitative under appropriate experimental designs [9]. Because the process of pyrolysis followed by mass spectrometry of the network polymer of natural heterogeneous biopolymers produces mass spectra with peaks which tend to be highly correlated, the interpretation problems created by the vast number of, and the overlap of, the masses are solved through multivariate analysis of the mass spectral profiles. A data system for a pyrolysis mass spectrometer would be of limited utility without multivariate statistical methods [2].

To test the hypothesis that an expert systems approach to data reduction is potentially helpful and feasible, an expert systems driver was implemented to mimic the data analysis portion of a Rocky Mountain Coal study done at Biomaterials Profiling Center in Utah. Detailed results of the original study can be found in [5,6] and of the numerical methods used in [2]. Briefly, 102 Rocky Mountain Coals were analyzed in quadruplicate by pyrolysis mass spectrometry. The pyrolysis profiles were added to a preexisting data set containing conventional measurements on the same coal samples (see table 1). After normalization of the profiles by the method of Eshwis et al. [10], average mass specta were analyzed using multivariate analysis techniques.

Figure 2 is a minimal design for an expert systems driven data acquisition and analysis system for pyrolysis mass spectrometry. Figure 3 shows the design of the data bases for the present application. As diagrammed in figure 2, each

Table 1. Conventional measurement contained in the "old data" data base for the Rocky Mountain coals.

| Conventional Measurements | |
| --- | --- |
| Vitrinite | % Potassium |
| Fusinite | % Phosphorus |
| Semifusinite | % Moisture |
| Macrinite | % Pyritic sulfur |
| Liptinite | % Mineral matter |
| Vitrinite Reflectance | % Volatiles |
| % Silicon | % Organic sulfur |
| % Aluminum | Calorific value |
| % Titanium | % Organic carbon |
| % Magnesium | % Organic hydrogen |
| % Calcium | % Organic sulfur |
| % Sodium | % Organic nitrogen |
| | % Organic oxygen |

operation of the instrument requires an expert systems driver and decision module for its automation. A rudimentary expert systems was built to mimic the decision making used for statistical analysis of coal pyrolysis patterns. This system demonstrates the problems and pitfalls associated with intelligent instrumental development.

Normalization is rarely an option in pyrolysis techniques since the size of the sample undergoing electron impact is determined by the quantity of pyrolsates actually making their way to the ion source of the mass spectrometer. For this reason, spectra are normalized to place each sample on a relative quantitative basis. Furthermore, when replicates of a single sample are analyzed in detail, it is found that some peaks replicate better than others. For example, if an organic solvent is used during the sample preparation, the mass peaks due to this solvent replicate poorly. The same is true of contaminants absorbed to the sample matrix. On the other hand, because the sample size in terms of total ion counts is a variable in these experiments, often one or more replicate spectra will exhibit outlying tendencies when compared to other replicates of the same sample.

NORMA [2], developed by Meuzelaar's group at Utah, is designed to select peaks with stable variance characteristics for inclusion in the normalization process. With this routine an expert interacts with the computer in a loop of peak deletions and replicate spectra deletions until a set of peaks is defined which stabilizes the normalization process.

An expert systems approach to the software interaction represents peaks by their variances over the samples and sample replicates and spectra by Euclidean distances between replicates over the mass units. A library data base for normalization is established which contains spectra patterns for commonly used solvents and commonly encountered contaminates as well as the background spectrum from the mass spectrometer. The expert sysems are initialized by computation of peak variances. The variances are ordered high to low and the shape of the plot of ordered variances is

Figure 2—Expert systems driven pyrolysis mass spectrometer. Each stage of data acquisition, analysis, and interpretation has decision protocols based on a knowledge base of an expert. Note that the interaction between mathematical methods and the expert system chemometrician.



Figure 3—Types of data bases required for an application oriented expert system for instrumental analysis and interpretation. Knowledge representation takes two forms: 1. knowledge which when taken as a whole or in parts can be represented as a similarity or correlation; and 2. knowledge for which antecedent clauses must be satisfied in order for the interpretative/decision making process to occur.

analyzed. The library is searched under expert systems guidance to account for the peaks exhibiting high relative variance. For example, background spectrum is placed under consideration only when the total ion counts of the sample spectrum is less than a factor α of the background. After examination of the library, a decision is made as to which peaks will be deleted based on their expected contribution to the peak variation and with restrictions on the total number of peaks that can be deleted at this stage of analysis. This first deletion is a permanent deletion of peaks. None of the casual base peak deletions is reexamined at later stages.

The distance matrix of replicates is generated using the peaks remaining in the analysis and samples are deleted based on a comparison of their distances from the expected distance generated as a mean distance over the samples with variance $\sigma_d$. (Note that such a formulation may ignore systematic error among the sample replicates, and won't perform at the expert level under such a condition).

The next stage on analysis involves computation of both peak variances and sample distances. The rank order of the peak variances is compared. If variance reduction seems to be exhibited by a small set of peaks upon deletion of the

457

samples, the peaks involved are temporarily deleted, the samples previously deleted are brought back into analysis and the distances reexamined. The expert systems decide which will be deleted at this stage, peaks or samples based on the recomputed distances. The iterative decision making—statistical computation continues until a key set of peaks remain activated for normalization over all spectra. A diagram of the proposed normalization experts system is shown in figure 4.

The system, as described, does not completely emulate an expert for all applications of pyrolysis mass spectrometry. It has already been noted that the data evaluation does not address errors which arise from time dependent or systematic error. In addition to this, the decision making process, while designed to emulate the human decision process based on statistical results, does not necessarily operate on a one-to-one correspondence with a human expert. The problem is that two experts working on the same set of data may arrive at approximately the same results via slightly different procedural routes. The same is true of the expert system when compared to a human expert. The major deviations from a

human procedure found when working with this system is the lack of "intuition" or "fuzzy logic" that is used by the expert. The expert systems converges in more steps than is necessary by human interaction with the statistical algorithms and for some data bases, has trouble defining convergence at the solution. For example, the optimal cut off parameters for variance and distance change between data sets. These and other problems are best solved by training the expert systems to recognize the structure of a good special solution in addition to the structure of good statistics.

## Correlation Based Hypothesis Testing

Perhaps the most often asked question of large data sets involves finding relationships between the variables or between the variables and an external parameter. Table 2 lists the form of the problem statements included in our data system for this demonstration. The term "relate" invokes one of several multivariate algorithms for the study of correlations in the data. The possible responses are the Pearson correlation coefficient, linear regression, factor analysis or canonical correlation analysis.

Consider the problem statement: What is the relationship of peak 34 ($H_2S$) in the mass spectrum of coal to the total organic sulfur from the conventional data matrix. "Relate *current data, mass 34* to *old data, organic sulfur* over *samples, all*" results in a computation of the correlation coefficients, a estimate of significance and the confidence interval about the correlation. "Relate *current data, mass 34* to *old data, organic, sulfur over samples, all*; Interprete using *old data, organic sulfur*" results in the additional computation of (organic sulfur)=a (mass 34)+b with residuals $E_i$. The residual pattern is tested for randomness. A failure results in a search through the library for a reference residual pattern with similarities to the computed residual pattern.

The next relate function asks for a study of the relationships among variables in the current data set. "Relate *current data, all* with *current data, all* over *samples, all*" results in the factor analysis of the data correlation matix. The loadings of the factor analysis are interpreted for the variable relationships seen along the orthogonal axes of the original rotation. This interpretation is an experts systems based analysis of the major peak series and will be discussed later.

Interpretation of the factor score is a difficult problem. Consider the two dimensional factor score projection of these data. Figure 5a is a projection without labels. Figure 5b is the same projection after a human expert has assigned sample labels corresponding to the geological source of the samples and an interpretation. Figure 5a is representative of the information about the factor scores stored by the computer. The addition of labels is readily accomplished but, without training, patterns formed by the sample labels
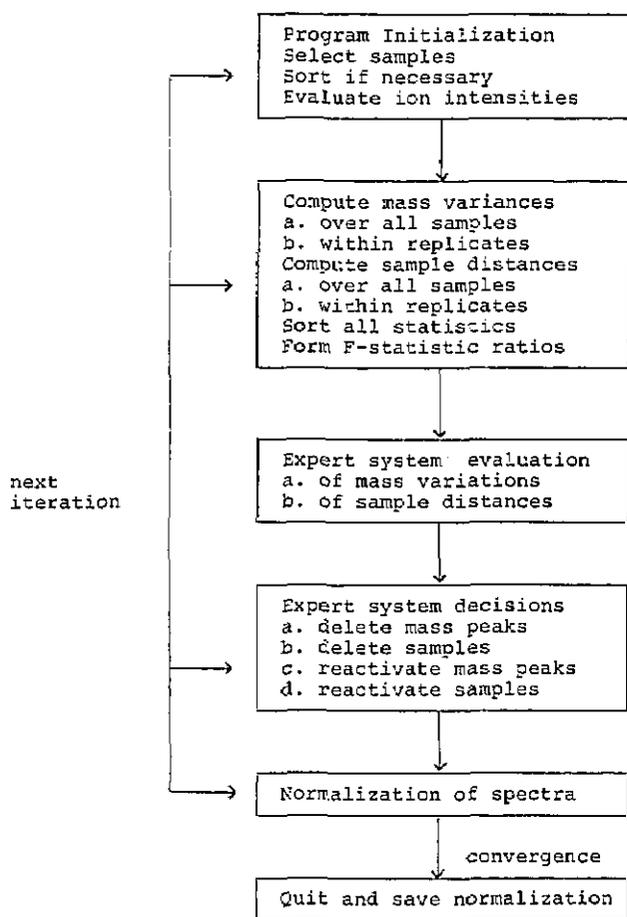


Figure 4—Interaction of expert systems and statistical computation for a rational normalization process for pyrolysis mass spectrometry.

458

**Table 2.** Variations of the expert system commands RELATE and INTERPRETE USING. The FORTRAN subroutine calling protocols are based on the data types used in each variable location of the command and on the command sequence. Note that the same command strings and rules can accomodate computation on questions about the data base transpose matrix when "...OVER samples ..." is replaced by "...OVER variables..."

RELATE data base 1, variable list 1 TO data base 2, variable list 2 OVER samples, sample list

**Examples:**

1. RELATE current data, mass 34 TO old data, organic sulfur OVER samples, match
   (Results: correlation coefficient and significance test)

2. RELATE current data, mass 34 TO old data, organic sulfur OVER samples, match; INTERPRETE USING old data, organic sulfur
   (Results: correlation computed at least squares fit, significance test, and residual pattern evaluation)

3. RELATE current data, all TO current data, all OVER samples, all
   (Results: Factor analysis of current data and loading interpretation)

4. RELATE current data, all TO old data, calorific value OVER samples, match, active; INTERPRETE USING old data, calorific value
   (Results: Factor analysis of current data followed by target rotation to calorific value and interpretation of loadings)

5. RELATE current data, all TO old data, all OVER samples, match
   (Results: Canonical correlation analysis of two data bases and interpretation of mass spectral loadings)

within the space have little meaning. A generalized solution to this problem does not seem likely. Each study would require an elaborate knowledge base specific to the samples in order to interpret trends seen in this picture.

For the coal data, the old data matrix of conventional measurements provides a more easily implemented route for hypothesis testing on the pyrolysis mass spectral factors. Consider once more the "relate" command. After matching the two data sets by the logic used in [5], factor analysis of the PY/MS set followed by "Relate *current data 2, all* with *old data, calorific value* using *samples, active*; interpret using *old data, calorific value*" results in the regression analysis of calorific value$=\Sigma$ $W$ (factor scores)$_{py-ms}+b$ producting both the variable and the sample relationships of the rotation of the Py/MS data to calorific value. The results are given in figure 6 and discussed in [6].

The last example of a correlation based statement is "Relate *current data, all* with *old data, all* over *samples, all*." This results in a cannonical correlation analysis factoring both data matrics in such a way that the overlap in information between the sets is maximized. Four factors are extracted. The first two of these are shown in figures 7 through 9. This analysis showed these data sets to be approximately 80% correlated over the coal samples. The major chemistry of the conventional measurement trends are given above the spectral representation form the Py/MS factors.



**Figure 5**—Comparison of computer representation of factor scores (a.) to the same plot after interpretation by an expert (b.) Symbols represent sample specific relationships deduced by the expert. Dotted lines are a separation of the samples into classes not available to the computer in the pyrolysis data base. The knowledge base required to mimic the expert would form a reference book of information about the samples. The generation of such a knowledge base would be a formitable task.

**459**

**Figure 6**—Chemical interpretation of the mass spectral peaks associated strongly with a targeted least squares rotation of pyrolysis factor scores to the external parameter calorific value. Peak interpretations were accomplished by the system described in figure 10. Results on chemical interpretation are identical to those of the original study, demonstrating that chemical information is more readily implimented as an expert system than sample specific information for this instrumental method.



**Figure 7**—First canonical variate loadings for the rotation of pyrolysis mass spectra of coal samples to the conventional data matrix described in table1. The correlation of the data bases to the derived variate (Z) is 0.99 and to each other is 0.95. $A_cX_c$ and $A_pX_p$ are the linear composites of the conventional and pyrolysis data bases respectively. Only the signs of the strongly correlated variables of the conventional data are given. Loadings for the pyrolysis data are given as positive and negative. Interpretations of mass peak were accomplished using the system described in figure 10.

460

$$R(Z, A_c X_c) = R(Z, A_p X_p) = 0.96$$
$$R(A_c X_c, A_p X_p) = 0.86$$

CONVENTIONAL MEASUREMENT SET

| SEMIFUSINITE | − |
| SODIUM | + |
| VOLATILES | + |
| REFLECTANCE | − |

PYROLYSIS MASS SPECTRAL DATA SET

**Figure 8**—Second canonical variate loadings for the rotation of pyrolysis mass spectra to the conventional data matrix described in table 1. See figure 7 for an explanation of the symbols used in this figure. Interpretation by the system described in figure 10 failed for the positive loadings labeled as isoprenoid signals?. These were assigned by the authors.

**Figure 9**—Canonical variate scores on the first two axes (described in figs. 7 and 8). Dotted lines separating classes were inserted by the authors and not interpreted by the expert system (see fig. 5 for an explanation of implimentation problem for sample interpretation.)



461

The assignment of chemical interpretations for mass spectra and for factor loadings is accomplished by an expert systems intepreter that can be used for any study in which the samples are coal. The data base for the interpretation includes commonly encountered chemical species in coal, the major peaks expected from their presence and, given two chemical species with similar patterns, the probability of their contribution as a major component to a coal pattern. Also included is a routine for generating molecular species from C,N,O, and S given the base peak molecular weight and the ion series. The operation of the expert system along with the results of each iteration are given in figure 10. Because Py-MS spectra contain many low intensity masses of questionable interpretation, and because factor loadings are rarely pure, only the most intense ion series patterns are interpreted. The program is initialized by setting an initial threshold limit (TL3). The ions ($m/z$) above the threshold are collected, sorted and given temporary chemical assignments. The threshold limit is lowered each time by a lesser factor until it encounters the "grassy" region of the spectrum. At this point, permanent ion series interpretations are assigned. The "?" appearing in the figure for Mass 104 means that no library interpretation of this peak was found and that the number of peaks in the ion series was below the limit set for generation of hypothetical molecular species. The entire system is diagramed in figure 11.

## Discussion

The correlation work described in this paper along with similar considerations of other algorithms seem to support the possibility that, for a given instrument, an expert systems driver for data analysis can be developed which is independent of the nature of the chemical problem. The way to accomplish this is to build expert systems drivers to interpret the problem statement and to interpret the data analysis results. Viewed in this manner, an instrument-dependent expert system can generate the experimental design and optmization from the problem statement and pass to the data analysis driver the statistical elements necessary for its decision on the proper data reduction protocol. Learning is initiated when the data analysis system receives an unrecognizable set of elements. Otherwise, the expert system selects from the knowledge base the algorithm sequence for data reduction.

We are currently extending these concepts to the construction of an expert system, EXMAT, for experimental design, optimization, data reduction and interpretation of measurements made on a sample by a variety of thermal analysis instrumental techniques. TIMM® (General Research Corp., Mclean, VA) a FORTRAN-based expert system generator, has been enabled development of a heuristically-linked set of expert systems for material analy-



Figure 10—Operation of the expert system used to interprete the pyrolysis mass spectra and spectral loadings in the demonstration. The rules used to sort and interprete the peaks are based on ion series produced by mass differences of 14 ($CH_2$). Base peaks help terminate series for resolution of multiple interpretations. ? is a peak not present in the data base.

sis. The attributes of the TIMM® system are listed in table 3. In order to accomplish the analytical goals of this project, we have combined the concepts of specific instrumental intelligence with the goals set forth in figure 1a to provide a data system capable of experimental designs utilizing one or several analysis techniques. Each instrument retains its own control, design, preprocessing and interpretative expert systems unit but the data analysis unit has, of necessity, been generalized for analysis of data from a single or from

462

Figure 11—Steps in the expert system interpretation on pyrolysis mass spectra of coals.

**Table 3.** Attributes of the TIMM[R] expert system. TIMM[R] was adapted to our application because addition of FORTRAN subroutine library calls is more readily accomplished than in other systems and because the decision logic can already use computationally based rules as well as chaining logic rules.

TIMM[R]: A FORTRAN - BASED EXPERT SYSTEMS APPLICATIONS GENERATOR, General Research Corporation

Forward/backward chaining using analog rather than propositional representation

Knowledge base divided into two sections: declarative knowledge and knowledge body

Pattern-matching using a nearest neighbors search algorithm to compare current situation with antecedent clauses

Unique similarity metric computed form order information in declarative knowledge giving distance metric over all classes

Decision structure and knowledge body readily developed and modified by expert in any domain

Heuristically-linked expert system using implicit and explicit method permitting processing of "microdecisions" that are part of "macrodecisions"

Each system independently built, trained, exercised, checked for consistancy and completeness and then generalized

multiple measurement techniques. Table 4 gives a outline of the decision making process for the experimental design expert system, and for data interpretation. Rule generation under the system is demonstrated in figure 12. The expert system strategy for the chemometrics portion of the system is similar to that described previously in the Py-MS example with the added feature that the system generate the protocol of analysis based on the initial problem statement and the available data. The system is in its earliest stages of development so any or all aspects of the proposed design are subject to modifications as experience is gained in training the system. Nevertheless, we feel that an expert systems such as ours offers a strategy for automation of laboratory instrumentation and interpretation under an expert systems approach.

```
Rule_1
  If:
      SCOPE              IS   R&D
      SAMPLE  AMT        IS   TRACE
      SAMPLE  FORM       IS   POWDER
      SAMPLE  PROCESS    IS   *
      SAMPLE  HISTORY    IS   *
      INSTR.  AVAIL      IS   ALL
  Then:
      ANAL  STRATEGY     IS   SPECTRMS(100)
```

Figure 12—Example of rule generation in EXMAT using the TIMM® expert system. Rule is from the expert system for analytical strategy.

**Table 4.** Example taken from overall organization structure of EXMAT, an expert systems for materials analysis.

| ANALYTICAL STRATEGY | DATA INTERPRETATON |
|---|---|
| **CHOICES:** | **CHOICES:** |
| CHROMGC | PATTERN RECOGNITION |
| CHROMLC | SAMPLE ID |
| SPECTRFTIR | GROUP ID |
| SPECTRMS | PEAKS ID |
| THERMTA | NO ID |
| ELEMEL | COMBINE DB |
| | EXTEND DB |
| | CORRELATION |
| **FACTORS:** | **FACTORS:** |
| 1. SCOPE | 1. DATA GENERATE |
| QUAL | FTIR DB |
| QUANT | MS DB |
| PURITY | LC DB |
| QUAL/QUANT | TA DB |
| TIME/FUND LIMIT | GC-FTIR DB |
| TRACE | GC-MS DB |
| R&D | GC DB |
| CORRELATION | 2. GC DB TREATMENT |
| SCREEN | DIRECT COMPARE |
| | CHEMOMETRICS |
| | DATA SET ID |
| | 3. FTIR DB TREATMENT |
| | DIRECT COMPARE |
| | PAIRS SEARCH |
| | 4. ETC. FOR EACH |
| | DATA BASE (DB) |

463

## References

[1] Harper, A.M., Polymer Characterization Using Gas Chromatography and Pyrolysis, S. Liebman and E.J. Levy, Marcel Dekker, Inc. (1985).

[2] Harper, A.M.; H.C.L. Meuzelaar, G.S. Metcalf, and D.L. Pope, Analytical Pyrolysis Techniques and Applications, K. Voorhees, ed. 157-195 (1984).

[3] Owens, P.M.; R.B. Lamb and T.L. Isenhour, Anal. Chem., 54 2344 (1982).

[4] McMurry, J.E.; J.A. Pino, P.C. Jurs, B. Lavine, and A.M. Harper, Anal. Chem., vol. 57, 295-302 (1985).

[5] Meuzelaar, H.L.C., and A.M. Harper, Characterization and Classification of Rocky Mountain Coals by Curie-Point Mass Spectrometry, Fuels, vol. 63, 639-652 (1984).

[6] Harper, A.M.; H.L.C. Meuzelaar, and P.H. Given, Fuels, vol. 63, 793-799 (1984).

[7] Deming, S.N.; S.L. Morgan, and M.R. Willcott, American Laboratory, October, 13 (1976).

[8] Meuzelaar, H.L.C.; J. Haverkamp, and F.D. Hileman, Curie Point Pyrolysis Mass Spectrometry of Recent and Fossil Biomaterials, Elsevier, Amsterdam (1982).

[9] Windig, W.; P.G. Kistemaker, and J. Haverkamp, J. Anal. Appl. Pyrol., 3, 199 (1982).

[10] Eshwis, W.; P.G. Kistemaker and H.L.C. Meuzelaar, Analytical Pyrolysis, CER, Jones, ed., Elsevier, Amsterdam, 151-166 (1977).

# DISCUSSION
of the Harper-Liebman paper, Intelligent Instrumentation

## Richard J. Beckman
Los Alamos National Laboratory

There has been an instrumentation revolution in the chemical world which has changed the way both chemists and statisticians think. Instrumentation has lead chemists to multivariate data—much multivariate data. Gone are the days when the chemist takes three univariate measurements and discards the most outlying.

Faced with these large arrays of data the chemist can become somewhat lost in the large assemblage of multivariate methods available for the analysis of the data. It is extremely difficult for the chemist—and the statistician for that matter—to form hypotheses and develop answers about the chemical systems under investigation when faced with large amounts of multivariate chemical data.

Professor Harper proposes an intelligent instrument to solve the problem of the analysis and interpretation of the data. This machine will perform the experiments, formulate the hypotheses, and "understand" the chemical systems under investigation.

What impact will such an instrument have on both chemists and statisticians? For the chemist, such an instrument will allow more time for experimentation, more time to think about the chemical systems under investigation, a better understanding of the system, and better statistical and numerical analyses. There would be a chemometrician in every instrument! For the statistician, the instrument will mean the removal of outliers, trimmed data, automated regressions, and automated multivariate analyses. Most important, the entire model building process will be automated.

There are some things to worry about with intelligent instruments. Will the chemist know how the data have been reduced and the meaning of the analysis? Instruments made today do some data reduction, such as calibration and trimming, and the methods used in this reduction are seldom known by the chemist. With a totally automated system the chemist is likely to know less about the analysis than he does with the systems in use today.

The statistician when reading the paper of Professor Harper probably asks what is the role of the statistician in this process? Will the statistician be replaced with a microchip? Can the statistician be replaced with a microchip? In my view the statistician will be replaced by a microchip in instruments such as those discussed by Professor Harper. This will happen with or without the help of the statistician, but it is with the statistician's help that good statistical practices will be part of the intelligent instrument.

Professor Harper should be thanked for her view of the future chemical laboratory. This is an exciting time for both the chemist and the statistician to work and learn together.

464

# The Regression Analysis
# of Collinear Data

### John Mandel

### National Bureau of Standards, Gaithersburg, MD 20899

This paper presents a technique based on the intuitively-simple concepts of Sample Domain and Effective Prediction Domain, for dealing with linear regression situations involving collinearity of any degree of severity. The Effective Prediction Domain (EPD) clarifies the concept of collinearity, and leads to conclusions that are quantitative and practically useful. The method allows for the presence of expansion terms among the regressors, and requires no changes when dealing with such situations.

## Introduction

The scientists' search for relations between measurable properties of materials or physical systems can be effectively helped by the statistical technique known as multiple regression. Even when limited to linear regression, the technique is often of great value, as we shall see below. Often, however, difficulties in interpretation arise because of a condition called collinearity. This condition, which is inherent in the structure of the design points (the $X$ space) of the regression experiment, is often treated, at least implicitly, as a sort of disease of the data that is to be remedied by special mathematical manipulations of the data.

We consider collinearity not as a disease but rather as additional information provided by the data to the data analyst, warning him to limit the use of the regression equation as a prediction tool to specific subspaces of the $X$ space, and telling him precisely what these subspaces are. Thus, collinearity is an indication of limitations inherent in the data. The statistician's task is to detect these limitations and to express them in a useful manner. If this viewpoint is adopted, there is no need for remedial techniques. All that is required is a method for extracting the additional information from the data. We will present such a method.

---

**About the Author:** John Mandel is a statistical consultant serving with NBS' National Measurement Laboratory.

---

## The Model

We assume that measurements $y$ have been made at a number of "$x$-points," each point being characterized by the numerical values of a number of "regressor-variables" $x_j$. We also assume that $y$ is a linear function of the $x$-variables. The mathematical model, for $p$ regressors, is:

$$y = \beta_1 x_1 + \beta_2 x_2 + \ldots + \beta_j x_j + \ldots + \beta_p x_p + \epsilon \qquad (1)$$

where $\epsilon$ is the error in the $y$ measurement. We denote by N the number of points, or "design points", i.e., the combinations of the $x$'s at which $y$ is measured.

*Usually, the variable $x_1$ is identically equal to "one" for all N points*, to allow for the presence of a constant term. Then the expected value of $y$, denoted $E(y)$, is equal to $\beta_1$ when all the other $x$'s are zero. This point, called the origin, is seldom one of the design points and is, in fact, quite often far removed from all design points. In many cases this point is even devoid of physical meaning.

## First Example:
## Firefly Data

We present the problem in terms of two examples of real data. The first data set (Buck [1][1]) is shown in table 1. It consists of 17 points and has two regressors, in addition to

---

[1]Figures in brackets indicate literature references.

**Table 1.** Data for firefly study.

| $x_1$ | $x_2$ | $x_3$ | $y$ |
|---|---|---|---|
| 1 | 26 | 21.1 | 45 |
| 1 | 35 | 23.9 | 40 |
| 1 | 40 | 17.8 | 58 |
| 1 | 41 | 22.0 | 50 |
| 1 | 45 | 22.3 | 31 |
| 1 | 55 | 23.3 | 52 |
| 1 | 55 | 20.5 | 54 |
| 1 | 56 | 25.5 | 38 |
| 1 | 70 | 21.7 | 40 |
| 1 | 75 | 26.7 | 28 |
| 1 | 79 | 25.0 | 38 |
| 1 | 87 | 24.4 | 36 |
| 1 | 100 | 22.3 | 36 |
| 1 | 100 | 25.5 | 46 |
| 1 | 110 | 26.7 | 40 |
| 1 | 130 | 25.5 | 31 |
| 1 | 140 | 26.7 | 40 |

Definition of Variables

$y$ = time of first flash (number of minutes after 6:30 p.m.)

$x_2$ = light intensity (in metercandles, mc)

$x_3$ = temperature (°C)

a constant term ($x_1 \equiv 1$). The measurement is the time of the first flash of a firefly, after 6:30 p.m. It is studied as a function of ambient light intensity ($x_2$) and temperature ($x_3$).

Figure 1 is a plot of $x_3$ versus $x_2$. There is obviously a trend: $x_3$ increases as $x_2$ increases. The existence of a rela-

tion of this type between some of the regressor variables often causes difficulties in the interpretation of the regression analysis. To deal with the problem in a general way we propose a method based on two concepts. The first of these we shall call the "sample domain."

For our data, the sample domain consists of the rectangle formed by the vertical straight lines going through the lowest and highest $x_2$ of the experiment, respectively, and by the horizontal straight lines going through the lowest and highest $x_3$, respectively (See Fig. 1). The concept is readily generalized to an $X$ space of any number of dimensions, and becomes a hypercube in such a space. Note that the vertex $B$ of the sample domain is relatively far from any of the design points. This has important consequences.

The regression equation

$$\hat{y} = \hat{\beta}_1 \cdot x_1 + \hat{\beta}_2 \cdot x_2 + \hat{\beta}_3 \cdot x_3 \qquad (2)$$

allows us to estimate $y$ at any point ($x_1, x_2, x_3$) (we recall that $x_1 = 1$) and to estimate the variance of $\hat{y}$ at this point. The point can be inside or outside the sample domain. Obviously the variance of $\hat{y}$, which we denote by Var ($\hat{y}$), will tend to become larger as the point for which the prediction is made is further away from the cluster of points involved in the experiment. Therefore Var ($\hat{y}$) at the point $B$ may be considerably larger than at points $A$, $C$, and $D$. Such a condition is associated with the concept of "*collinearity*." We define collinearity, in a semi-quantitative way, as the condition that arises when for at least one of the vertices of the sample



Figure 1—Sample domain.

domain, Var $(\hat{y})$ is considerably larger than for the other vertices. The concept will become clearer as we proceed.

At any rate, the larger variance at one of the vertices of the sample domain is generally the lesser of two concerns, the other being that the regression equation, for which validity may have been reasonably firmly established in the vicinity of the cluster of experimental points, may no longer be valid at a more distant point. It is important to note that the evidence from the data alone cannot justify inferences at such distant points. In order to validate prediction at such points, it is necessary to introduce either additional data or additional assumptions.

For these reasons, we seek to establish a region in the X-space for which prediction is reasonably safe *on the basis of the experiment alone*. We call this the *Effective Prediction Domain*, or EPD.

The EPD is the second concept required for our treatment of collinear data. It is closely related to the first concept, the sample domain, as will be shown below.

## Establishing the EPD

Our procedure consists of two steps, involving two successive transformations of the coordinate system. The original coordinate system in which the x-regressors are expressed is referred to as the *X-system*.

### 1. The Z System

The first step consists in a *translation* of the X-system (parallel to itself) to a different origin, located centrally within the cluster of experimental points (*centering*); and simultaneously by a *rescaling* of each x to a standard scale. The new system, called the *Z-system*, is given by the equations[2]

For $j = 1$: $z_1 = K$ (a constant) (3a)

For $j > 1$: $z_j = \dfrac{x_j - C_j}{R_j}$ (3b)

For $C_j$ and $R_j$ we consider two choices, which we call the Correlation Scale Transformation (CST) and the Range Midrange Transformation (RMT). We discuss first the Correlation Scale Transformation defined by the choice

$$C_j = \bar{x}_j, \quad R_j = \sqrt{\sum_i (x_{ij} - \bar{x}_j)^2}$$ (4)

where $i = 1$ to $N$.

It easily follows from (3b) that

---

[2]We assume that in the X-system, the regressor $x_1$ is identically equal to *unity*, to allow for an independent term.

$$\bar{z}_j = 0, \quad \sum_i z_{ij}^2 = 1$$ (5)

It is then reasonable to choose a value $K$ in (3a) equal to

$$K = 1/\sqrt{N}$$ (6)

so as to make $\displaystyle\sum_i z_{i1}^2 = 1$

The values of $C_j$ and $R_j$ for the firefly data are given in table 2. Contrary to statements found in the literature (see discussion at end of this paper), the centering and rescaling defined by the Correlation Scale Transformation have no effect whatsever on collinearity. The location of the sample domain relative to the design points remains unchanged, though it is expressed in different coordinates.

To arrive at an EPD, a second operation is necessary, viz. a *rotation* of the Z-coordinate system to a new coordinate system, which we shall call the *W-system* (of coordinates).

### 2. The W-System

The rotation from Z to W is accomplished by the method of Principal Components, or its equivalent, the Singular Value Decomposition (SVD). For a discussion of this method the reader is referred to Mandel [2]. Here we merely recall a few facts. Each w-coordinate is a linear combination of all z-coordinates given by the matrix equation:

$$W = Z V$$ (7)

where $V$ is an orthogonal matrix.

In algebraic notation, eq (7) becomes

$$w_{ik} = \sum_j z_{ij} v_{kj} \qquad \begin{array}{l} i = 1 \text{ to } N \\ j = 1 \text{ to } p \end{array}$$ (8)

where the $v_{kj}$ are the elements of the $V$ matrix. The $v_{kj}$, for a given $k$, are simply the direction cosines of the $w_k$ axis with respect to the Z-system. Consequently,

$$\sum_j v_{kj}^2 = 1$$ (9)

Table 2. Firefly data—parameters for correlation scale transformation.

| j | C | R |
|---|---|---|
| 1 | 0 | 4.123106 |
| 2 | 73.176471 | 135.264447 |
| 3 | 23.582353 | 10.073962 |

Since the rotation is orthogonal, any two distinct $w$-axes, say $w_k$ and $w_{k'}$, are orthogonal and consequently:

$$\sum_j v_{kj} \cdot v_{k'j} = 0 \qquad \text{for } k \neq k' \qquad (10)$$

For the firefly data, the $V$ matrix is shown in table 3, and the complete set of $z$ and $w$ coordinates is given in table 4.

Note that row 2, as well as column 1, in table 3 consists of the element "one" in one cell and zeros in all others cells. This is a consequence of the orthogonality of $z_1$ with respect to all $z_j$ with $j > 1$. This orthogonality is in turn due to the nature of the Correlation Scale Transformation, as expressed by eq (4).

At the bottom of the $w$ columns we find values labeled $\lambda_j$. They are simply the sums of squares of all $w$-values in that column.

$$\lambda_j = \sum_i w_{ij}^2 \qquad (11)$$

**Table 3.** Firefly data—V matrix.

| | *j* | | |
|---|---|---|---|
| *k* | *1* | *2* | *3* |
| 1 | 0 | .7071 | .7071 |
| 2 | 1.000 | 0 | 0 |
| 3 | 0 | −.7071 | .7071 |

**Table 4.** Firefly data—$z$ and $w$ coordinates (CST).[1]

| Point | $z_2$ | $z_3$ | $w_1$ | $w_3$ |
|---|---|---|---|---|
| 1 | −.3488 | −.2464 | −.4216 | .0724 |
| 2 | −.2822 | .0315 | −.1780 | .2219 |
| 3 | −.2453 | −.5740 | −.5800 | −.2324 |
| 4 | −.2379 | −.1571 | −.2800 | .0572 |
| 5 | −.2083 | −.1273 | −.2381 | .0573 |
| 6 | −.1344 | −.0280 | −.1156 | .0753 |
| 7 | −.1344 | −.3060 | −.3121 | −.1213 |
| 8 | −.1270 | .1904 | .0440 | .2245 |
| 9 | −.0235 | −.1869 | −.1495 | −.1155 |
| 10 | .0135 | .3095 | .2276 | .2094 |
| 11 | .0431 | .1407 | .1292 | .0691 |
| 12 | .1022 | .0812 | .1289 | −.0148 |
| 13 | .1983 | −.1273 | .0495 | −.2302 |
| 14 | .1983 | .1904 | .2741 | −.0055 |
| 15 | .2722 | .3095 | .4106 | .0264 |
| 16 | .4201 | .1904 | .4309 | −.1624 |
| 17 | .4940 | .3095 | .5674 | −.1304 |
| | | | $\lambda_1 = 1.6549$ | $\lambda_3 = .3451$ |

[1] $z_1 = 1/\sqrt{17} = .2425$ for all $i$
$w_2 = 1/\sqrt{17} = .2425$ for all $i$, $\lambda_2 = 1.0000$

The $\lambda_j$ are also the *eigenvalues* of the $Z'Z$ matrix which, for our choice of $C_j$ and $R_j$, is the correlation matrix of the regressors $x$. Note that $w_2$ is the constant $= 1/\sqrt{N}$. Consequently

$$\lambda_2 = N\left(\frac{1}{\sqrt{N}}\right)^2 = 1 \ .$$

We need to consider $w_1$ and $w_3$ only. A similar situation applied to the $z$ coordinates, where $z_1 = 1/\sqrt{N}$ for all $i$. Figure 2 shows both the $z$-coordinates ($z_2$ and $z_3$) and the $w$-coordinates ($w_1$ and $w_3$) for the firefly data. The order of the $w$-coordinates ($w_1$, $w_2$, $w_3$) is that of the corresponding $\lambda$-values, in decreasing order.

## 3. The Effective Prediction Domain (EPD)

The EPD is simply the sample domain corresponding to the $W$-system of coordinates. Thus, straight lines parallel to the $w_3$-axis are drawn through the smallest and largest $w_1$, respectively, and lines parallel to the $w_1$-axis are drawn through the smallest and largest $w_3$. Here again generalization is readily made to a $p$-dimensional $W$-space. The EPD for the firefly data is also shown in figure 2.

The interpretation of EPD is straightforward. Unlike the sample domain in either the $X$-system or the $Z$-system, the EPD excludes points that are distant from the cluster of regressor points. This has two advantages. In the first place, the use of the regression equation is justified for all points inside, and on the periphery of the EPD. And accordingly, the variance of the predicted value $\hat{y}$ for any such point will not be unduly large. These statements require more detailed treatment. To this effect we introduce the concept of *variance factor* (VF).

## 4. The Variance Factor (VF)

From regression theory we know that the variance of any linear functon, say $L$, of the coefficient estimates $\hat{\beta}_j$ is of the form:

$$\text{Var } (L) = f(X) \cdot \sigma_\epsilon^2 \qquad (12)$$

where $\sigma_\epsilon^2$ is the variance of the experimental errors $\epsilon$ of the $y$ measurements. The multiplier $f(X)$ is independent of the $y$ and depends only on the $X$ matrix and on the coefficients in the $L$ function. We call this multiplier the *variance factor*, VF.

Thus, we have:

$$\text{Var } (\hat{\beta}_j) = \text{VF}(\hat{\beta}_j) \cdot \sigma_\epsilon^2 \qquad (13)$$

and

Figure 2—EPD for firefly data.

*Light Intensity (X₂)*

$$\text{Var } (\hat{y})=\text{VF}(\hat{y})\cdot\sigma_\epsilon^2 \qquad (14)$$

In eq (14), $\hat{y}$ is the estimated, or *predicted* $y$ value at any chosen point in $X$-space. VF $(\hat{y})$ is of course a function of the location of this point.

Returning now to our statements above, it is well-known that a regression equation can show excellent (very small) residuals and yet be very poor for certain prediction purposes. The small residuals merely mean that a good fit has been obtained *at the points used in the experiment*. This is no guarantee that the fit is good at other points. However, if the regression equation is scientifically reasonable, it is likely that the experimental situation underlying it will also be valid for points *that are close* to the cluster of the regressor points used in the experiment. Every point in the EPD satisfies this requirement.

Furthermore, the variance of prediction, measured by the VF, will also be reasonably small for all points of the EPD,

simply because they are geometrically close to the design points.

The calculation of VF $(\hat{y})$ is quite simple, once the V-matrix and the $\lambda$ values have been calculated. It is based on the equation

$$\text{VF}(\hat{y})=\sum_k u_k^2 \qquad (15)$$

where $u_k$ is defined as:

$$u_k=\frac{w_k}{\sqrt{\lambda_k}} \qquad (16)$$

Combining eqs (8) and (16), we obtain

$$u_{ik}=\sum_j z_{ij}\frac{v_{kj}}{\sqrt{\lambda_k}} \qquad (17)$$

469

and hence:

$$VF(\hat{y})=\sum_k \frac{\left(\sum_j z_{ij}v_{kj}\right)^2}{\lambda_k} \qquad (18)$$

Figure 3 shows the VF values at the vertices of the original sample domain and of the EPD. Interpreting these results, we see that the collinearity of our data is reflected in the rejection of an appreciable portion of the sample domain for purposes of safe prediction. This does *not* mean that prediction outside the EPD is impossible, or unacceptable. It merely means that such prediction cannot be justified on the basis of the data alone. Of course, the risk of predicting outside the EPD increases with the distance from the EPD. It will generally be reasonably safe to use the regression equation even outside the EPD, as long as the point for which prediction is made is reasonably close to the borders of the EPD. Using eq (18), the VF for any contemplated prediction point is readily calculated and can serve as a basis for decision.

## Second Example: Calibration for Protein Determination

The instructive and intuitively satisfying graphical display of the EPD becomes impossible when the number of regressors, including the independent term, exceeds 3. We must then replace the graphical procedure by an analytical one, as will now be shown in the treatment of our second example.

The data were presented by Fearn [3], in a discussion of Ridge Regression. They represent the linear regression of percent protein, in ground wheat samples, on near-infrared reflectance at six different wavelengths.

For reasons of simplicity in presentation, we include here only three of the six wavelengths, a change that has a rather small effect on the final outcome of the analysis: it turns out that the regression equation based on these 3 wavelengths is very nearly as precise as that based on 6 wavelengths.

The data, displayed in table 5, are a very good example of the use of regression equations: the regression equation is indeed to be used as a "calibration curve" for the analysis of protein, using the rapid spectrometry instead of the far more time-consuming Kjeldahl nitrogen determination. Our data have an $N$ value of 24, and $p$ (including the independent term) is 4.

Table 6 exhibits the correlation matrix of the 24 design points. It is very apparent that the $x$ values at all three wavelengths are highly correlated with each other, thus indicating a high degree of collinearity. At a first glance one would be very skeptical about such a set of data, and suspect that the $X$ matrix shows such a high degree of redundancy as to make the regression useless for prediction purposes. Fearn explains that the correlations are more a reflection of particle size variability than of protein content. Our analysis will confirm that, properly interpreted, the data lead to a very satisfactory calibration procedure.

We will find it useful to introduce a slightly different $Z$ transformation, which we call the *Range-Midrange Transformation*.



**Sample Domain**

| Vertex | VF |
|--------|------|
| A | .39 |
| B | 1.71 |
| C | .69 |
| D | .30 |

**EPD**

| Vertex | VF |
|--------|------|
| a | .41 |
| b | .41 |
| c | .41 |
| d | .40 |

Figure 3—VF at vertices of sample domain and of EPD.

470

**Table 5.** Protein Calibration Data[*]

| Point | Reflectance $x_2$ | $x_3$ | $x_4$ | % Protein $y$ |
|---|---|---|---|---|
| 1 | 246 | 374 | 386 | 9.23 |
| 2 | 236 | 386 | 383 | 8.01 |
| 3 | 240 | 359 | 353 | 10.95 |
| 4 | 236 | 352 | 340 | 11.67 |
| 5 | 243 | 366 | 371 | 10.41 |
| 6 | 273 | 404 | 433 | 9.51 |
| 7 | 242 | 370 | 377 | 8.67 |
| 8 | 238 | 370 | 353 | 7.75 |
| 9 | 258 | 393 | 377 | 8.05 |
| 10 | 264 | 384 | 398 | 11.39 |
| 11 | 243 | 367 | 378 | 9.95 |
| 12 | 233 | 365 | 365 | 8.25 |
| 13 | 288 | 415 | 443 | 10.57 |
| 14 | 293 | 421 | 450 | 10.23 |
| 15 | 324 | 448 | 467 | 11.87 |
| 16 | 271 | 407 | 451 | 8.09 |
| 17 | 360 | 484 | 524 | 12.55 |
| 18 | 274 | 406 | 407 | 8.38 |
| 19 | 260 | 385 | 374 | 9.64 |
| 20 | 269 | 389 | 391 | 11.35 |
| 21 | 242 | 366 | 353 | 9.70 |
| 22 | 285 | 410 | 445 | 10.75 |
| 23 | 255 | 376 | 383 | 10.75 |
| 24 | 276 | 396 | 404 | 11.47 |

[*] $x_1 = 1$

**Table 6.** Protein calibration data—correlation matrix of $x_1$ through $x_4$.

| 1 | 0 | 0 | 0 |
|---|---|---|---|
|  | 1 | .9843 | .9337 |
|  |  | 1 | .9545 |
|  |  |  | 1 |

## The Range-Midrange Transformation

The Range-Midrange Transformation (RMT) is defined as follows:

$$\text{For } j = 1: \quad z_1 = 1 \tag{19a}$$

$$\text{For } j > 1: \quad z_j = \frac{x_j - C_j}{R_j} \tag{19b}$$

but now $C_j$ is defined as the *midrange* of the $N$ values of $x_j$ and $R_j$ is *one-half the range* of these values. With these definitions, it is clear that the smallest $z$-value, for any regressor, is $(-1)$ and the largest $z$-value is $(+1)$. It is because of this $-1$ to $+1$ scale that this transformation was introduced. The benefits of this scale will become apparent in the following section.

## EPD for the Protein Data

The EPD resulting from the Singular Value Decomposition based on the Range-Midrange Transformaton will not be he same as the EPD we would have obtained using the Correlation Scale Transformation, but we will see that those features of the EPD that are of importance for us, in establishing the limitations of the regression equation, are practically unaffected.

Table 7 shows the $C$ and $R$ values for the four regressors and table 8 exhibits the $V$ matrix and the $\lambda$ values obtained from the Singular Value Decomposition. The latter, it may be recalled, simply expresses the rotation of the $Z$ coordinate system to the $W$ system.

For each $w_k$ coordinate, there are 24 values, corresponding to the 24 regressor points.

Table 9 shows the smallest and the largest $w_k$ value, for each of the four $k$.

According to table 9, we must have, in the EPD:

$$-1.9282 \leqq w_1 \leqq .6181 \tag{20}$$

with similar statements for $w_2$, $w_3$, and $w_4$. Applying now eq

**Table 7.** Protein calibration data—parameters for Z transformation (RMT).

| $j$ | $C$ | $R$ |
|---|---|---|
| 1 | 0 | 1 |
| 2 | 296.5 | 63.5 |
| 3 | 418.0 | 66.0 |
| 4 | 432.0 | 92.0 |

**Table 8.** Protein calibration data—V matrix and $\lambda$ values (RMT).

| $k$ | 1 | 2 | 3 | 4 | $\lambda$ |
|---|---|---|---|---|---|
| 1 | −.6665 | .4845 | .4217 | .3784 | 43.7810 |
| 2 | .7365 | .3299 | .3797 | .4523 | 8.3782 |
| 3 | −.1096 | −.5491 | −.2509 | .7896 | .3758 |
| 4 | −.0332 | −.5958 | .7843 | −.1698 | .06624 |

**Table 9.** Protein calibration data—limits defining the EPD.

| Coordinate (k) | Smallest w | Largest w |
|---|---|---|
| 1 | −1.9282 | .6181 |
| 2 | −.4097 | 1.8989 |
| 3 | −.1669 | .3158 |
| 4 | −.0801 | .1324 |

471

(8), this double inequality can be written:

$$-1.9282 \leqq -.6665\ z_1 + .4845\ z_2 + .4217\ z_3 + .3784\ z_4$$

$$\leqq .6181$$

Since $z_1$ is constant and $= 1$, this double inequality becomes:

$$-1.2617 \leqq .4845\ z_2 + .4217\ z_3 + .3784\ z_4 \leqq 1.2846\ .\quad (21a)$$

With the RMT, the value of any $z_k$ is, for any $k > 1$, between $(-1)$ and $(+1)$. Thus the expression in the middle has, for all design points, a value between $-1.2846$ and $1.2846$, where $1.2846$ is the sum of the absolute values of the three coefficients. Therefore, the double inequality expressed by eq (21a) holds, essentially, for every point in the original sample domain. Thus, $w_1$, the first coordinate of the EPD, which represents its largest dimension, imposes essentially no restrictions on the sample domain.

Doing the same calculations for the three other $w$-coordinates (see table 9), we obtain, respectively:

$$-1.1462 \leqq .3299\ z_2 + .3797\ z_3 + .4523\ z_4 \leqq 1.1619 \quad (21b)$$

$$-0.568 \leqq -.5491\ z_2 - .2509\ z_3 + .7896\ z_4 \leqq .4254 \quad (21c)$$

$$-.0469 \leqq -.5958\ z_2 + .7843\ z_3 - .1698\ z_4 \leqq .1656. \quad (21d)$$

We see that $w_2$ too, imposes only very light restrictions on the sample domain. On the other hand, $w_3$ and $w_4$ do imply limitations that eliminate appreciable portions of the sample domain from the EPD.

We could readily convert eqs (21c) and (21d) to $x$ coordinates by means of table 7 and eqs (19a) and (19b), but the $z$-coordinates, using the Range-Midrange Transformation, are more readily interpreted in terms of the severity of collinearity than the $x$-coordinates.

Thus, the sum of the absolute values of the coefficients in the middle terms of (21c) and (21d) are $1.5896$ and $1.5499$, respectively. Points for which these linear combinations take the valves $\pm 1.5896$ and $\pm 1.5499$ exist in the original sample domain. The EPD, on the other hand, limits these functions to intervals with much narrower limits.

## Effect of Type
## of Z Transformation

We have used two different $Z$ transformations, the Correlation Scale, and the Range-Midrange. It is proper to ask how our results would have been affected in the Protein Calibration Data, had we used Correlation Scale, instead of the Range-Midrange Transformation. We show the com-

| $w$ coordinate | Z Transf. | Inequalities |
|---|---|---|
| 1 | CST | $-3.034 \leq 1.021\ z_2 + 1.061\ z_3 + z_4 \leq 3.082$ |
| | RMT | $-3.334 \leq 1.280\ z_2 + 1.114\ z_3 + z_4 \leq 3.395$ |
| 2 | CST | |
| | RMT | $-2.534 \leq .729\ z_2 + .840\ z_3 + z_4 \leq 2.569$ |
| 3 | CST | $-.075 \leq -.686\ z_2 - .321\ z_3 + z_4\ .535$ |
| | RMT | $-.072 \leq -.695\ z_2 - .318\ z_3 + z_4 \leq .539$ |
| 4 | CST | $-.278 \leq -3.531\ z_2 + 4.640\ z_3 - z_4 \leq .980$ |
| | RMT | $-.276 \leq -3.509\ z_2 + 4.619\ z_3 - z_4 \leq .975$ |

[1]All inequalities are expressed in RMT $z$ coordinates.

parison in table 10. Let us recall that with the CST, one of the $w$ coordinates yields a $\lambda$-value of unity, and a constant $w$ value for all points. Therefore, we obtain for CST, only three sets of inequalities, as compared to the four sets for RMT. To allow the comparison between the two transformation to be made, we have multiplied eqs (21a) through (21d) by positive constants, so as to make the coefficient of $z_4$ equal to $\pm 1$. The same was done for the corresponding inequalities obtained by the Correlation Scale Transformation.

Of course, since the $z$ coordinates are different for the two transformations, the inequalities for the CST, expressed in the CST $z$-units, had to be converted to RMT $z$-units, for a meaningful comparison. As can be seen from table 10, the two smallest dimensions of the EPD are practically the same for the two transformations. Thus, even though the method of principal components is not invariant with respect to linear transformations of scale, our analysis leads, in this case, to very similar results for the small dimensions of the EPD. We believe that this is generally true for all situations in which collinearity is noticeable, i.e., for all situations in which the EPD eliminates considerable portions of the original sample domain. For situations in which this does not apply, i.e., totally non-collinear cases, the inequalities do not matter, since they impose no restrictions on the sample domain.

It is interesting to contrast the remarkable similarity between the inequalities for $w_3$ and $w_4$ for the two transformations in table 10, with the behavior of a commonly advocated measure of collinearity (Belsley, Kuh, and Welsch [4], the *condition-number*.

The SVD resulting from the CST yields the following eigenvalues: $2.9151, 1.0000, .07176, .01312$. The condition number is defined as the ratio of the largest to the smallest eigenvalue. In this case:

condition number$= 2.9151/.01312 = 222.2$

On the other hand, the SVD resulting from the RMT on the same data yields the eigenvalues: $43.7810, 8.3782, .37575, .066244$. This time we have:

472

Thus the condition number varies considerably when the data are subjected to different standardizing transformations. It is not clear what useful information can be derived from the condition number.

By contrast, the treatment of collinearity we advocate has a useful and readily understood interpretation: the EPD is that part of the $X$ space in which, and near which, prediction is safe. It also indicates what portions of the original sample domain are inappropriate for prediction *on the basis of the given data alone*. It fulfills this function in a way which is practically invariant with respect to *intermediate* transformations of scale. We use the qualifier "intermediate" because collinearity has meaning only in terms of a given original coordinate system (the $X$ system). This system, which determines the original sample domain, must be considered fixed. On the other hand, transformations of this system prior to calculating the EPD can be defined in different ways without affecting the practical inferences drawn from the data on the basis of the final EPD derived form the standardizing transformation.

## Cross-Validation

We can take advantage of the availability of a second set of protein calibration data, also given in Fearn [3]; to verify the correctness of our approach. Fearn lists 26 additional points for which the reflectance measurements, as well as the Kjeldahl nitrogen determination, were made. We applied the $Z$ transformation obtained above (RMT on first set of 24 points) to each of these 26 points, and noted every point for which at least one of the four sets of inequalities (21a) through (21d) failed to be satisfied. We found 14 such points. This means that 14 "future points" obtained under the same test conditions were outside the EPD established on the basis of the original 24 points. However, as we observed above, as long as the point is not far from the EPD, prediction at that point is likely to be valid. We tested "predictability" at these 14 points by calculating the VF value for each of them, and by comparing the predicted protein value with the measured one. The results are shown in table 11. It is apparent that all VF are relatively small, indicating that even though these 14 points are outside the EPD calculated from the original set, they are not far from that EPD. This is confirmed by the good agreement between the observed and predicted values. The standard deviation of fit for the original set of 24 points was 0.23; the standard deviation for a single measurement derived from the 14 differences in table 11 is 0.30.

## Expansion Terms

Quite frequently, a regression equation contains $x$ variables that are non-linear functions of one or more of the

Table 11. Protein calibration data—cross-validation of analysis.

| Point[1] | % Protein Observed | Predicted | VF |
|---|---|---|---|
| 1 | 8.66 | 9.53 | .281 |
| 4 | 11.77 | 11.97 | .416 |
| 6 | 10.46 | 10.96 | .193 |
| 9 | 12.03 | 11.47 | .212 |
| 10 | 9.43 | 9.54 | .762 |
| 11 | 8.66 | 8.15 | .454 |
| 12 | 14.44 | 13.99 | .881 |
| 14 | 10.41 | 10.17 | .468 |
| 16 | 11.69 | 11.24 | .472 |
| 17 | 12.19 | 11.83 | .390 |
| 18 | 11.59 | 11.39 | .314 |
| 20 | 8.60 | 8.39 | .201 |
| 22 | 9.34 | 8.93 | .151 |
| 26 | 10.89 | 10.94 | .741 |

[1]Point in additional set (Fearn [3]) with its number designation in that set.

other $x$ variables, such as $x_2^2$, $x_2 \cdot x_3$, etc. Polynomial regressions are necessarily of this type. Since the $x$ variables are non-stochastic in the usual regression models, the least squares solution for the regression equation is not affected by the presence of such "expansion terms." On the other hand, collinearity can be introduced, or removed, or modified by them.

In our treatment the expansion terms cause no additional problems. Consider for example, the regression

$$y = \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_2^2 + \epsilon \qquad (22)$$

with $x_1 \equiv 1$.

Here we have $p = 3$. Using RMT, followed by a singular value decomposition, we obtain an EPD of three dimensions, leading to the inequalities.

$$A_1 \leq w_1 \leq B_1, \quad A_2 \leq w_2 \leq B_2, \quad A_3 \leq w_3 \leq B_3 . \qquad (23)$$

Expressing the $w$ as functions of the $z$, this leads to three double inequalities governing the $z$, of the form

$$A_1 \leq f_1(z) \leq B_1, \quad A_2 \leq f_2(z) \leq B_2, \quad A_3 \leq f_3(z) \leq B_3, \qquad (23)$$

Now, since $x_3 = x_2^2$, we have

$$z_3 = \frac{x_3 - C_3}{R_3} = \frac{x_2^2 - C_3}{R_3} = \frac{(R_2 z_2 + C_2)^2 - C_3}{R_3} .$$

Hence:

$$z_3 = \frac{C_2^2 - C_3}{R_3} z_1 + \frac{2\,C_2 R_2}{R_3} z_2 + \frac{R_2^2}{R_3} z_2^2 . \qquad (24)$$

Because of this relation the functions $f_1(z)$, $f_2(z)$, $f_3(z)$ become functions of $z_1$, $z_2$ (and $z_2^2$) only. Using this fact, we

interpret the three sets of inequalities (23) exactly as we have interpreted eqs (21a) through (21d) by determining which of these inequalities, if any, impose restrictions on the use of the original sample domain.

To illustrate this procedure, consider the small set of artificial data shown in table 12, for which the model is given at the bottom of the table. The term $x_3 = x_2^2$ introduces a high correlation between $x_2$ and $x_3$ and consequently also considerable collinearity.

The inequalities characterizing the EPD based on a Range-Midrange Transformation and converted to the $z$-scales, are shown in table 13. Applying eq (24) to express $z_3$ in terms of $z_2$, the three double-inequalities become:

for $w_1$: $-.8431 \leq 1.1284\ z_2 + .2853\ z_2^2 \leq 1.4137$
for $w_2$: $-.6928 \leq .8541\ z_2 + .1612\ z_2^2 \leq 1.0153$
for $w_3$: $.0077 \leq .0003\ z_2 + .3218\ z_2^2 \leq .3221.$

It is readily verified that of these six inequalities, all but one are satisfied for all $z_2$ values between $-1$ and $+1$. The last one, involving the left side of the third set, is satisfied for all $z_2$ values except for the interval: $-.156 \leq z_2 \leq .155$. This corresponds to an $x_2$ interval between 2.1 and 2.8, or between the design points $x_2 = 2.1$ and $x_2 = 3.6$ (see table 12). The interpretation of this finding is that while all design points are of course inside the EPD, a small portion of the curve $x_2^2$ versus $x_2$ falls slightly outside the EPD. This is of no practical significance since the VF for these points, even though they are outside the EPD, does not exceed 0.58. By comparison, the smallest VF value along the curve, for the range $x_2 = .2$ to $x_2 = 4.7$, is of the order of 0.26. Thus we see that the serious collinearity in this data set is merely a consequence of the presence of the expansion term $x_3 = x_2^2$.

Table 12. An artificial quadratic example[1].

| Point | $x_2$ | $x_3$ | $y$ |
|-------|-------|-------|-----|
| 1 | .2 | .04 | 28.3 |
| 2 | .4 | .16 | 27.5 |
| 3 | 1 | 1.00 | 25.6 |
| 4 | 2.1 | 4.41 | 28.7 |
| 5 | 3.6 | 12.96 | 46.4 |
| 6 | 4.7 | 22.09 | 69.8 |

[1] $y = \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \epsilon$; $\beta_1 = 30$, $\beta_2 = 8$, $\beta_3 = 3.5$, $\sigma_\epsilon = 0.2$ $x_1 = 1$.
Note that $x_3 = x_2^2$.

Table 13. Quadratic example—inequalities for EPD.

| w-coordinate | Inequalities |
|--------------|--------------|
| $w_1$ | $-1.1281 \leq -.5070\ z_2 + .6214\ z_3 \leq 1.1284$ |
| $w_2$ | $-.8541 \leq .5029\ z_2 + .3511\ z_3 \leq .8541$ |
| $w_3$ | $-.3141 \leq -.7005\ z_2 + .7008\ z_3 \leq .0003$ |

Any point in X space, in order to be acceptable, must lie on the curve $x_3 = x_2^2$. An $x_3$ with any other value is obviously not valid and our analysis of the data, through the EPD, calls attention to this fact: in the direction of $w_3$, the width of the EPD is only .31 as compared with widths of 2.26 and 1.71 for $w_1$ and $w_2$.

## Discussion

The common mathematical definition of collinearity is the existence of at least one linear relation between the $x$'s, of the form

$$\sum_j c_j x_{ij} = 0 \qquad j = 1 \text{ to } p \qquad (25)$$

where the $c_j$ are not all zero, and such that eq (25) holds with the same $c_j$ values, for all $i$. This defines what we shall call "exact collinearity." Geometrically, it means that all design points lie in an hyperplane of the $x$-space, going through the origin of the coordinate system. Equation (25) also implies that the matrix $X'X$ is singular, and consequently that the estimates of the $\beta$ coefficients are not uniquely defined.

Exact collinearity seldom occurs in real experimental situations; indeed, if the $X$ matrix is not the result of a designed experiment, it is highly improbable that a relation such as eq (25) would hold exactly. If, on the other hand, the experiment is designed, care would generally have been taken to avoid a situation of exact collinearity.

While exact collinearity is practically of little concern, near-collinearity is a frequent occurrence in real-life data. This occurs when an equation such as (25) is "approximately" true for all $i$. Many attempts have been made to define more closely the concept of near-collinearity, but while these endeavors have led to a number of proposals for measuring collinearity, they are of little practical use to the experimenter confronted with the task of interpreting his data.

It is not our intention to discuss here the pros and cons of the various attempts made by a number of authors to "remedy" a near-collinear situation. The best-known of these remedial procedures is Ridge Regression. We merely repeat what we have said in the body of the paper: any attempt to remedy collinearity must necessarily be based on additional assumptions, unless it consists of making additional measurements. The latter alternative is of course logical and valid, but the making of assumptions invented specifically for the purpose of removing collinearity does not appear to us to be a recommendable policy in data analysis.

One easily recognizable condition leading to collinearity is the existence of at least one high correlation coefficient among the non-diagonal elements of the correlation matrix

474

of the $x$'s. This has given rise to the *concept of the Variance Inflation Factor* (VIF). The VIF for $\hat{\beta}_j$ is defined (Draper and Smith [5]), as:

$$\text{VIF}(\hat{\beta}_j) = \frac{1}{1 - R_j^2} \qquad (26)$$

where $R_j$ is the multiple correlation coefficient of $x_j$ on all other regressors. If $d$ represents a residual in this regression, the usual formula for $R_j$ is given by

$$R_j^2 = 1 - \frac{\Sigma d^2}{\sum_i (x_{ij} - \bar{x}_j)^2} . \qquad (27)$$

Now, Snee and Marquardt (Belsley [6], "comments") make, implicitly, a distinction between the two "models":

$$y = \beta_1 x_1 + \beta_2 x_2 + \cdots + \beta_p x_p + \epsilon \qquad (28a)$$

with $x_1 \equiv 1$, and

$$y - \bar{y} = \beta_2(x_2 - \bar{x}_2) + \cdots + \beta_p(x_p - \bar{x}_p) + \epsilon \qquad (28b)$$

where (28b) is called the "centered" model. For (28b), Snee and Marquardt use eq. (27), but for (28a) they appear to use the definition:

$$R_j^2 = 1 - \frac{\Sigma d^2}{\sum_i x_{ij}^2} . \qquad (29)$$

Equation 29, in which the denominator of the last term is not centered, is not explicitly given by Snee and Marquardt, but is implied by their statement:

"If the domain of prediction includes the full range from the natural origin through the range of the data, then collinearity diagnostics should not be mean-centered," and confirmed by the VIF values given in their table 1. In this table, "no centering" results in VIF values of 200,000 and 400,000, while the VIF for the "centered" data are unity. The quoted statement occurs in a section entitled "Model building must consider the intended or implied domain of prediction." The basic idea underlying the section in question is that the analysis of the data, based on the "collinearity diagnostics" (specifically: the VIF values), is goverened by the location of the points were one *wishes* to make predictions and, more specifically, on whether the origin ($x_1 = 1$, $x_2 = x_3 \cdots = 0$) is such a point. The VIF values which, according to Snee and Marquardt's formu-

las, depend heavily on whether or not this origin is included, will then indicate the quality of the predicted values.

A more reasonable approach, and one more consistent with the procedures commonly used by scientists, is to limit prediction to the vicinity of where one made the measurements, *unless additional information is available* that justifies extrapolation of the regression equation to more distant points of the samples space. The vicinity of the measured points is determined by the EPD which, in the case of collinearity, may be considerably smaller than the sample domain. In this view, it is the location of the *design points*, rather than that of the *intended points of prediction*, that determines predictability. The latter is measured, not by VIF values, but rather by the more concrete VF values, for any desired point of prediction.

The view advocated by Snee and Marquardt sometimes results in an enormous difference in the VIF values between the centered and non-centered forms. Equation 29 serves no useful purpose and is, in fact, unjustified and misleading. It is unjustified because it not only includes the origin ($x_1 = 1$, $x_k = 0$ for $k > 1$) in the correlation and VIF calculations, but moreover, gives this point infinite weight in these calculatons. Yet, no measurement was made at that point. Equation 29 is also misleading because it leads to very large VIF values for some non-centered regressions, implying that severe "ill-conditioning" exists, even when the $X$ matrix is except for some trivial coding, completely orthogonal (cf. [6]).

The ill-conditioning exists only in terms of the large VIF value. It is an artifact arising from the desire to make the two forms of the regression equation into two distinct "models".

The two forms, eqs 28a and 28b lead to identical estimates for the $\beta_j$, including $\beta_1$, and for their standard errors. They also lead to identical values and variances for an estimated (predicted) $\hat{y}$, at any point of the $X$ space. There seems to be no valid reason for the two distinct equations for the VIF. They only lead to the false impression that centering can reduce or even remove collinearity.

Our viewpoint in this paper is that the usefulness of a regression equation lies in its abilty to "predict" $y$ for *interesting combinations of the $x$'s. We also take the position* that inferences *from the data alone* should be confined to $x$ points that are in the general geometric vicinity of the cluster of design points. An inference for points that are well outside this domain (i.e., outside a suitably defined EPD) is, in the absence of additional information, only a tentative conclusion, and not a valid scientific inference. Such conclusions may however, be very useful, provided their tentative character is recognized, and provided they are subsequently subjected to further experimental verification.

Daniel and Wood [7] discuss briefly the relation between the variance of $\hat{y}$ and the location of the point at which the prediction is made. However, their discussion is in the con-

475

text of selecting the best subset of regressors from among the entire set of regressors, a subject different from the one dealt with in this paper.

Another publication that deals explicitly with predictability is a paper by Willan and Watts [8]. These authors define a "Region of Effective Predictability" (REP$_A$) as that portion of the $X$ space in which the variance of the predicted $\hat{y}$ does not exceed twice the variance of $\hat{y}$ predicted at the centroid of the $X$ matrix. The volume of the region is then compared with that of a similarly defined REP, denoted REP$_0$. The latter refers to a "fictitious orthogonal reference design" of "orthogonal data with the same N and the same rms values as the actual data." The ratio of the volume of REP$_A$ to that of REP$_0$ is taken as "an overall measure of the loss of predictability volume due to collinearity".

This concept, apart from its artificial character, suffers from other shortcomings. Like so many other treatments, it attempts to provide a *measure of collinearity*. But the practitioner who is confronted with a collinear $X$ matrix does not need a measure of collinearity: he needs a way to use the data for the purpose for which they were obtained. Furthermore, this measure loses its meaning when expansion variables are present. For example, for the artificial quadratic set of table 12, Willan and Watts' measure would indicate a high degree of collinearity which, while literally true, is totally misleading since the collinearity in no way reduces the usefulness and predicting power of the regression equation, as long as the meaning of the expansion term is taken into account. But even in cases without expansion terms, the measure in question may be misleading. Thus when applied to the protein calibration data of table 5, it may well lead the analyst to give up on these data as a hopelessly highly-collinear set, whereas, as we have seen, there is nothing wrong with this set and it can indeed be used very effectively for the calibration of a method for protein determination based on reflectance measurements.

Finally, a few words about estimating the $\beta$-coefficients considered as rates of change of $y$ with changes in the individual $x_j$. As pointed out by Box [9], this is generally *not* a desirable use of regression equations. If, however, it

is the major purpose of a particular experiment, then this experiment should be designed accordingly, which means: essentially with an orthogonal $X$ matrix. A collinear $X$ matrix leads to the ability to estimate certain linear combinations of the $\beta$'s much better than the $\beta$'s themselves. The experimenter can calculate the VF values, not only for any point of $X$ space, but also for any $\beta$ or combination of $\beta$'s, *and he can do this without making a single measurement*, i.e., in the planning stages of the experiment. If the experimenter does not take advantage of this opportunity, he may be in for considerable disappointment, after having spent time, money, and effort on inadequate experimentation. We believe that he advocacy of remedial techniques, such as Ridge Regression for collinear data is unwise. One of the most important tasks of a data analyst is to detect, and to call attention to, limitations in the use and interpretation of the data.

# References

[1] Buck, J.B., Studies on the Firefly, Part I: The Effects of Light and Other Aspects on Flashing in Photinus Pyralic, with Special Reference to Periodicity and Diurnal Rhythm, Physiological Zoology, 10, 45-58 (1937).

[2] Mandel, J., Use of the Singular Value Decomposition in Regression Analysis, The American Statistician, 36, 15-24 (1982).

[3] Fearn, T., A Misuse of Ridge Regression in the Calibration of a Near Infrared Reflectance Instrument, Applied Statistics, 32, 73-79 (1983).

[4] Belsley, D.A., E. Kuh, and R.E. Welsh, Regression Diagnostics: Identifying Influential Observations and Sources of Collinearity, Wiley, NY (1980).

[5] Draper, N.R. and H. Smith, Applied Regression Analysis, Wiley, 2nd Edition, NY (1981).

[6] Belsley, D.A., Demeaning Conditioning Diagnostics Through Centering, The American Statistician, 38, 73-77, and "Comments", 78-93 (1984).

[7] Daniel, C., and F.S. Wood, Fitting Equations to Data, Wiley, 2nd Edition, NY (1980).

[8] Willan, A.R. and D.G. Watts, Meaningful Multicollinearity Measures, Technometrics, 20, 407-12 (1978).

[9] Box, G.E.P., Use and Abuse of Regression, Technometrics 8 625-29 (1966).

# DISCUSSION

of the John Mandel paper, The Regression
Analysis of Collinear Data

## R. W. Gerlach

Monsanto Agricultural Products Co.

I fully agree with Mandel in that one's model can (usually) only be assigned a degree of validity in the region spanned by the data used to generate the model. Though the range for all variables may be quite large, collinearity effectively restricts the model to a particular subregion. One should be aware of these restrictions so as not to misapply the model to those regions not represented by the data. I want to point out the need to carry out an additional initial operation; one should always examine the dataset for outliers. Otherwise the suggestion for using the largest and smallest values on each principal coordinate to examine the constraints may lead to overstating the region for a valid calibration. In fact, this could happen anyway if the shape formed by the data vectors was peculiar, perhaps occupying two disjoint regions for instance.

Additional constraints are frequently available to the analytical chemist. Minimum and maximum values along the original variables as well as conditions upon functions of these variables are frequently encountered. The location of the effective predictive domain within this potentially allowed domain could also be useful. A comparison might lead the researcher to conclude that more effort should be spent gathering additional data so that the calibration equation was valid over the desired region.

The variance factor (VF), resulting from the propagation of errors through the transformations, is a good method for observing how well characterized the model is at any location. Though principal components regression has been in the literature of analytical chemistry for some time [1,2][1] a paper dealing with the region of applicability for the model has only recently been published [3]. In this case the authors used as their criteria the expected mean square error. Hopefully, the propagation of error in this and related techniques will become more commonplace in analytical chemistry.

I think that the comparison of the measures advocated in this paper to the condition number is somewhat misdirected. The condition number can be used to provide a measure of how sensitive a model could be to variations in the data matrix. However, it would certainly not be appropriate to consider a condition number for the complete data matrix if one is dealing with only a subset of its dimensions in the principal component regression. The condition number

assists one in interpreting the sensitivity of the model given all the original variables (or any orthogonal transformation). The condition number for the rotated coordinate system of the principal coordinates will be the same as for the original coordinate system. In the original coordinate system a large condition number signaled that the regression coefficients were not all well known. In the rotated eigenvector coordinate system this same condition number reflects the fact that coefficients for the eigenvectors with small eigenvalues will not be estimated accurately. However, since only the eigenvectors with significant eigenvalues will be considered in the principal component regression, the condition number for the entire matrix is not an appropriate parameter to consider. In fact, the only thing we can say is that one expects large condition numbers every time a principal component regression is the method of choice.

It should also be pointed out that other aspects of collinearity are frequently encountered by analytical chemists. While this paper deals with collinearity as it affects the region for applicability of the model in terms of predictability, it doesn't address questions as to the reliability of the model coefficients. Also, instead of generating a calibration or predictive equation, one might wish to evaluate possible models in which the independent factors behave somewhat similarly. What limitations are placed on the results of the traditional regression analysis? I want to mention that statisticians have already developed several appropriate techniques [4], such as methods to estimate confidence regions and the effective sample size. Hopefully, these and other measures to test the validity of the proposed model will be more widely used.

The propagation of errors through a constrained correlated regression would also be an appropriate technique for investigating the significance of the terms in a proposed model. As mentioned above, often there are known constraints, yet this information is commonly overlooked. A recent comparison of multivariate techniques applied to source apportionment of aerosols in which collinearity was an important factor showed that the known constraints were mostly ignored [5]. Mathematical techniques which deal with these extra conditions [6], though more complex numerically, should be investigated for their potential benefits to areas of analytical chemistry and brought into more common use.

---

[1] Figures in brackets indicate literature references.

# References

[1] Bos, M., and G. Jasink, The Learning Machine in Quantitative Chemical Analysis, *Analytica Chimica Acta* **103**, pp 151–165 (1978).

[2] Martens, H., Factor Analysis of Chemical Mixtures, *Analytica Chimica Acta* **112**, pp 423–442 (1979).

[3] Fredericks, P. M.; Lee, J. B.; Osborn, P. R., and D. A. J. Swinkels, Materials Characterization Using Factor Analysis of FT-IR Spectra. Part 2: Mathematical and Statistical Considerations, *Applied Spectroscopy* **39**, pp 311–316 (1985).

[4] Willan, A. R., and D. G. Watts, Meaningful Multicollinearity Measures, *Technometrics* **20**, pp. 407–412 (1978).

[5] Currie, L. A.; Gerlach, R. W., and C. W. Lewis, *et al*, Interlaboratory Comparison of Source Apportionment Procedures: Results for Simulated Data Sets, *Atmospheric Environment* **18**, pp 1517–1537 (1984).

[6] Rust, B. W., and W. R. Burrus, Mathematical Programming and the Numerical Solution of Linear Equations, Elsevier (1972).

478

# Optimization

## Stanley N. Deming

### University of Houston–University Park, Houston, TX 77004

Most research and development projects require the optimization of a system response as a function of several experimental factors. Familiar chemical examples are the maximization of product yield as a function of reaction time and temperature; the maximization of analytical sensitivity of a wet chemical method as a function of reactant concentration, pH, and detector wavelength; and the minimization of undesirable impurities in a pharamaceutical preparation as a function of numerous process variables. The "classical" approach to research and development involves answering the following three questions in sequence:

   1) What are the important factors? (Screening)

   2) In what way do these important factors affect the system? (Modeling)

   3) What are the optimum levels of the important factors?

As R. M. Driver has pointed out, when the goal of research and development is optimization, an alternative strategy is often more efficient:

   1) What is the optimum combination of *all* factor levels? (Optimization)

   2) In what way do these factors affect the system? (Modeling *in the region of the optimum*)

   3) What are the important factors?

The key to this alternative approach is the use of an efficient experimental design strategy that can optimize a relatively large number of factors in a small number of experiments. For many chemical systems involving continuously variable factors, the sequential simplex method has been found to be a highly efficient experimental design strategy that gives improved response after only a few experiments. It does not involve detailed mathematical or statistical analysis of experimental results. Sequential simplex optimization is an alternative evolutionary operation (EVOP) technique that is not based on traditional factorial designs. It can be used to optimize several factors (not just one or two) in a single study. Some research and development projects exhibit multiple optima. A familiar analytical chemical example is column chromatography which often possesses several sets of locally optimal conditions. EVOP strategies such as the sequential simplex method will operate well in the region of one of these local optima, but they are generally incapable of finding the global or overall optimum. In such situations, the "classical" approach can be used to estimate the general region of the global optimum, after which EVOP methods can be used to "fine tune" the system. For example, in chromatography the Laub and Purnell "window diagram" technique can often be applied to discover the general region of the global optimum, after which the sequential simplex method can be used to "fine tune" the system, if necessary. The theory of these techniques and applications to real situations will be discussed.

Key words: optimization; screening; simplex.

## 1. Introduction

Most research and development projects require the *optimization* of a system response (dependent variable) as a function of several experimental factors (independent variables). Familiar chemical examples are:

**About the Author:** Stanley N. Deming is with the Department of Chemistry at the University of Houston–University Park.

1) *re-establishing* acceptable product yield as a function of reaction time and reaction temperature after a design change in a chemical process;

2) *maximizing* the analytical sensitivity of a wet chemical method as a function of reactant concentration, pH, and detector wavelength;

3) *tuning-up* a nuclear magnetic resonance spectrometer by adjusting eleven highly interacting shim coil controls to produce optimum peak shape.

4) *finding* a combination of values for eluent variables that

will give adequate separation in high performance liquid chromatography.

Although "optimization" is often taken literally to mean making something "as perfect, effective, or functional as possible" [1][1], in chemical practice it usually means making something "acceptable" or "adequate," as in examples one and four above. Optimization in chemistry usually involves adjusting a system until it is brought to some desired threshold of performance.

The dual purposes of this short paper are to discuss several strategies for the optimization of chemical systems and to discuss strengths, weaknesses, and appropriate settings for each approach. The intent of these comments is not to suggest rigid guidelines for the proper uses of optimization methods, but rather to stimulate discussion directed toward a better understanding of how these methods can be used in practice.

## 2. Classical Experimental Designs

The "classical" approach to optimization in research and development involves answering the following three questions in sequence:

1) What are the important factors? (SCREENING)
2) In what way do these important factors affect the system? (MODELING)
3) What are the optimum levels of the important factors? (OPTIMIZATION)

Classical experimental designs (e.g., fractional factorial designs and central-composite designs [2,3]) can be used to screen factors and to acquire data for modeling the system as a function of the most important variables. The resulting model can then be used to predict the treatment combination (experimental conditions) giving the optimum response [4-6]. The statistical literature is rich in examples showing how statistically designed experiments have been used in this way to solve significant chemical problems (e.g., [7]).

### A. Modeling

The critical part of the classical approach is the second step, modeling. At the very least, a model that fits reasonably well over a limited region of the factor space can be used to predict a direction to move to obtain improved response (as in evolution operation, or EVOP [8]). A model that fits well over a larger region of factor space is, of course, even more useful. However, if the (usually empirical) model contains more than a few factors, then the number of experiments required to fit the model will be impractically large. For example, if a full second-order polynomial model containing $k$ factors is used to model the system, the number of model parameters will be equal to $(k+1)(k+2)/$

[1]Figures in brackets indicate literature references.

2; for five, six, seven, and eight factors the numbers of model parameters are 21, 28, 36, and 45, respectively. At least this many experiments must be carried out to provide data for the estimation of the parameter values; typically, central composite designs are used which require $2^k + 3k + 1$ experiments (plus three replicates to estimate "pure error") for a total of 46, 80, 146, and 276 experiments for five, six, seven, and eight factors, respectively.

Thus, a desire to avoid extraordinarly large numbers of experiments becomes a strong driving force for limiting (typically to only thee or four) the number of factors to be investigated by classical experimental designs. Hence, the need for the initial screening of factors to choose only the most important ones.

### B. Screening

There are problems with screening experiments. For example, most screening experiments are based on first-order models which assume no interactions. If interactions do exist, then factors which truly have a significant affect on the system might not appear to be statistically significant and would be discarded by the screening process.

A second problem with screening experiments can occur if the effect of a factor depends upon its own level "self interaction"). If screening experiments are carried out in a region where the response is "flat" with respect to the factor of interest (a stationary region), that factor will not appear to a be very significant when, in fact, at different levels of that factor, the effect on response might be considerable.

As a final example of difficulties in screening for significant factors, the "wrong" statistical test is usually used when screening factors for their significance. It is true that if a factor is "significant at the 95% level of confidence," then it is probably an important factor and should be retained for further investigation. However, if a factor is "not significant at the 95% level of confidence," it does not mean that it is an unimportant factor and can be neglected. It might, for example, be significant at the 94.73% level of confidence, not enough to exceed the common threshold of 95% confidence, but still highly significant nonetheless. Ideally, the question that should be asked while carrying out screening experiments is not "which factors are significant at some high level of confidence," but rather, "which factors are insignificant at some equally high level of confidence." Unfortunately, the type of experimentation required to answer this second question is extensive and expensive. An alternative approach is to increase the risk (alpha) of stating that a factor is significant when in fact it is not, so that fewer truly significant factors are rejected [9].

### C. Comments on the use of classical experimental design for optimization

Classical experimental designs appear to have been successful in the past for "optimizing" many existing chemical systems largely because these systems are usually run not at the true optimum but rather are operated at some "threshold

480

of acceptability." A response surface view of this would be that the system is being run at a point on the *side* of a hill; not at the *top* of the hill, but far enough up on the side that the system gives acceptable performance. As long as the response surface maintains its shape and position, and as long as the factor levels are kept in statistical control, the system will perform acceptably.

However, if the response surface changes its shape or "moves" slightly (as a result, for example, of scale buildup in heat exchangers, or different suppliers of feed stocks), then the previous set of factor levels might no longer produce adequate performance from the system: the same setpoint will now correspond to some slightly lower position on a changed response surface. In situations like this, small screening experiments (such as saturated fractional factorial designs [2] or Plackett-Burmann designs [10]) are not too much affected by factor interactions and are likely to give nearly true estimates of the first-order factor effects (e.g., the effect of temperature, or the effect of increased amounts of feed stocks). Similarly, a first-order model offers a good approximation to the true shape of the response surface over a limited region. Thus, the application of screening experiments to choose the most significant factors is usually successful in this application. When these most significant factors are used in a model of the system that is first-order with *interactions* (fitted, say, to the results of a two-level factorial design), then the fitted model will usually suggest an appropriate direction to move. Changing the factor levels in this direction will usually move "up" the hill to a point lying above the threshold of performance and once again achieve adequate ("optimum") response from the system.

## 3. Sequential Simplex Optimization

As R.M. Driver has pointed out [11], when the goal of research and development is optimization, an alternative strategy is often more efficient. This alternative strategy asks essentially the same questions as the classical approach to optimization, but it asks the questions in reverse order:

1) What is the optimum combination of *all* factor levels? (OPTIMIZATION)
2) In what way do these factors affect the system? (MODELING *in the region of the optimum*)
3) What are the important factors? (SCREENING for effects *in the region of the optimum*)

The key to this alternative approach is the use of an efficient experimental design strategy that can optimize a relatively large number of factors in a small number of experiments. Once in the region of the optimum, classical experimental designs can be used to full advantage to model the system and determine factor importance in a limited region of the total factor space.

### A. Ignoring Initial Screening Experiments and Avoiding Models

For many chemical systems involving continuously variable factors and relative short times for each experiment, the sequential simplex method [12-38] has been found to be a highly efficient experimental design strategy that gives improved response after only a few experiments. It is a logically-driven algorithm that does not involve detailed mathematical or statistical analysis of experimental results. Sequential simplex optimization is an alternative evolutionary operation (EVOP) technique that is not based on traditional factorial designs.

There are two reasons for the efficiency of the sequential simplex method. The first reason is the number of experiments required in the experimental design itself. A simplex is a geometric figure containing a number of vertexes equal to one more than the number of dimensions of the factor space. Each vertex locates a treatment combination in factor space. Thus, the number of experiments required for a simplex is $k+1$ where, again, $k$ is the number of factors. Thus, a five, six, seven, or eight factor system would require only 6, 7, 8, or 9 experiments to define a simplex.

The second reason for the efficiency of the sequential simplex method is that it takes only one or two additional experiments to move the experimental design into an adjacent region of factor space. This is independent of the number of factors involved. When classical experimental designs are used in this type of "evolutionary operation" mode, a larger number of experiments (at least half of the factor combinations in the pattern) is usually required to move the experimental design into an adjacent region of factor space.

In our experience with the simplex, systems of up to 11 factors can be brought into the region of the optimum in only 15 or 20 experiments after the initial simplex has been constructed.

### B. Limitations

The simplex does have its limitations, however. The system must be in "statistical control" if the simplex is to be used—that is, the system should have only a small amount of purely experimental uncertainty ("pure error"). It is recommended that after the initial simplex has been evaluated and before the first simplex move is begun, the vertex giving the worst response and the vertex giving the best response be repeated two more times each to evaluate the reproducibility of the system. If the reproducibility is good, then the simplex will progress well; if the reproducibility is poor, then the simplex will tend to wander. In this latter case, steps should be taken to improve the purely experimental uncertainty of the system; if this is not possible, then classical experimental designs offer advantages because of their noise-reducing capabilities [39].

The system should not drift with time. However, changes with time can often be detected and corrected for by running periodic experiments at a standard treatment combination.

The time of any one experiment must be relatively short. It has been suggested that the reason factorial experiments were developed before the sequential simplex was because of the experimental environment, specifically the improvement of agricultural crop yields. In this context, factorial experiments offer a great advantage in that several experiments can be carried out simultaneously and many results can be obtained after only one growing season. If the sequential simplex were to be used to improve agricultural production, only one move could be carried out each year and it could take several decades to optimize production.

Finally, the simplex is most powerful for continuous ("quantitative") variables. It can be used for discrete variables where there are several levels—perhaps at least five or six—and the levels can be logically ranked. It can not be used for unranked discrete ("qualitative") variables.

## 4. Systems Possessing Multiple Optima

Some research and development projects exhibit multiple optima. A familiar analytical chemical example is column chromatography which often possesses several sets of locally optimal conditions [40]. The reason for the existence of multiple optima is related to the phenomenon of changes in the order of elution with changing chromatographic conditions. EVOP strategies such as the sequential simplex method will operate well in the region of one of these local optima, but they are generally incapable of finding the global or overall optimum [23]. In such situations, classical factorial-type experiments can be used to fit models which in turn can be used to estimate the general region of the global optimum, after which EVOP methods can be used to "fine tune" the system. For example, in chromatography the Laub and Purnell "window diagram" technique [40] can often be applied to discover the general region of the global optimum, after which the sequential simplex method can be used to "fine tune" the system, if necessary [41-49].

## References

[1] Webster's New Collegiate Dictionary, G. & C. Merriam Company, Springfield, MA (1977), p. 806.

[2] Box, G.E.P.; Hunter, W.G. and J.S. Hunter, Statistics for Experimenters: An Introduction to Design, Data Analysis, and Model Building, Wiley, New York, NY (1978).

[3] Deming, S.N. and S.L. Morgan, The Use of Linear Models and Matrix Least Squares in Clinical Chemistry, Clin. Chem, 25, 840 (1979).

[4] Box, G.E.P., and K.B. Wilson, On the Experimental Attainment of Optimum Conditions, J. Roy. Stat. Soc., Ser. B, 13, 1 (1951).

[5] Box, G.E.P., The Exploration and Exploitation of Response Surfaces: Some General Considerations and Examples, Biometrics, 10, 16 (1954).

[6] Box, G.E.P., and P.V. Youle, The Exploration and Exploitation of Response Surfaces: An Example of the Link Between the Fitted Surface and the Basic Mechanism of the System, Biometrics, 11, 287 (1955).

[7] Davies, O.L., Ed., Design and Analysis of Industrial Experiments, 2nd ed., Hafner, New York, NY (1971).

[8] Box, G.E.P., and N.R. Draper, Evolution Operation: A Statistical Method for Process Improvement, Wiley, New York, NY (1969).

[9] Wilson, E.B., Jr., An Introduction to Scientific Research, McGraw-Hill, New York, NY (1952), p. 59.

[10] Plackett, R.L., and J.P. Burman, "The Design of Optimum Multifactorial Experiments," Biometrika, 33, 305 (1946).

[11] R.M. Driver, Chem. Brit., 6, 154 (1970).

[12] Spendley, W.; Hext, G.R., and F. R. Himsworth, Sequential Application of Simplex Designs in Optimization and Evolutionary Operation, Technometrics, 4, 441 (1962).

[13] Nelder, J.A., and R. Mead, A Simplex Method for Function Minimization, Computer J., 7, 308 (1965).

[14] Deming, S.N., and S.L. Morgan, Simplex Optimization of Variables in Analytical Chemistry, Anal. Chem., 45, 278A (1973).

[15] Deming, S.N. and P.G. King, Computers and Experimental Optimization, Research/Development, 25(5), 22 (1974).

[16] King, P.G., Automated Development of Analytical Methods, Ph.D. Dissertation, Emory University, Atlanta, GA (1974).

[17] Morgan, S.L., and S.N. Deming, Simplex Optimization of Analytical Methods, Anal. Chem., 46, 1170 (1974).

[18] King, P.G., and S. N. Deming, UNIPLEX: Single-Factor Optimization of Response in the Presence of Error, Anal. Chem., 46, 1476 (1974).

[19] Yarbro, L.A., and S.N. Deming, Selection and Preprocessing of Factors for Simplex Optimization, Anal. Chim. Acta, 73, 391 (1974).

[20] King, P.G.; Deming, S.N., and S.L. Morgan, Difficulties in the Application of Simplex Optimization to Analytical Chemistry, Anal. Lett., 8, 369 (1975).

[21] Parker, L.R., Jr., Morgan, S.L., and S.N. Deming, Simplex Optimization of Experimental Factors in Atomic Absorption Spectrometry, App. Spectrosc., 29, 429 (1975).

[22] Dean, W.K.; Heald, K.J., and S.N. Deming, Simplex Optimization of Reaction Yields, Science, 189, 805 (1975).

[23] Morgan, S.L., and S.N. Deming, Optimization Strategies for the Development of Gas-Liquid Chromatographic Methods, J. Chromatogr., 112, 267 (1975).

[24] Olansky, A.S., and S.N. Deming, Optimization and Interpretation of Absorbance Response in the Determination of Formaldehyde with Chromotropic Acid, Anal. Chim. Acta, 83, 241 (1976).

[25] Deming, S.N., Morgan, S.L., and M.R. Willcott, Sequential Simplex Optimization, Amer. Lab., 8(10), 13 (1976).

[26] Turoff, M.L.H. and S.N. Deming, Optimization of the Extraction of Iron (II) from Water into Cyclohexane with Hexafluoroacetylacetone and Tri-n-Butyl Phosphate," Talanta, 24, 567 (1977).

[27] Deming, S.N. and S.L. Morgan, Advances in the Application of Optimization Methodology in Chemistry, Chapter 1 in B.R. Kowalski, Ed., Chemometrics: Theory and Application, ACS Symposium Series 52, American Chemical Society, 1977, p. 1.

[28] Deming, S.N., Optimization of Experimental Parameters in Chemical Analysis, Chapter 5 in J.R. DeVoe, Ed., Validation of the Measurement Process, ACS Symposium Series 63, American Chemical Society, 1977, p. 162.

[29] Olansky, A.S., Parker, L.R., Jr., Morgan, S.L., and S.N. Deming, Automated Development of Analytical Chemical Methods. The Determination of Serum Calcium by the Cresolphthalein Complexone Method, Anal. Chim. Acta, 95, 107 (1977).

[30] Deming, S.N. and L.R. Parker, Jr., A Review of Simplex Optimization in Analytical Chemistry, CRC Crit. Rev. Anal. Chem., 7, 187 (1978).

[31] Deming, S.N., Optimization of Methods, Chapter 2 in R.F. Hirsch, Ed., Proceedings of the Eastern Analytical Symposium on Principles of Experimentation and Data Analysis, Franklin Institute Press, 1978, p. 31.

[32] Olansky, A.S. and S.N. Deming, Automated Development of a Kinetic Method for the Continuous-Flow Determination of Creatinine, *Clin. Chem.*, **24**, 2115 (1978).

[33] Shavers, C.L., Parsons, M.L., and S.N. Deming, Simplex Optimization of Chemical Systems, *J. Chem. Educ.*, **56**, 307 (1979).

[34] Deming, S.N., The Role of Optimization Strategies in the Development of Analytical Chemical Methods, *American Laboratory*, **13**(6), 42 (1981).

[35] Deming, S.N. and S.L. Morgan, Teaching the Fundamentals of Experimental Design, *Anal. Chim. Acta.*, **150**, 183 (1983).

[36] Nickel, J.H. and S.N. Deming, Use of the Sequential Simplex Optimization Algorithm in Automated Liquid Chromatographic Methods Development, LC, 414 (1983).

[37] Golden, P.J. and S.N. Deming, Sequential Simplex Optimization with Laboratory Microcomputers, *Laboratory Microcomputers*, **3**(2), 44 (1984).

[38] Walters, F.H. and S.N. Deming, A Two-Factor Simplex Optimization of a Programmed Temperature Gas Chromatographic Separation, *Anal. Lett.*, **17**, 2197 (1984).

[39] Mendenhall, W. *Introduction to Linear Models and the Design and Analysis of Experiments*, Duxbury, Belmont, CA (1968).

[40] Laub, R.J. and J.H. Purnell, *J. Chromatogr*, **112**, 71 (1975).

[41] Morgan, S.L. and S.N. Deming, Experimental Optimization of Chromatographic Systems, *Sep. Purif. Methods*, **5**, 330 (1976).

[42] Deming, S.N. and M.L.H. Turoff, Optimization of Reverse-Phase Liquid Chromatographic Separation of Weak Organic Acids, *Anal. Chem.*, **50**, 546 (1978).

[43] Price, W.P., Jr.; Edens, R., Hendrix, D.L., and S.N. Deming, Optimized Reverse-Phase High-Performance Liquid Chromatographic Separation of Cinnamic Acids and Related Compounds, *Anal. Biochem.*, **93**, 233 (1979).

44] Price, W.P., Jr., and S.N. Deming, Optimized Separation of Scopoletin and Umbelliferone and *cis-trans* Isomers of Ferulic and *p*-Coumaric Acids by Reverse-Phase High-Performance Liquid Chromatography, *Anal. Chim. Acta*, **108**, 227 (1979).

[45] Kong, R.C., Sachok, B., and S.N. Deming, Combined Effects of pH and Surface-Active-Ion Concentration in Reversed-Phase Liquid Chromatography, *J. Chromatogr.*, **199**, 307 (1980).

[46] Sachok, B., Kong, R.C., and S.N. Deming, Multifactor Optimization of Reversed-Phase Liquid Chromatographic Separations, *J. Chromatogr.*, **199**, 317 (1980).

[47] Sachok, B., Stranahan, J.J., and S.N. Deming, Two-Factor Minimum Alpha plots for the Liquid Chromatographic Separation of 2,6-Disubstituted Anilines, *Anal. Chem.*, **53**, 70 (1981).

[48] Nickel, J.H., and S.N. Deming, Use of Window Diagram Techniques in Automated LC Methods Development, *Amer. Lab.*, **16**(4), 69 (1984).

[49] Deming, S.N., Bower, J.G., and K.D. Bower, Multifactor Optimization of HPLC Conditions, Chapter 2 in J.C. Giddings, Ed., *Advances in Chromatography*, **24**, 35 (1984).

# DISCUSSION

of the Stanley N. Deming paper, Optimization

## C. K. Bayne

Computing and Telecommunications Division,
Oak Ridge National Laboratory.

I appreciate the opportunity to make comments on Dr. Deming's paper. I will confine my comments to three areas: 1) optmization applications; 2) strategies for screening experiments; and 3) the steepest ascent method.

## 1. Optimization Applications

In 1971, Rubin, Mitchell, and Goldstein [1][1] surveyed the previous 25 years of major English language journals of analytical chemistry under the index heading of "statistics." This survey uncovered few papers in which experiments were statistically designed. Similar results were found by Morgan and Deming in 1974 [2] in their literature search under the heading "Optim/" in *Chemical Abstract* and

*Chemical Titles* covering eight previous years. Nine years later, Deming and Morgan [3] found 189 titles for the years 1962-1982 listed in *Chemical Abstracts* and *Science Citation Index* related to sequential simplex optimization. About 156 papers in this search are direct applications to chemical problems. In a recent survey by Rubin and Bayne [4] for the years 1974-1984, 65 applications of optimizations and response surface methods were found to be related just to analytical chemistry. These recent literature surveys indicate that statistically designed experiments are becoming an important part of chemical experiments.

Dr. Deming deserves a large share of credit for this increased use of statistically designed experiments in chemistry. He has promoted experimental design by his many publications, seminars, and lectures. The fact that he is a chemist who has championed the statistical cause is to be admired.

---

[1]Figures in brackets indicate literature references.

## 2. Screening Experiments

The strategy for using screening experiments is influenced by the cost of an experimental run. Here cost can be interpreted as price of material, time required for an experimental run, etc. For a high cost experiment, the strategy is to run a screening experiment to identify the important factors followed by an optimization experiment using these factors. An additional screening experiment is sometimes performed to confirm that the important factors have been properly identified. Dr. Deming is correct in stating that the error rate for rejecting significant factors should have more emphasis than the error rate for accepting non-significant factors in the initial screening experiment. This conservative approach can be accomplished by testing at a 0.20 or 0.25 significance level rather than the usual 0.05 significance level.

The low cost experimental run situations may require a different strategy. Dr. Deming advocates that first an optimization using all factors be performed using a sequential simplex method followed by a screening experiment to identify important factors at optimum conditions. For additional factors, he points out that only a small increase is required for the number of initial experiments which require $K+1$ experiments for $K$ factors, and no increase in the number of experiments is needed to move into an adjacent region of the factor space. However, the number of experiments for convergence may increase rapidly with an increase in the number factors. Nelder and Mead [6] reported that the mean number of experiments needed for convergence increases as a second-order function of the number of factors. Even for low cost experimental runs, the total number of experiments required for optimization may be impractical.

When dealing with a large number of factors, two adjustments to the sequential methods are suggested. First, use small screening experiments for initial experiments; and secondly, discard more then one vertex when moving into a different factor region.

Using small screening designs for initial experiments is suggested because execution of the initial regular simplex can be tedious. For example, the step size fraction required for one of the vertices in a four-factor regular simplex is (0.500, 0.289, 0.204, 0.791) [6]. In practice, running an experiment at the simplex vertices may either be difficult or impossible. By allowing the factor levels to be either a low or high level, initial screening designs can be easily run. Screening experiments for initial designs are given in table 1.

Changing only one vertex per experimental move may be too slow when there are many factors. For these cases, more than one vertex can be discarded to increase the rate of convergence. The rule for simplex moves is modified to delete more than one vertex by:

Table 1. Screening Experiments for Initial Designs.

| Factors | Runs | Initial Design |
|---|---|---|
| 2 or 3 | 3 or 4 | Simplex |
| $4 \leq K \leq 7$ | 8 | Fractional Factorial |
| $8 \leq K \leq 11$ | 12 | Plackett-Burman |
| $12 \leq K \leq 15$ | 16 | Fractional Factorial |
| $K > 16$ | | Fractional Factorial or Plackett-Burman |

$2 \times$(average of best vertices)$-$worst *vertices*.

To decide which are the best vertices and the worst vertices, first rank the response from lowest to highest value. Next, divide the responses into two groups with the lowest values in one group and the highest values in the second group. This division can be done $K$ ways for $K$-factors. For each division, calculate the average response for each group and then take their difference. The division that has the largest difference will indicate the best and the worst vertices. By this method, the difference between the average responses for the best and the worst vertices is maximized.

## 3. Steepest Ascent

In a literature search by Rubin and Bayne [4], few applications were found of the steepest ascent method advocated by Box and Wilson [7] for optimization. This method maximizes the gain in the reponse and performs better than fixed step size simplex method [8]. The main arguments against the steepest ascent method are: the initial experiments require too many experiments, and the calculations are too complex.

Although the steepest ascent method does require more initial experiments, the total number of experiments may not be as many as those for the simplex method. Worksheets similar to those used for the simplex methods can be used to alleviate calculation difficulties. Because steepest ascent has the potential of out performing the simplex method, its use should be encouraged.

## References

[1] Rubin, I.B., T.J. Mitchell, and G. Goldstein, A Program of Statistical Designs for Optimizing Specific Transfer Ribonucleic Acid Assay Conditions, Analytical Chemistry, 43 717-721 (1971).

[2] Morgan, S.L., and S.N. Deming, Simplex Optimization of Analytical Chemical Methods, Analytical Chemistry, 46 1170-1181 (1974).

[3] Deming, S.N., and S.L. Morgan, Teaching the Fundamentals of Experimental Design, Analytica Chemica Acta, 150 183-198 (1983).

[4] Rubin, I.B., and C.K. Bayne, unpublished communication (1984).

[5] Nelder, J.A., and R. Mead, A simplex Method for Function Minimiza-

tion, Computer Journal, **7** 308-313 (1965).

[6] Long, D.E., Simplex Optimization of the Response from Chemical Systems, *Analytical Chemical. Acta.* **46** 195-206 (1969).

[7] Box, G.E.P., and K.B. Wilson, On the Experimental Attainment of Optimum Conditions, *Journal of the Royal Statistical Society,* Series B, **13** 1-45 (1951).

[8] Spendley, W., G.R. Hext and F.R. Himsworth, *Sequential Application of Simplex Designs in Optimization and Evolutionary Operation, Technometrics 4 441-451 (1962).*

# Strategies for the Reduction and Interpretation of Multicomponent Spectral Data

## Isiah M. Warner, S. L. Neal, and T. M. Rossi
### Emory University, Atlanta, GA 30322

Fluorescence data can be rapidly acquired in the form of an emission-excitation matrix (EEM) using a novel fluorometer called a video fluorometer (VF). An EEM array of 4096 data points composed of fluorescence intensity measured at 64 different emission wavelengths and excited at 64 different excitation wavelengths can be acquired in less than one second. The time-limiting factor in using this information for analytical measurement is the interpretation step. Consequently, sophisticated computer algorithms must be developed to aid in interpretation of such large data sets. For "$r$" number of components, the EEM data matrix, $\mathbf{M}$, can be conveniently represented as

$$\mathbf{M} = \sum_{i=1}^{r} \alpha_i \mathbf{x}(i)\mathbf{y}(i)^t$$

where $\mathbf{x}(i)$ and $\mathbf{y}(i)^t$ are the observed excitation and emission spectra of the $i^{th}$ component and $\alpha_i$ is a concentration dependent parameter. Such a data matrix is readily interpreted using linear algebraic procedures.

Recently a new instrument has been described which rapidly acquires fluorescence detected circular dichroism (FDCD) data for chiral fluorophores as a function of multiple excitation and emission wavelengths. The FDCD matrix is similar in form to EEM data. However, since the FDCD matrix may have legitimate negative entries while the EEM is theoretically non-negative, different assumptions are required. This paper will describe the mathematical algorithms developed in this laboratory for the interpretation of the EEM in various forms. Particular emphasis will be placed on linear algebraic and two-dimensional Fourier Transform procedures.

Key words: circular dichroism; eigenvectors; fluorescence; pattern recognition.

## Introduction

Advances in computer technology and developments in multiparametric detection devices have had a profound effect on developments in chemical analysis. These developments have made it possible to expand the applicability of several analytical methods to more complex systems. Multicomponent analysis by fluorescence spectroscopy is one example of such a method.

In the conventional fluorescence experiment, the sample solution is irradiated with monochromatic light which produces molecules in the excited state when absorbed. As these molecules return to the ground state, light with a characteristic wavelength distribution is emitted. This distribution of the emitted light is known as the fluorescence emission spectrum. The fluorescence excitation spectrum is the fluorescence intensity as a function of absorption wavelength. At low absorbance ($<0.01$), the intensity of the fluorescence emitted at a given wavelength is directly proportional to the amount of light absorbed and therefore to the concentration of the analyte in the sample solution.

These characteristic spectra make fluorescence spectroscopy inherently more selective than absorption methods and provide qualitative as well as quantitative measurements. For example, when the excitation spectrum of the analyte only partially overlaps the excitation spectrum of the

other components in the sample, the sample can be irradiated with light of a wavelength which is only absorbed by the analyte. Since only the analyte absorbs the incident light, fluorescence will only be emitted by the analyte, in the absence of synergistic effects such as energy transfer. This technique is called selective excitation [1][^1] and can be used qualitatively by acquiring the complete emission spectrum of the analyte or quantitatively by determining the analyte concentration using the emission of a sample of known concentration as a standard.

This approach to multicomponent analysis can be expanded by acquiring fluorescence spectra at several excitation and emission wavelengths. Acquiring multiple spectra with conventional instrumentation is a time-consuming process, even with a microprocessor-controlled instrument. It can require more than one hour to acquire the complete emission spectrum of a sample at 64 excitation wavelengths using conventional instrumentation. Many samples would undergo significant photodecomposition over such a period of time. The development of the videofluorometer (VF) which rapidly acquires two dimensional fluorescence data has resolved this problem [2]. This instrument uses polychromatic excitation and a silicon intensified target vidicon detector (television camera) to acquire data in matrix format without mechanical scanning. The VF can acquire 64 emission spectra generated at 64 excitation wavelengths in less than 1 second. Therefore, data processing of this emission-excitation matrix (EEM) is now the time-limiting step in the analysis and requires the use of a computer for data reduction.

The rows of the EEM of a pure sample are multiples of the emission spectrum of that compound, while the columns are multiples of the excitation spectrum, indicating that the EEM is of bilinear form. Such matrices are ideally suited for data reduction techniques such as factor analysis [3] and pattern recognition [4].

Another technique whose applicability is expanded by multiparametric detection is fluorescence detected circular dichroism (FDCD) [5]. The FDCD spectrum is the difference in the intensity of flurnescence produced by excitation with left and right circularly polarized light recorded as a function of the wavelength of the exciting light. Chiral molecules preferentially absorb one form of circularly polarized light. The FDCD spectrum reflects the structure of the chiral fluorophore in a solution. A two-dimensional rapid scanning FDCD spectrometer has been developed which measures the FDCD at several excitation and emission wavelengths [6]. When these data are collected in matrix format, they are also of bilinear form. However, the algorithms designed to resolve the spectra of components from the EEM cannot be applied to this ellipticity matrix since it may contain legitimate negative values. This

[^1]Figures in brackets indicate literature references.

manuscript will discuss the strategies developed for the qualitative reduction of multicomponent EEMS as well as the alternative techniques developed for multicomponent ellipticity matrices.

## Eigenvector Analysis [7]

When the absorbance of a sample is less than 0.01, the intensity of the fluorescence, $I_f$, can be approximated by the expression

$$I_f = 2.303 \, I_0 \phi_f \epsilon bc \tag{1a}$$

where $I_0$ is the intensity of the incident radiation, $\phi_f$ is the fluorescence quantum efficiency (the fraction of absorbed photons emitted as fluorescence), $\epsilon$ is the molar extinction coefficient, $b$ is the thickness of sample cell and $c$ is the concentration of the fluorophore in the sample solution. Each element, $m_{ij}$, of the emission-excitation matrix, $M$, represents the fluorescence intensity at wavelength $\lambda_j$ that was generated by excitation at wavelength $\lambda_i$. Therefore, each of these elements can be generally expressed as

$$m_{ij} = 2.303 \phi_f I_0(\lambda_i) \epsilon(\lambda_i) \gamma(\lambda_j) \delta(\lambda_j) bc \tag{2a}$$

where $\gamma(\lambda_j)$ reflects the dependence of $I_f$ on the monitored emission wavelength and $\delta(\lambda_j)$ is a parameter which incorporates instrumental artifacts like sensitivity and signal collection geometry. Combining these terms based on excitation and emission wavelength related variables results in a simpler expression:

$$m_{ij} = \alpha x_i y_i \tag{3a}$$

where $\alpha$ is a scalar that equals $2.303 \phi_f bc$, $x_i$ is the excitation term given by

$$x_i = I_0(\lambda_i) \epsilon(\lambda_i) \tag{4a}$$

and $y_j$ is the emission term which is expressed

$$y_j = \gamma(\lambda_j) \delta(\lambda_j) \tag{5a}$$

When the $x_i$ are properly sequenced, the array of $x_i$ is a representation of the excitation spectrum and can be denoted x in vector notation. Likewise, when the $y_j$ are properly sequenced, the vector y represents the emission spectrum. Since the emission profile is independent of the exciting wavelength, and the excitation profile is independent of the monitored emission wavelength, the matrix M can clearly be expressed as the vector product of x and y, multiplied by a scalar concentration term, i.e.

$$M = \alpha x y^T \tag{6a}$$

488

In a sample containing "n" fluorescent compounds, the matrix $M$ is the sum of the EEMs of the individual components provided synergistic effects are negligible. Thus, the n component matrix can be expressed as

$$M = \sum_{k=1}^{n} \alpha_k x_k y_k^T .$$ (7a)

A more convenient notation for $M$ is

$$M = XY$$ (8a)

in which the columns of the matrix $X$ are the excitation spectra, $x_k$, of the $n$ components, and the rows of $Y$ are the emission spectra, $y_k$, of the components. The concentration term can be considered to be absorbed into either of the matrices.

Qualitative analysis of the matrix $M$ requires a determination of the number of independently emitting compounds, i.e., the rank of the matrix, and the set of basis vectors $x_k$ and $y_k$ which are the excitation and emission spectra of the components. Eigenanalysis plays a significant role in both determinations. Therefore, it is appropriate to preface the discussion of rank estimation and spectral resolution with a brief presentation of pertinent eigenanalysis principles.

An eigenvector is defined as any vector $x$ which is a solution of the equation

$$Ax = \lambda x$$ (9a)

in which $\lambda$ is a scalar called the eigenvalue. The magnitude of the eigenvalue is a consequence of the importance of the information reflected in the eigenvector to the data in the matrix. When the eigenvalue is large, then the factor represented by the eigenvector makes a large contribution to the data. If it is small, then the contribution of the factor is small. Therefore, when a matrix has a rank $n$ which is greater than 1, it has $n$ eigenvectors and $n$ eigenvalues.

Since the matrix $M$ is bilinear, the covariance matrices, $M^T M$ and $MM^T$ can be used to generate the eigenvectors. The eigenvalues of the covariance matrices are the squares of the eigenvalues of $M$, but apart from this small detail, it is more expedient to use the covariance matrices to generate the eigenvectors since the covariance matrices are always square and symmetric.

As the preceding discussion indicates, rank estimation of an ideal matrix (one that is noise-free) is straightforward, simply determined by the number of non-zero eigenvalues. However, experimental data matrices are not free of noise. They can contain systematic or random errors superimposed on the signal. For these matrices, there are several methods of rank estimation [3,8]. When the signal-to-noise ratio is high, it is possible to differentiate the eigenvalues associated with fluorescence (primary eigenvectors) from those associated with noise (secondary eigenvectors) by a direct comparison of their magnitudes. As the signal-to-noise ratio decreases, this approach becomes more difficult. Another method which will be described in a later section of this manuscript is differentiation of eigenvectors based on their frequency distributions, realizing that random noise in spectral data usually will have higher frequency than the signal.

The other part of the analysis of $M$ is the resolution of the basis vectors $x_k$ and $y_k$ from the matrix. An infinite number of basis vectors exists for a given matrix, and resolving the spectral vectors from the matrix is usually not possible without a priori knowledge of the components. However, the eigenvectors are an orthonormal basis for $M$ and are easily generated. Since the eigenvectors reproduce the matrix, $M$ can be expressed as

$$M = UV$$ (10a)

in which the columns of $U$ are the excitation eigenvectors, $u_k$, and the rows of $V$ are the emission eigenvectors, $v_k$. The eigenvalues have been absorbed into one of the matrices. Since the eigenvectors are an orthonormal basis, they often contain negative elements even though $M$ is theoretically a non-negative matrix. Emission and excitation spectra are also theoretically non-negative. Therefore, the eigenvectors can be transformed to possible spectral vectors by transforming them to a non-negative basis set.

To perform this transformation, the values of the matrix $K$ and its inverse, $K^{-1}$, which transform $U$ and $V$ to non-negative matrices, must be found. This is algebraically sound, since it is equivalent to multiplying $U$ and $V$ by the identity matrix. This transformation is mathematically expressed as

$$M = UKK^{-1}V .$$ (11a)

This condition also ensures that the transformed vectors are also a basis for $M$. The transformed excitation vectors are the columns of the matrix $U'$, and the transformed emission vectors are the rows of the matrix $V'$. These matrices are given by the equations

$$U' = UK \geq 0$$ (12a)

and

$$V' = K^{-1}V \geq 0 .$$ (13a)

The values of the elements of $K$ and $K^{-1}$ can be found from the expressions for the elements of $U'$ and $V'$. In the two component case, the elements are given by

489

$$u'_{1i} = k_{11}u_{1i} + k_{21}u_{2i} \geq 0 \ , \tag{14a}$$

$$u_{2i} = k_{22}u_{2i} + k_{12}u_{1i} \geq 0 \ , \tag{15a}$$

$$v'_{1j} = 1/|k|(k_{22}v_{1j} - k_{12}v_{2j}) \geq 0 \ , \tag{16a}$$

and

$$v'_{2j} = 1/|k|(k_{11}v_{2j} - k_{21}v_{1j}) \geq 0 \ . \tag{17a}$$

It can be assumed that $k_{11}=k_{22}=1$ without loss of generality. The values for $k_{12}$ and $k_{21}$ can be found by solving these expressions. It's clear from these expressions why this method is only applicable to two component matrices since the elements of $\mathbf{K}^{-1}$ become non-linear for more than two components. The boundaries for the values of $k_{12}$ and $k_{21}$ given by these expressions are

$$\min_{v_{2j}>0}\ \frac{v_{1j}}{v_{2j}} \geq k_{12} \geq \max_{u_{1j}>0}\ \frac{-u_{2i}}{u_{1i}} \tag{18a}$$

and

$$\min_{v_{1j}>0}\ \frac{v_{2j}}{v_{1j}} \geq k_{21} \geq \max_{u_{2i}>0}\ \frac{-u_{1i}}{u_{2i}} \tag{19a}$$

The accuracy of transformations performed with values meeting these criteria has been shown to be a function of the overlap of the spectra of the components. This is illustrated in the ambiguity table in figure 1 which summarizes the results of the transformation for the 16 possible spectral overlap combinations for a binary mixture. It should be noted that in 7 out of the 16 possible cases, at least one spectrum of each component is correctly resolved. Close inspection of the table shows that in cases where the spectra of neither component are enveloped in both dimensions, at least one of the spectra of each component is resolved unambiguously. Clearly, if both the emission and excitation spectra of the two compounds have only partial overlap, all four spectra will be resolved unambiguously from their mixture matrix.

## Emission Overlap

| comp. 1 / comp. 2 | | | |
|---|---|---|---|
| $\min_{v_{2i}>0}\frac{v_{1i}}{v_{2i}} = k_{12} = \max_{u_{1i}>0}\frac{-u_{2i}}{u_{1i}}$ <br> $\min_{v_{1i}>0}\frac{v_{2i}}{v_{1i}} = k_{21} = \max_{u_{2i}>0}\frac{-u_{1i}}{u_{2i}}$ <br> **All spectra certain.** | $\min_{v_{2i}>0}\frac{v_{1i}}{v_{2i}} = k_{12} = \max_{u_{1i}>0}\frac{-u_{2i}}{u_{1i}}$ <br> $\min_{v_{1i}>0}\frac{v_{2i}}{v_{1i}} \geq k_{21} > \max_{u_{2i}>0}\frac{-u_{1i}}{u_{2i}}$ <br> 1 spectrum of ea. comp. certain. 1 spectrum of ea. component given by extreme k. | $\min_{v_{2i}>0}\frac{v_{1i}}{v_{2i}} \geq k_{12} > \max_{u_{1i}>0}\frac{-u_{2i}}{u_{1i}}$ <br> $\min_{v_{1i}>0}\frac{v_{2i}}{v_{1i}} = k_{21} = \max_{u_{2i}>0}\frac{-u_{1i}}{u_{2i}}$ <br> 1 spectrum of ea. comp. certain. 1 spectrum of ea. component given by extreme k. | $\min_{v_{2i}>0}\frac{v_{1i}}{v_{2i}} \geq k_{12} > \max_{u_{1i}>0}\frac{-u_{2i}}{u_{1i}}$ <br> $\min_{v_{1i}>0}\frac{v_{2i}}{v_{1i}} \geq k_{21} > \max_{u_{2i}>0}\frac{-u_{1i}}{u_{2i}}$ <br> All spectra uncertain. Both spectra of ea. component given by extreme k's. |
| $\min_{v_{2i}>0}\frac{v_{1i}}{v_{2i}} > k_{12} \geq \max_{u_{1i}>0}\frac{-u_{2i}}{u_{1i}}$ <br> $\min_{v_{1i}>0}\frac{v_{2i}}{v_{1i}} = k_{21} = \max_{u_{2i}>0}\frac{-u_{1i}}{u_{2i}}$ <br> 1 spectrum of ea. comp. certain 1 spectrum of ea. component given by extreme k. | $\min_{v_{2i}>0}\frac{v_{1i}}{v_{2i}} > k_{12} \geq \max_{u_{1i}>0}\frac{-u_{2i}}{u_{1i}}$ <br> $\min_{v_{1i}>0}\frac{v_{2i}}{v_{1i}} \geq k_{21} > \max_{u_{2i}>0}\frac{-u_{1i}}{u_{2i}}$ <br> All spectra uncertain. Both spectra of ea. component given by extreme k's. | $\min_{v_{2i}>0}\frac{v_{1i}}{v_{2i}} > k_{12} \geq \max_{u_{1i}>0}\frac{-u_{2i}}{u_{1i}}$ <br> $\min_{v_{1i}>0}\frac{v_{2i}}{v_{1i}} = k_{21} = \max_{u_{2i}>0}\frac{-u_{1i}}{u_{2i}}$ <br> 1 spectrum of ea. comp. certain. 1 spectrum of ea. component given by intermediate k. | $\min_{v_{2i}>0}\frac{v_{1i}}{v_{2i}} > k_{12} \geq \max_{u_{1i}>0}\frac{-u_{2i}}{u_{1i}}$ <br> $\min_{v_{1i}>0}\frac{v_{2i}}{v_{1i}} \geq k_{21} > \max_{u_{2i}>0}\frac{-u_{1i}}{u_{2i}}$ <br> All spectra uncertain. 1 spectrum of ea. comp. given by ext. k, other by int. k. |
| $\min_{v_{2i}>0}\frac{v_{1i}}{v_{2i}} = k_{12} = \max_{u_{1i}>0}\frac{-u_{2i}}{u_{1i}}$ <br> $\min_{v_{1i}>0}\frac{v_{2i}}{v_{1i}} \geq k_{21} > \max_{u_{2i}>0}\frac{-u_{1i}}{u_{2i}}$ <br> 1 spectrum of ea. comp. certain. 1 spectrum of ea. component given by extreme k. | $\min_{v_{2i}>0}\frac{v_{1i}}{v_{2i}} = k_{12} = \max_{u_{1i}>0}\frac{-u_{2i}}{u_{1i}}$ <br> $\min_{v_{1i}>0}\frac{v_{2i}}{v_{1i}} > k_{21} > \max_{u_{2i}>0}\frac{-u_{1i}}{u_{2i}}$ <br> 1 spectrum of ea. comp. certain. 1 spectrum of ea. component given by intermediate k. | $\min_{v_{2i}>0}\frac{v_{1i}}{v_{2i}} \geq k_{12} > \max_{u_{1i}>0}\frac{-u_{2i}}{u_{1i}}$ <br> $\min_{v_{1i}>0}\frac{v_{2i}}{v_{1i}} > k_{21} \geq \max_{u_{2i}>0}\frac{-u_{1i}}{u_{2i}}$ <br> All spectra uncertain. Both spectra of ea. component given by extreme k's. | $\min_{v_{2i}>0}\frac{v_{1i}}{v_{2i}} \geq k_{12} > \max_{u_{1i}>0}\frac{-u_{2i}}{u_{1i}}$ <br> $\min_{v_{1i}>0}\frac{v_{2i}}{v_{1i}} > k_{21} \geq \max_{u_{2i}>0}\frac{-u_{1i}}{u_{2i}}$ <br> All spectra uncertain. 1 spectrum of ea. comp. given by ext. k, other by int. k. |
| $\min_{v_{2i}>0}\frac{v_{1i}}{v_{2i}} > k_{12} \geq \max_{u_{1i}>0}\frac{-u_{2i}}{u_{1i}}$ <br> $\min_{v_{1i}>0}\frac{v_{2i}}{v_{1i}} > k_{21} \geq \max_{u_{2i}>0}\frac{-u_{1i}}{u_{2i}}$ <br> All spectra uncertain. Both spectra of ea. component given by extreme k's. | $\min_{v_{2i}>0}\frac{v_{1i}}{v_{2i}} > k_{12} \geq \max_{u_{1i}>0}\frac{-u_{2i}}{u_{1i}}$ <br> $\min_{v_{1i}>0}\frac{v_{2i}}{v_{1i}} > k_{21} > \max_{u_{2i}>0}\frac{-u_{1i}}{u_{2i}}$ <br> All spectra uncertain. 1 spectrum of ea. comp. given by ext. k, other by int. k. | $\min_{v_{2i}>0}\frac{v_{1i}}{v_{2i}} > k_{12} \geq \max_{u_{1i}>0}\frac{-u_{2i}}{u_{1i}}$ <br> $\min_{v_{1i}>0}\frac{v_{2i}}{v_{1i}} > k_{21} \geq \max_{u_{2i}>0}\frac{-u_{1i}}{u_{2i}}$ <br> All spectra uncertain. 1 spectrum of ea. comp. given by ext. k, other by int. k. | $\min_{v_{2i}>0}\frac{v_{1i}}{v_{2i}} > k_{12} > \max_{u_{1i}>0}\frac{-u_{2i}}{u_{1i}}$ <br> $\min_{v_{1i}>0}\frac{v_{2i}}{v_{1i}} > k_{21} > \max_{u_{2i}>0}\frac{-u_{1i}}{u_{2i}}$ <br> All spectra uncertain. Both spectra of ea. component given by intermediate k's. |

*Excitation Overlap* (row axis)

**Figure 1**—Ambiguity Table of Two-Component Excitation-Emission Matrix.

490

## Rank Estimation by Frequency Analysis of the Eigenvectors [9]

The need to correctly determine the rank of a matrix for qualitative analysis has been demonstrated previously. Most rank estimation methods are statistical in nature and depend on the correct evaluation of the variance in the data. The method of rank analysis presented here uses Fourier Transform Image Analysis and is based on the assumption that the primary eigenvectors of an EEM (i.e., those associated with the fluorescence) will contain information that is weighted toward the low frequency Fourier coefficients since spectra tend to be broad-banded. Similarly, the transforms of the secondary eigenvectors will have larger high frequency coefficients. However, this manuscript is not a suitable medium for a comprehensive presentation of the theory of Fourier transforms. The reader is referred to any introductory text on the subject of Fourier analysis [10].

The Fourier transform is fairly simple to implement since fast Fourier transform algorithms are routinely available for use even on small computers. Since any continuous function can be reproduced by addition of a series of sine and cosine functions with various frequencies, amplitudes, and phases, the forward Fourier transform is a method of determining these frequencies, amplitudes, and phases from the function. The inverse transform reconstructs the time domain function from the frequencies, amplitudes, and phases. The discrete Fourier transform equation used to transform the eigenvector to the frequency domain is

$$V(u) = 1/N \sum_{x=0}^{N-1} v(x)\exp-2\pi i(xu)/N \qquad (1b)$$

where $u$ is the frequency domain coordinate, $V(u)$ is the Fourier transform of $v(x)$, and $N$ is the number of points in the discrete approximation of the function. The complex frequency domain function, $V(u)$, is frequently represented by the Fourier spectrum which is given by

$$|V(u)| = \left[V(u)^2_{real}+V(u)^2_{imag}\right]^{1/2} . \qquad (2b)$$

The area of the Fourier spectrum, $T$, is the sum of all the frequency coefficients and is given by

$$T = \sum_{u=0}^{N-1} |V(u)| . \qquad (3b)$$

The segment of the Fourier spectrum bounded by $+/-$ulim is denoted $A_{ulim}$ which is defined as

$$A_{ulim} = \sum_{u=u\text{lim}}^{u=u\text{lim}} |V(u)| \qquad (4b)$$

which is simply the sum of the frequency coefficients in the section of the Fourier spectrum bound by ulim and $-$ulim. The relative importance of this frequency region in reproducing the time domain eigenvector, $v(x)$, can be expressed by calculating the percent of $T$ which lies in this frequency range. The parameter which represents the importance of the range from ulim to $-$ulim is called %$A_{ulim}$ and is calculated from the expression

$$\%A_{ulim} = A_{ulim}/T \times 100 . \qquad (5b)$$

When the value of ulim is well chosen, there is a marked drop in the %$A_{ulim}$ for secondary eigenvectors. Table 1 shows a table comparing the rank estimation by frequency analysis to four statistical methods for matrices with ranks greater than or equal to 3. The frequency analysis method was the most accurate on the data tested here.

## Eigenvector Ratioing [11]

This method was developed to resolve the spectra of compounds from the ellipticity matrices of two component mixtures. The presence of legitimate negative values in this matrix produced a need for a different algorithm to analyze this matrix despite its similarity to the EEM. The ellipticity matrix, $F$, is also bilinear and can be expressed as

$$\mathbf{F=ST} \qquad (1c)$$

where the columns of the matrix $\mathbf{S}$ are the circular dichroism (CD) spectra, $s_k$ of the fluorophores in the sample, and the rows of the matrix $\mathbf{T}$, symbolized by the vectors $t_k$, are the emission spectra of those chiral fluorophores.

The CD and emission eigenvectors of $\mathbf{F}$ are also an orthonormal basis which span the matrix, $\mathbf{F}$. The CD eigenvectors, $q_k$, may be represented as the columns of the matrix $\mathbf{Q}$ and the emission eigenvectors, $p_k$, as the rows of the matrix $\mathbf{P}$. Thus, $\mathbf{F}$ is also given by the equation

$$\mathbf{F=OP} \qquad (2c)$$

assuming that the eigenvalues have been absorbed into one of the matrices.

These eigenvectors must also be transformed to possible spectral vectors. However, the CD eigenvectors should not be transformed to a non-negative vector because the CD spectral vectors are not always non-negative. However, since the matrix $\mathbf{F}$ contains a finite amount of information, if the correct emission spectral vectors are found, the corre-

| Mixture No. | No. of Components Known | REFAE | Stat1 | Stat2 | Stat3 | Stat4 |
|---|---|---|---|---|---|---|
| 1 | 3 | 3* | 3* | 2 | 3* | 4 |
| 2 | 3 | 3* | 2 | 2 | 3* | 4 |
| 3 | 3 | 3* | 2 | 2 | 2 | 2 |
| 4 | 2 | 2* | 1 | 1 | 2* | 2* |
| 5 | 2 | 2* | 1 | 1 | 2* | 3 |
| 6 | 2 | 2* | 2* | 2* | 2* | 5 |
| 7 | 3 | 3* | 3* | 2 | 2 | 5 |
| 8 | 4 | 4* | 3 | 3 | 3 | 5 |
| 9 | 5 | 5* | 3 | 3 | 3 | 5 |
| 10 | 6 | 5 | 3 | 3 | 4 | 6* |

"*" denotes a correctly estimated rank
Stat1 = Eigenvalue perturbation
Stat2 = Average error method
Stat3 = Chi squared test
Stat4 = Standard error test

sponding CD spectral vectors are fixed and easily obtained.

In this algorithm, the elements of the matrix, $K$, which transforms the emission eigenvectors to non-negative vectors are sought. Then, the possible CD spectral vectors are generated using the transformed emission vectors. The transformed emission vectors are the rows of the matrix $P'$ which is given by

$$P' = KP \geq 0 . \tag{3c}$$

For a two component mixture the values of the elements of $K$ can be determined by expressing eq (3c) in terms of the elements of the matrices and solving the resulting expressions for the elements of $K$ algebraically:

$$P_{1i}' = k_{11}p_{1i} + k_{12}p_{2i} \geq 0 , \tag{4c}$$

$$P_{2i}' = k_{22}p_{2i} + k_{21}p_{1i} \geq 0 . \tag{5c}$$

Again, assuming that $k_{11} = k_{22} = 1$, the elements of $K$ are within the ranges defined by

$$\min_{p_{2i}<0} \frac{-p_{1i}}{p_{2i}} \geq k_{12} \geq \max_{p_{2i}>0} \frac{-p_{1i}}{p_{2i}} \tag{6c}$$

$$\min_{p_{1i}<0} \frac{-p_{2i}}{p_{1i}} \geq k_{21} \geq \max_{p_{1i}>0} \frac{-p_{2i}}{p_{1i}} \tag{7c}$$

These ranges are generated because the sense of an inequality changes when both sides of the inequality are divided by a negative number.

The possible CD spectral vectors are generated by solving the following equation:

$$Q = FP'^T(P'P'^T)^{-1} . \tag{8c}$$

This is valid since

$$F = Q'P' \tag{9c}$$

and

$$I = P'P'^T(P'P'^T)^{-1} . \tag{10c}$$

This means that eq (8c) is another representation of the equation

$$Q' = Q'I . \tag{11c}$$

Due to the weaker constraints on the transformation elements in this algorithm, non-converging ranges are generated for the values of the transformation elements. However, it was found that if either of the components is a sole emitter in the monitored emission range, the value of the transformation element needed to transform an eigenvector to the spectra of the other component is usually given at an extreme of one of the ranges in eqs (6c) and (7c). This is because the values from the regions of sole emission best meet the criteria expressed in eqs (6c) and (7c). Therefore, except when the spectra are totally coincident at the base-lines at least one spectrum of each component is often retrievable using this technique. Figure 2 shows an ambiguity table that was generated to illustrate the usefulness of this algorithm.

This algorithm has also been tested on EEMs and it was found that this technique can be used with either set of eigenvectors: the emission eigenvectors (as it is with the ellipticity matrix) or the excitation eigenvectors. The results can be summarized by an ambiguity table similar to the one in figure 2 which had excitation overlap on the vertical axis. The ambiguity table using only one axis is more ambiguous than the earlier table (fig. 1) using both axes. These results do not conflict with those found using eigenvector analysis on the EEM. "Multiplying" the two ambiguity tables generated by this algorithm yields the table generated for eigenvector analysis, verifying the validity of both methods.

**Emission Overlap**



| $\min \dfrac{-p_{2i}}{p_{1j}<0 \;\; p_{1j}} \geq k_{12} \geq \max \dfrac{-p_{2i}}{p_{1i}>0 \;\; p_{1i}}$ | $\min \dfrac{-p_{2i}}{p_{1j}<0 \;\; p_{1j}} \geq k_{12} \geq \max \dfrac{-p_{2i}}{p_{1i}>0 \;\; p_{1i}}$ | $\min \dfrac{-p_{2i}}{p_{1j}<0 \;\; p_{1j}} \geq k_{12} \geq \max \dfrac{-p_{2i}}{p_{1i}>0 \;\; p_{1i}}$ | $\min \dfrac{-p_{2i}}{p_{1j}<0 \;\; p_{1j}} \geq k_{12} \geq \max \dfrac{-p_{2i}}{p_{1i}>0 \;\; p_{1i}}$ |
|---|---|---|---|
| $\min \dfrac{-p_{1i}}{p_{2j}<0 \;\; p_{2j}} \geq k_{21} \geq \max \dfrac{-p_{1i}}{p_{2i}>0 \;\; p_{2i}}$ | $\min \dfrac{-p_{1i}}{p_{2j}<0 \;\; p_{2j}} \geq k_{21} \geq \max \dfrac{-p_{1i}}{p_{2i}>0 \;\; p_{2i}}$ | $\min \dfrac{-p_{1i}}{p_{2j}<0 \;\; p_{2j}} \geq k_{21} \geq \max \dfrac{-p_{1i}}{p_{2i}>0 \;\; p_{2i}}$ | $\min \dfrac{-p_{1i}}{p_{2j}<0 \;\; p_{2j}} \geq k_{21} \geq \max \dfrac{-p_{1i}}{p_{2i}>0 \;\; p_{2i}}$ |
| All spectra uncertain Both spectra of ea. component given by extreme k's | All spectra uncertain. 1 spectrum of ea. comp. given by ext. k, other by int. k. | All spectra uncertain. 1 spectrum of ea. comp. given by ext. k, other by int. k | All spectra uncertain. Both spectra of ea. component given by intermediate k's. |

Figure 2—Ambiguity Table of Two-Component Ellipticity Matrix.

## Conclusions

This manuscript has provided an overview of qualitative analysis techniques developed for matrix formatted fluorescence data. Qualitative analysis of matrices was shown to generally consist of three basic procedures: rank estimation, determination of unknown component spectra, and screening of expected compounds. The techniques outlined here only addressed the first two phases of the problem but they represent only a portion of the methods that have been developed to fill these requirements. The methods presented here successfully attack the stated problems within the framework of the limitations described.

Rank estimation by frequency analysis can sometimes be more accurate than statistical methods for evaluation of spectral data. It would be useful to develop a criterion where the algorithm will automatically select a useful range for differentiating the secondary eigenvectors from the primary.

In the present approach, eigenvector analysis and eigenvector ratioing are limited to binary mixtures. Few real samples are binary. These methods must be extended to higher order mixtures.

These techniques have been developed for use with fluorescence data, but are generally applicable to other forms of matrix formatted data. Some of the algorithms require that the data matrix be bilinear in form; however this is a characteristic of many types of data. For example, diode array detection of liquid chromatography and absorption kinetic data using a diode array detector are bilinear in form. Clearly, there is a need for further development in this area.

## References

[1] Parker, C.A., Photoluminescence of Solutions with Applications to Photochemistry and Analytical Chemistry, American Elsevier, New York (1968).

[2] Warner, I.M.; J.B. Callis, E.R. Davidson, M. Gouterman, and G.D. Christian, Anal. Lett., 8, 665 (1965).

[3] Malinowski, E.R., and D.G. Howery, Factor Analysis in Chemistry, John Wiley, New York (1980).

[4] Jurs, P.C., and T. L. Isenhour, Chemical Applications of Pattern Recognition, John Wiley, New York, NY, 1975.

[5] Tinoco, I. Jr. and D.H. Turner, J. Am. Chem. Soc., 98(21), 6453 (1976).

[6] Thomas, M.; G. Patonay and I.M. Warner, work in progress.

[7] Warner, I.M., "The Analysis of Matrix Formatted Multicomponent Data," Contemporary Topics in Analytical and Clinical Chemistry, D.M. Hercules, G.M. Hieftje, L.R. Snyder and M.A. Evenson, eds., Plenum, New York (1982).

[8] Fukunaga, K., Introduction to Statistical Pattern Recognition, Academic Press, New York (1972).

[9] Rossi, T.M., and I.M. Warner, "Rank Determination of Emission Excitation Matrices via Fourier Transformation of Eigenvectors," American Chemical Society, 187th National Meeting, April 8-13, 1984, St. Louis, MO.

[10] Bloomfield, P., Fourier Analysis of Time Series: An Introduction, John Wiley and Sons, New York (1976).

[11] Neal, S.L.; C.-N. Ho and I.M. Warner, "Qualitative Analysis of Multicomponent Fluorescence Detected Circular Dichroism Data," American Chemical Society, 187th National Meeting, April 8-13, 1984, St. Louis, MO.

# Some New Ideas
# in the Analysis of Screening Designs

**George Box**
University of Wisconsin-Madison, Madison, WI 53705
and
**R. Daniel Meyer**
Lubrizol Corp., Wickliffe, OH 44092

Consideration of certain aspects of scientific method leads to discussion of recent research on the role of screening designs in the improvement of quality. A projective rationale for the use of these designs in the circumstances of *factor sparsity* is advanced. In this circumstance the possibility of identification of sparse *dispersion* effects as well as sparse *location* effect is considered. A new method for the *analysis of fractional factorial designs* is advanced.

Key words: factor sparsity; fractional factorial designs; screening designs; sparse dispersion; sparse location.

Humans differ from other animals most remarkably in their ability to learn. It is clear that although throughout the history of mankind technological learning has taken place, although until three or four hundred years ago change occurred very slowly. One reason for this was that in order to learn something - for example, how to make fire or champagne - two *rare events* needed to coincide: (a) an informative event had to *occur*, and (b) a person able to draw logical conclusion and to act on them had to be *aware* of that informative event.

Passive surveillance is a way of increasing the probability that the rare informative event will be constructively taken note of and is exemplified by quality charting methods. Thus a Shewhart chart is a means to ensure that possibly informative events are brought to the attention of those who may be able to discover in them an "assignable cause" [1][1] and act appropriately.

---

**About the Author, Paper:** George Box, who is with the Research Center for Quality and Productivity Improvement at the University of Wiconsin-Madison, and R. Daniel Meyer are both Statisticians. The work they describe was sponsored by U.S. Army Contract DAAC 29-80-C-0041 and National Science Foundation Grant DMS-8420968.

---

[1]Figures in brackets indicate literature references.

Active intervention by experimentation aims, in addition, to increase the probability of an informative event *actually occurring*. A designed experiment conducted by a qualified experimenter can dramatically increase the probability of learning because it increases simultaneously the probability of an informative event occurring and also the probability of the event being constructively witnessed. Recently there has been much use of experimental design in Japanese industry particularly by Genichi Taguchi [2] and his followers. In off-line experimentation he has in particular emphasized the use of highly fractionated designs and orthogonal arrays and the minimization of variance.

In the remainder of this paper we briefly outline some recent research on the use of such screening designs.

## 1. Use of Screening Designs to Identify "Active" Factoring

Table 1 shows in summary a highly fractionated two-level factorial design employed as a screening design in an off-line welding experiment performed by the National Railway Corporation of Japan [2]. In the column to the right of the table is shown the observed tensile strength of the weld, one of several quality characteristics measured.

The design was chosen on the assumption that in addition to main effects only the two-factor interactions AC, AG, AH, and GH were expected to be present. On that supposition, all nine main effects and the four selected two-factor interactions can be separately estimated by appropriate orthogonal contrasts, the two remaining contrasts corresponding to the columns labelled $e_1$ and $e_2$ measure only experimental error. Below the table are shown the grand average, the 15 effect contrasts, and the effects plotted on a dot diagram. The effects plotted on normal probability paper suggested that, over the ranges studied, only factors B and C affect tensile location by amounts not readily attributed to noise.

If this conjecture is true, then, at least appoximately, the 16 runs could be regarded as four replications of a $2^2$ factorial design in factors B and C only. However, when the results are plotted in figure 1 so as to reflect this, inspection suggests the existence of a dramatic effect of a different kind—when factor C is at its plus level the spread of the data appears much larger[3] than when it is at its minus level. Thus, in addition to detecting shift in location due to B and C, the experiment may also have detected what we will call a *dispersion effect* due to C. The example raises the general possibility of analyzing unreplicated designs for dispersion effects as well as for the more usual location effects.

---

[2]To facilitate later discussion, we have set out the design and labelled the levels somewhat differently from [2].

[3]Data of this kind might be accounted for by the effect of one or more variables other than B that affected tensile strength only at the "plus level" of C (only when the alternative material was used). Analysis of the eight runs made at the plus level of C does not support this possibility, however.

## 2. Rationales for Using Screening Designs

Before proceeding we need to consider the question, "In what situations are screening designs, such as highly fractionated factorials, useful?"

**2.1. Effect Sparsity.** A common industrial problem is to find from a rather large number of factors those few that are responsible for *large effects*. The idea is comparable to that which motivates the use in quality control studies of the "Pareto diagram." (See, for example, [3]). The situation is approximated by postulating that only a small proportion of effects will be *"active"* and the rest *"inert"*. We call this the postulate of *effect sparsity*. For studying such situations, highly fractionated designs and other orthogonal arrays [2,4,5,6] which can screen moderately large numbers of variables in rather few runs are of great interest. Two main rationalizations have been suggested for the use of these designs; both ideas rely on the postulate of effect sparsity but in somewhat different ways.

**2.2. Rationale Based on Prior Selection of Important Interactions.** It is argued (see for example [7]) that in some circumstances physical knowledge of the process will make only a few interactions likely and that the remainder may be assumed negligible. For example, in the welding experiment described above there were 36 possible two-factor interactions between the nine factors, but only four were regarded as likely, leaving 32 such interactions assumed
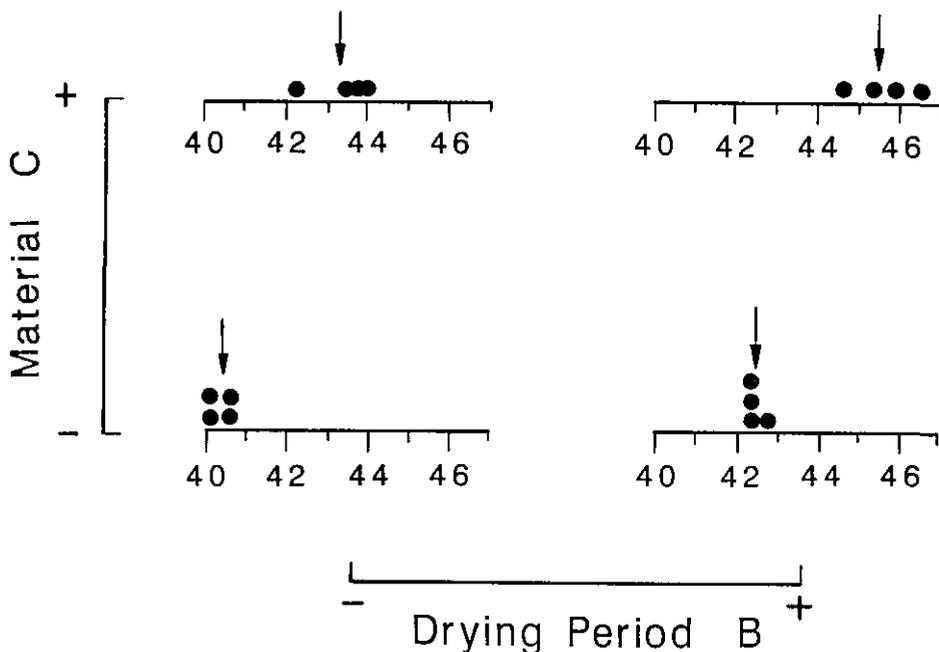


Figure 1—Tensile strength data as four replicates of a $2^2$ factorial design in factors B and C only.

negligible. The difficulty with this idea is that in many applications the picking out of a few "likely" interactions is difficult if not impossible. Indeed the investigator might justifiably protest that, in the circumstance where an experiment is needed to determine which *first order* (main) effects are *important, it is illogical that he be expected to guess in advance which effects of *second order* (interactions) are important.

**2.3. Projective Rationale Factor Sparsity.** A somewhat different notion is that of *factor* sparsity. Thus suppose that, of the $k$ factors considered, only a small subset of vaguely known size $d$, *whose identity is, however, unknown*, will be active in providing main effects and interactions within the subset. Arguing as in [8], a two-level design enabling us to study such a system is a fraction of resolution $R = d+1$ (or in the terminology of [6], an array of strength $d$) which produces complete factorials (possibly replicated) in every one of the $\binom{k}{d}$ spaces of $d = R - 1$ dimensions. For example, we have seen that on the assumption that only factors B and C are important, the welding design could be regarded as four replicates of a $2^2$ factorial in just those two factors. But because the design is of resolution $R = 3$ the same would have been true for any of the 36 choices of two out of the nine factors tested. Thus the design would be appropriate if it were believed that not more than two of the factors were likely to be "active".

For further illustration we consider again the 16-run orthogonal array of table 1 and, adopting a roman subscript to denote the resolutions of the design, we indicate in table 2 various ways in which that array might be used. It may be shown that

(a) If we associated the 15 contrast columns of the design with 15 factors, we would generate a $2_{III}^{15-11}$ design providing two-fold replication of $2^2$ factorials in every one of the 105 two-dimensional projections.

(b) If we associated only columns 1, 2, 4, 7, 8, 11, 13, and 14 with eight factors we would generate a $2_{IV}^{8-4}$ design providing two-fold replication of $2^3$ factorials in every one of the 56 three-dimensional projections.

(c) If we associated only columns 1, 2, 4, 8, and 15 with five factors we would generate a $2_V^{5-1}$ design providing a $2^4$ factorial in every one of the four-dimensional projections.

(d) If we associated only columns 1, 2, 4, and 8 with four factors we would obtain the complete $2^4$ design from which this orthogonal array was in fact generated.

Designs (a), (b), and (c) would thus be appropriate for situations where we believed respectively that not more than 2, 3, or 4 factors would be active[4]. Notice that intermediate

---

[4] The designs give partial coverage for a larger number of factors. For example ([8] (1961)) 56 of the 70 four-dimensional projections of the $2_{IV}^{8-4}$ yield a full factorial in four variables.

---

**Table 1.** A fractional two-level design used in a welding experiment showing observed tensile strength and effects.

A: Kind of Welding Rods
B: Period of Drying
C: Welded Material
D: Thickness
E: Angle
F: Opening
G: Current
H: Welding Method
J: Preheating

| Factor<br>Column Number | 0 | D 1 | H 2 | e₁ 3 | G 4 | F 5 | GH 6 | AC 7 | A 8 | E 9 | AH 10 | e₂ 11 | AG 12 | J 13 | B 14 | C 15 | Tensile strength kg/mm² |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | + | − | − | + | − | + | + | − | − | + | + | − | + | − | − | + | 43.7 |
| 2 | + | + | − | − | − | − | + | + | − | − | + | + | + | + | − | − | 40.2 |
| 3 | + | − | + | − | − | + | − | + | − | + | − | + | + | − | + | − | 42.4 |
| 4 | + | + | + | + | − | − | − | − | − | − | − | − | + | + | + | + | 44.7 |
| 5 | + | − | − | + | + | − | − | + | − | + | + | − | − | + | + | − | 42.4 |
| 6 | + | + | − | − | + | + | − | − | − | − | + | + | − | − | + | + | 45.9 |
| 7 | + | − | + | − | + | − | + | − | − | + | − | + | + | + | − | + | 42.2 |
| 8 | + | + | + | + | + | + | + | + | − | − | − | − | − | − | − | − | 40.6 |
| 9 | + | − | − | + | − | + | + | − | + | − | − | + | − | + | + | − | 42.4 |
| 10 | + | + | − | − | − | − | + | + | + | + | − | − | − | − | + | + | 45.5 |
| 11 | + | − | + | − | − | + | − | + | + | − | + | − | − | + | − | + | 43.6 |
| 12 | + | + | + | + | − | − | − | − | + | + | + | + | − | − | − | − | 40.6 |
| 13 | + | − | − | + | + | − | − | + | + | − | − | + | + | − | − | + | 44.0 |
| 14 | + | + | − | − | + | + | − | − | + | + | − | − | + | + | − | − | 40.2 |
| 15 | + | − | + | − | + | − | + | − | + | − | + | − | + | − | + | − | 42.5 |
| 16 | + | + | + | + | + | + | + | + | + | + | + | + | + | + | + | + | 46.5 |
| Effect | 43.0 | .13 | −.15 | −.30 | −.15 | .40 | −.03 | .38 | .40 | −.05 | .43 | .13 | .13 | −.38 | 2.15 | 3.10 | |

(Run)

497

Table 2. Some alternative uses of the orthogonal array.

| Columns | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| (a) $2^{15-11}_{III}$ | • | • | • | • | • | • | • | • | • | • | • | • | • | • | • |
| (b) $2^{8-4}_{IV}$ | • | • | | • | | | • | • | | | • | | | • | • |
| (c) $2^{5-1}_{V}$ | • | • | | • | | | | • | | | | | | | • |
| (d) $2^{4}$ | • | • | | • | | | | • | | | | | | | |

values of $k$ could be accommodated by suitably omitting certain columns. Thus the welding design is a $2^{9-5}_{III}$ arrangement which can be obtained by omitting six columns from the complete $2^{15-11}_{III}$. Notice finally that for intermediate designs we can take advantage of both rationales by arranging, as was done for the welding design, that particular interactions are isolated.

A discussion of the iterative model building process [9] characterized three steps in the iterative data analysis cycle indicated below

┌→ identification ——→ fitting ——→diagnostic checking ←┐
└ ─ ─ ─ ─ ─ ─ ─ ─ ─ ─ ─ ─ ─ ─ ─ ─ ─ ─ ─ ─ ─ ─ ─ ─ ─ ┘

Most of the present paper is concerned with model identification - the search for a model worthy to be formally entertained and fitted by an efficient procedure such as maximum likelihood. The situation we now address concerns the analysis of fractional designs such as the welding design in the

above context when only a few of the factors are likely to have effects but these may include dispersion effects as well as location effects.

## 3. Dispersion Effects

We again use the design of table 1 for illustration. There are 16 runs from which 16 quantities—the average and 15 effect contrasts—have been calculated. Now if we were also interested in possible dispersion effects we could also calculate 15 variance ratios. For example, in column 1 we can compute the sample variance $s^2_{1-}$ for those observations associated with a minus sign and compare it with the sample variance $s^2_{1+}$ for observations associated with a plus sign to provide the ratio $F_1 = s^2_{1-} / s^2_{1+}$. If this is done for the welding data we obtain values for $\ln F_i$ given in figure 2(a).[5] It will be recalled that in the earlier analysis a large dispersion effect associated with factor C (column 15) was found, but in figure 2(a) the effect for factor C is not especially extreme, instead the dispersion effect for factor D (column 1) stands out from all the rest. This misleading indication occurs because we have not so far taken account of the aliasing of location and dispersion effects. Since 16 linearly independent location effects have already been calculated for the original data, calculated dispersion effects must be functions

[5]In this figure familiar normal theory significance levels are also shown. Obviously the necessary assumptions are not satisfied in this case, but these percentages provide a rough indication of magnitude.
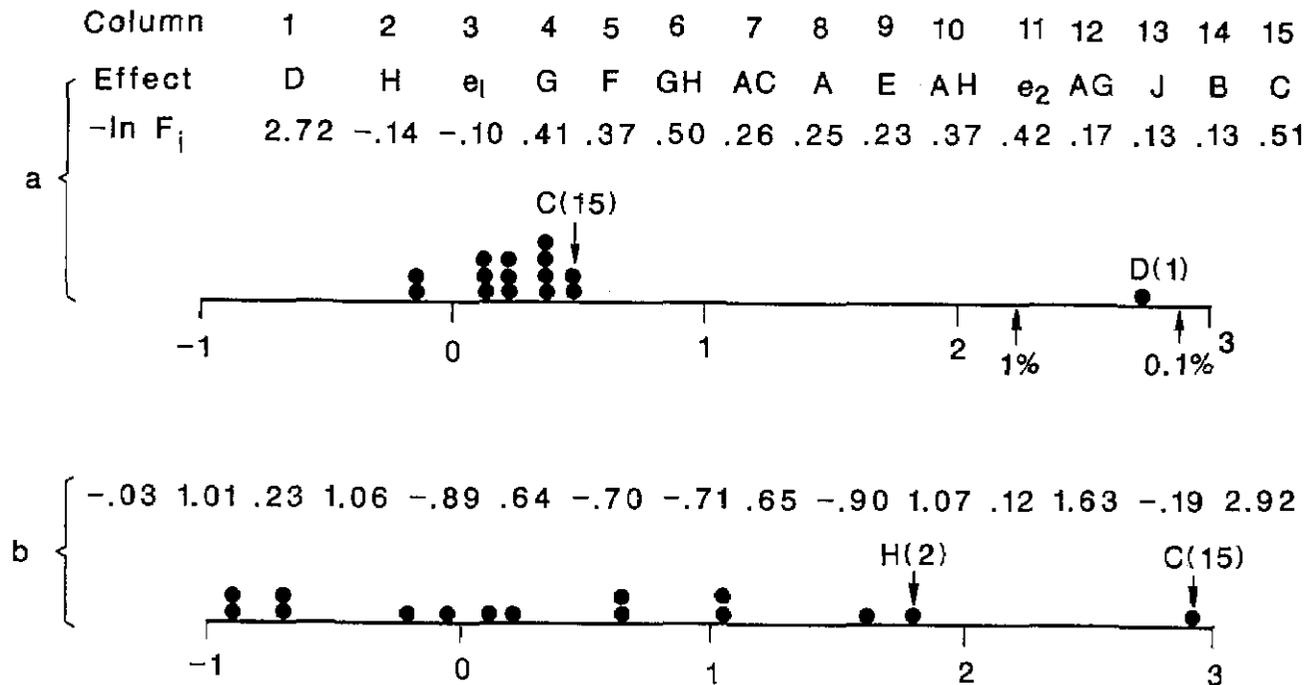
| Column | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Effect | D | H | $e_l$ | G | F | GH | AC | A | E | AH | $e_2$ | AG | J | B | C |
| $-\ln F_i$ | 2.72 | -.14 | -.10 | .41 | .37 | .50 | .26 | .25 | .23 | .37 | .42 | .17 | .13 | .13 | .51 |

a {

C(15)

D(1)

-1    0    1    2 ↑ 1%    ↑ 0.1%   3

b {

| -.03 | 1.01 | .23 | 1.06 | -.89 | .64 | -.70 | -.71 | .65 | -.90 | 1.07 | .12 | 1.63 | -.19 | 2.92 |

H(2)     C(15)

-1    0    1    2    3

Figure 2—Welding experiment log dispersion effects (a) before, and (b) after elimination of location effects for B and C.

498

of these. Recently [10] a general theory of location-dispersion aliasing has been obtained for factorials and fractional factorials at two levels. For illustration, in this particular example it turns out that the following identity exists for the dispersion effect $F_1$, that is the $F$ ratio associated with factor D and hence for column 1 of the design.

$$F_1 = \frac{(\hat{2}-\hat{3})^2+(\hat{4}-\hat{5})^2+(\hat{6}-\hat{7})^2+(\hat{8}-\hat{9})^2+(\hat{10}-\hat{11})^2+(\hat{12}-\hat{13})^2+(\hat{14}-\hat{15})^2}{(\hat{2}+\hat{3})^2+(\hat{4}+\hat{5})^2+(\hat{6}+\hat{7})^2+(\hat{8}+\hat{9})^2+(\hat{10}+\hat{11})^2+(\hat{12}+\hat{13})^2+(\hat{14}+\hat{15})^2}$$

(1)

Now (see table 1) $\hat{14}=\hat{B}=2.15$ and $\hat{15}=\hat{C}=3.10$ are the two largest location effects, standing out from all the others. The extreme value of $F_1$ associated with an apparent dispersion effect of factor D(1) is largely accounted for by the squared sum and squared difference of the location effects B and C which appear respectively as the last terms in the denominator and numerator of eq (1). A natural way to proceed is to compute variances from the residuals obtained after eliminating large location effects. After such elimination the alias relations of eq (1) remain the same except that location effects from eliminated variables drop out. That is, zeros are substituted for eliminated variables. Variance analysis for the residuals after eliminating effects of B and C are shown in figure 2(b). The dispersion effect associated with C (factor 15) is now correctly indicated as extreme. It is shown in the paper referenced above how, more generally, under circumstances of effect sparsity a location-dispersion model may be correctly identified when a few effects of both kinds are present.

## 4. Analysis of Unreplicated Fractional Designs

Another important problem in the analysis of unreplicated fractional designs and other orthogonal arrays concerns the picking out of "active" factors. A serious difficulty is that with unreplicated fractional designs no simple estimate of the experimental error variance against which to judge the effects is available.

In one valuable procedure due to Cuthbert Daniel [11,12] effects are plotted on normal probability paper. For illustration table 3 shows the calculated effects from a $2_{IV}^{8-4}$ design used in an experiment on injection molding [13, p. 379]. These effects are plotted on normal probability paper in figure 3.

An alternative Bayesian approach [14] is as follows: Let $T_1, T_2, \ldots, T_\nu$ be standardized[6] effects with

$$T_i = e_i \text{ if effect inert}$$

$$T_i = e_i + \tau_i \text{ if effect active}$$

---

[6]For three-level and mixed two and three level designs for example, this analysis is carried out after the effects are scaled so that they all have equal variances.

Table 3. Calculated effects from a $2_{IV}^{8-4}$ design showing alias structure assuming three factor and higher order interactions negligible, injection molding experiment.

| | | |
|---|---|---|
| $T_1 = -0.7 \rightarrow 1$ | mold temp. |
| $T_2 = -0.1 \rightarrow 2$ | moisture content |
| $T_3 = 5.5 \rightarrow 3$ | holding pressure |
| $T_4 = -0.3 \rightarrow 4$ | cavity thickness |
| $T_5 = -3.8 \rightarrow 5$ | booster pressure |
| $T_6 = -0.1 \rightarrow 6$ | cycle time |
| $T_7 = 0.6 \rightarrow 7$ | gate size |
| $T_8 = 1.2 \rightarrow 8$ | screw speed |

$$T_9 = T_{1.2} = -0.6 \rightarrow 1.2 + 3.7 + 4.8 + 5.6$$

$$T_{10} = T_{1.3} = 0.9 \rightarrow 1.3 + 2.7 + 4.6 + 5.8$$

$$T_{11} = T_{1.4} = -0.4 \rightarrow 1.4 + 2.8 + 3.6 + 5.7$$

$$T_{12} = T_{1.5} = 4.6 \rightarrow 1.5 + 2.6 + 3.8 + 4.7$$

$$T_{13} = T_{1.6} = -0.3 \rightarrow 1.6 + 2.5 + 3.4 + 7.8$$

$$T_{14} = T_{1.7} = -0.2 \rightarrow 1.7 + 2.3 + 6.8 + 4.5$$

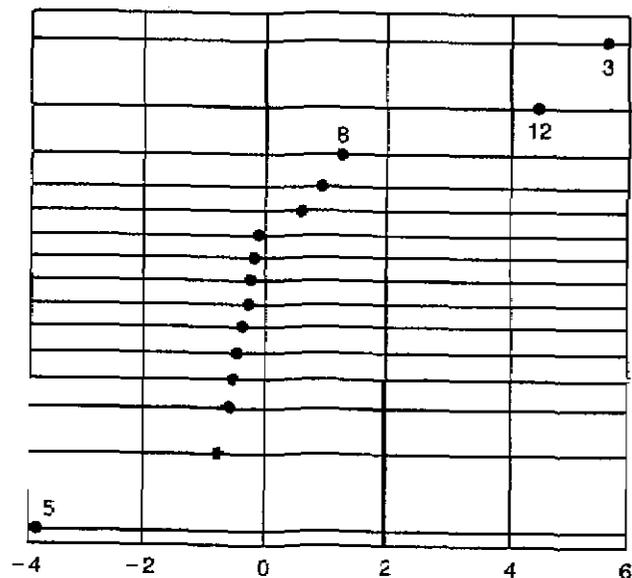$$T_{15} = T_{1.8} = -0.6 \rightarrow 1.8 + 2.4 + 3.5 + 6.7$$



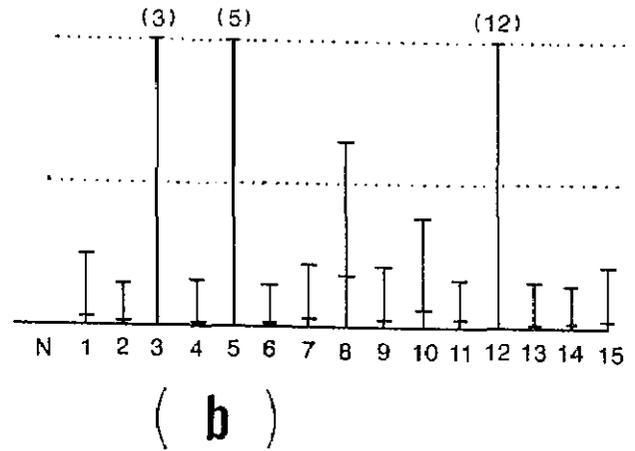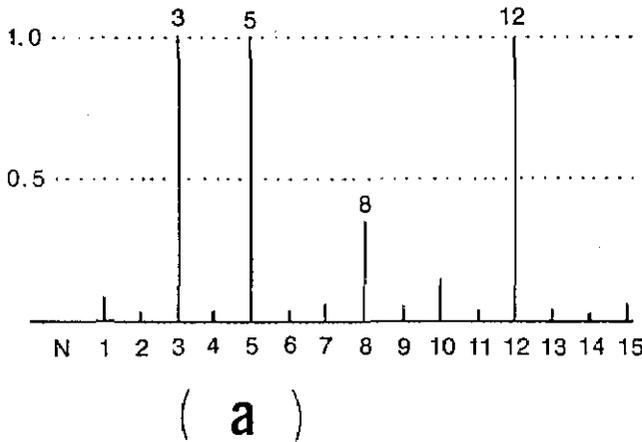Figure 3—Normal plot of effects. Injection molding experiment.

499

Figure 4—(a) Welding experiment. Posterior probability that factor $i$ is active (a = 0.30, k = 10). (b) Sensitivity analysis for posterior probability (a = .15 − .45, k = 5 − 15).

$$e_i \rightarrow N(0,\sigma^2), \quad \tau_i \rightarrow N(0,\sigma_\tau^2) \quad k^2 = \frac{\sigma^2 + \sigma_\tau^2}{\sigma^2} .$$

Suppose the probability that an effect is active is $\alpha$.

Let $a_{(r)}$ be the event that a particular set of $r$ of the $v$ factors are active, and let $T_{(r)}$ be the vector of estimated effects corresponding to active factors of $a_{(r)}$. Then, [15] with $p(\sigma) \propto 1/\sigma$ the posterior probability that $T_{(r)}$ are the only active effects is:

$$P[a_{(r)}|T,\alpha,k] \propto \left[\frac{\alpha k^{-1}}{1-\alpha}\right]^r \left\{1 - \left(1-\frac{1}{k^2}\right)\frac{S_{(r)}}{S}\right\}^{-\frac{v}{2}} ,$$

where $S_{(r)} = T'_{(r)}T_{(r)}$ and $S = T'T$. In particular the marginal probability that an effect $i$ is active give $T$, $\alpha$ and $k$ is proportional to

$$\sum_{\substack{a_{(r)} \\ i\ active}} \left[\frac{\alpha k^{-1}}{1-\alpha}\right]^r \left\{1 - \left(1-\frac{1}{k^2}\right)\frac{S_{(r)}}{S}\right\}^{-\frac{v}{2}} .$$

A study of the fractional factorials appearing in [7,12,13]. suggested that $\alpha$ might range from 0.15-0.45 while $k$ might range from 5 to 15. The posterior probabilities

computed with the (roughly average) values, $\alpha = 0.30$ and $k = 10$ are shown in figure 4(a) in which $N$ denotes the probability (negligible for this example) that there are no active effects. The results from a sensitivity analysis in which $\alpha$ and $k$ were altered to vary over the ranges mentioned above is shown in figure 4(b).

It will be seen that figure 4(a) points to the conclusion that active effects are associated with columns 3, 5 and 12 of the design and that column 8 might possibly also be associated with an active factor. Figure 4(b) suggests that this conclusion is very little affected by widely different choices for $\alpha$ and $k$. Further research with different choices of prior, with marginization with respect to $k$, and with different choices of the distribution assumptions is being conducted.

## 5. Allowance for Faulty Observations

Recent work [16] has shown how a double application of the scale-contamination model (both to the observations themselves as well as to the affects) can make it possible to allow for faulty observations in the analysis of unreplicated factorials or fractional factorials.

## References

[1] Shewhart, W.A., Economic Control of Quality of Manufactured Product, D. Van Nostrand Company, Inc. New York, NY (1931).

[2] Taguchi, G., and Y. Wu, Introduction to Off-Line Quality Control, Central Japan Quality Control Association, Nagoya, Japan (1980).

[3] Ishikawa, K., Guide to Quality Control, Asian Productivity Organization, Tokyo (1976).

[4] Finney, D.J., The Fractional Replication of Factorial Arrangements.

Annals of Eugenics, 12, 4 291-301 (1945).

[5] Plackett, R.L., and J.P. Burman, Design of Optimal Multifactorial Experiments. Biometrika, 23 305-325 (1946).

[6] Rao, C.R., Factorial Experiments Derivable from Combinatorial Arrangements of Arrays. J. Roy. Statist. Soc., 89 128-140 (1947).

[7] Davies, O.L., The Design and Analysis of Industrial Experiments. London: Oliver and Boyd (1954).

[8] Box, G.E.P., and J.S., Hunter, The $2^{k-p}$ Fractional Factorial Designs. Technometrics, 3 311-351, 449-458 (1961).

# Vijayan Nair

AT&T Bell Laboratories

and

# Michael Frenklach

Department of Materials Science and Engineering
Pennsylvania State University[1]

**Nair:** The paper by Box and Meyer deals with some interesting problems that arise in the off-line quality control methods introduced by Taguchi (see Taguchi and Wu [1][2]). My comments will be restricted to the first part of the paper, viz. estimating dispersion effects.

### 1) Estimating dispersion effects for quality control

In industrial experiments designed to detect important factors that affect the quality of a manufacturing/ production process, the estimation of dispersion effects is as important as the estimation of location effects. In fact in situations where there are readily identifiable signal factors (see [1]), the primary goal is in estimating dispersion effects. The location effects, in this case, play the role of nuisance parameters. León, Shoemaker and Kackar [2] offer an excellent discussion of the statistical formulation of the parameter design problem in industrial experiments.

### 2) Effect sparsity

The Box-Meyer techniques exploit the notion of effect sparsity to obtain "replicates" in an unreplicated experiment. It is likely that in most cases only a few factors are *highly* significant. However, in many situations, one could also expect many of the other factors included in the experiment to have sizeable effects. This is particularly true when a fair amount of the information about the process is known and is used in the selection of the factors. In such situations, one could not reasonably expect to estimate both location and dispersion effects from an unreplicated experiment.

### 3) When to log?

Box suggests using $\log[S^2(i-)/S^2(i+)]$ as a preliminary estimate of the dispersion effects. An alternative method would be to take log of the squared residuals and do ANOVA with an additive model for the log of the scale parameters. This is the type of analysis usually done in experiments with replications. Some efficiency calculations suggest that the Box-Meyer analysis is more efficient when there are only a few large dispersion effects and less efficient when there are many.

### 4) Iterating

It is possible that during the first step of the iteration (which does an unweighted analysis) some significant location effects are not detected. So after the dispersion effects are estimated, the location model should be refitted for all the factors.

### 5) Transformations

In replicated experiments, the quasi-likelihood models of Nelder and Pregibon [3] allow one to determine the transformations under which the location effects and the scale effects are approximately additive.

# References

[1] Taguchi, G., and Y. Wu, *Introduction to off-line quality control*. Central Japan Quality Control Association (1980).

[2] León, R. V.; A. C. Shoemaker and R. N. Kackar, Performance measures independent of adjustment: An alternative to Taguchi's signal to noise ratio. *AT&T Bell Laboratories Technical Memorandum* (1985).

[3] Nelder, J .A., and D. Pregibon, Quasi-likelihood models and data analysis. Biometrika (to appear).

---

[1] Michael Frenklach's contribution to the subject stems from work performed in the Department of Chemical Engineering, Louisiana State University.

[2] Figures in brackets indicate literature references.

**Frenklach:** Professor Box presented a method of analysis of factorial designs for detection of main effects and faulty observations. The approach is analytical and provides numerical measures for what previously has been approached graphically. The availability of the analytical algorithm is important for computerization of the analysis.

The central hypothesis of the method is what the authors call effect sparsity, which states that usually only a small number of input and control process variables would have a significant effect on the process response(s). This situation appears to be true not only in experimental environments but also in computer modeling of various industrial processes and natural phenomena. Mechanistic models usually take the form of differential equations for which no analytical solution is available. The model may contain a large number of (physical) parameters and it is not always obvious from a simple inspection of the computational results what effect each parameter has on a given response or responses. Sensitivity analysis has been used to reveal this information. Among other techniques, the use of screening factorial designs for sensitivity analysis of computer models has been suggested by Box et al. 1978; Frenklach 1984; Frenklach and Bornside, 1984; Miller and Frenklach, 1983; and Morris and Mitchell, 1983 [1-5].

The present experience with chemical kinetic modeling, for example, is as follows. Due to technical difficulties of instrumentation, there are only a few experimental responses available, typically one or two. The cases studied indicate that it is a very small number of chemical reactions, out of hundreds of reactions comprising the model, whose rate coefficient values, within their uncertainty intervals, have significant or "active" effects on the experimentally verifiable model responses. These are exactly the conditions of effect sparsity discussed above. Thus, the method presented by Box is well-suited not only for "real" experimentation, but also for computer modeling.

Computations, however, do not have random errors. Does this fact simplify and economize the analysis? Should special methods and designs be developed or existing ones modified for a most efficient use in screening analysis of computer models?

## References

[1] Box, G. E. P.; W. C. Hunter and J. S. Hunter, Statistics for Experimenters. An Introduction to Design, Data Analysis, and Model Building, Wiley, New York, pp. 429–432 (1978).

[2] Frenklach, M., Modeling, Combustion Chemistry (Edited by W. C. Gardiner, Jr.), Springer-Verlag, New York, Chap. 7 (1984).

[3] Frenklach, M., and D. E. Bornside, Shock-initiated ignition in methane-propane mixtures, Combust. Flame 56, 1–27 (1984).

[4] Miller, D., and M. Frenklach, Sensitivity analysis and parameter estimation in dynamic modeling of chemical kinetics, Int. J. Chem. Kinet. 15, 677–696 (1983).

[5] Morris, M. D., and T. J. Mitchell, Two-level multifactor designs for detecting the presence of interaction, Technometrics 25, 345–355 (1983).

# Polymers and Random Walks—Renormalization Group Description and Comparison With Experiment

**Karl F. Freed**

The University of Chicago, Chicago, IL 60637

Although real polymers involve the sequential addition of monomers having fixed bond lengths, fixed bond angles and some freedom of rotation about single bond, the properties of polymers over large length scales can be modeled by treating the polymer configuration as that of a random walk formed by the monomer units. Serious complications arise in the theoretical description of these polymers because of excluded volume constraints which prohibit different monomers from occupying the same position in space. This polymer excluded volume problem has been modeled in terms of a simple continuous random walk with short range repulsive interactions. The expansion of polymer properties in this repulsive interaction can readily be shown by dimensional analysis to involve an expansion in a large parameter, in the limit of long polymers. The renormalization group method is utilized as a systematic means for resuming this divergent perturbation expansion. The theory proceeds by analytically continuing the treatment to continuous range of spatial dimensionalities to expose and regularize the singularities in the analytically continued theory. The renormalization group approach is described from a heuristic physical standpoint and extensive comparisons are provided to show how it quantitatively reproduces vast amounts of dilute solution polymer properties with no adjustable parameters.

Key words: experiment comparison; modeling; monomer units; polymer properties; random walk; repulsive interaction.

The study of the configuration statistics of polymers in dilute solutions presents problems of interest to analytical chemists, chemical physicists, engineers and applied mathematicians. Roughly half of the American chemical industry is involved with polymers, and analytical chemists are concerned with the characterization of their properties. Dilute solution physical properties are used to determine the molecular weight and general shape and architecture of the polymer. Hence, the availability of a quantitative theoretical description of the dependence of dilute solution physical properties on solvent characteristics, molecular weight, temperature, and branching is an important aid in characterizing and understanding the properties of polymers.

The interest of polymers to applied mathematicians lies in their mathematical description in terms of interacting random walks.[1][^1] This mathematical representation of a polymer chain can be motivated by visualizing the polymer as a sequence of bonded monomer units. Each new monomer, added to the chain, is attached by a chemical bond of fixed length, generally having fixed bond angles with respect to the previous bond. However, there is considerable freedom of rotation about the individual bonds thereby generating a large number of configurations for this idealized random walk polymer model [2]. The individual bond angle rotations represent random variables describing the chain configurations, so that when the number of these random variables, corresponding to the possible bond-vectors, gets large enough, the central limit theorem requires that the probability distribution for a vector between the ends of the chain must tend to a limiting Gaussian [1,2]. The random walk

[^1]: Figures in brackets indicate literature references.

configurations of an idealized polymer permit two different segments of the polymer to occupy the same position in space, something which is not possible for real polymer molecules. Hence, real polymers are described in terms of interacting random walks with an excluded volume interaction prohibiting the multiple occupancy of monomers at the same place in space.

Experimental methods of polymer characterization include *light scattering, osmometry, sedimentation, viscometry*, etc. When the polymer concentration $c$ approaches zero; osmometry provides the determination of the molecular weight $M$ of the polymer, while small angle light scattering yields the polymers' radius of gyration $R_G$. The limiting slope of the light scattering intensity as a function of concentration in the zero angle limit provides the polymer second virial coefficient $A_2$ which measures the effective volume that a polymer excludes to other chains. The translational diffusion coefficient $D$ is written for $c \rightarrow 0$ using Stokes' and Einstein's laws in terms of hydrodynamic radius $R_H$ by $D = kT/6\pi\eta_0 R_H$ where $T$ is the absolute temperature, $k$ is Boltzmann's constant and $\eta_0$ is the solvent viscosity. If $\eta$ designates the viscosity of the polymer solution, the intrinsic viscosity $[\eta] = \lim_{c \to 0}(\eta - \eta_0)/c\eta_0$ gives another measure of the volume occupied by a single polymer chain.

All of these large scale observables for polymers provide different measures of the overall size and shape of the polymer. These properties are often found to vary with the polymer molecular weight in the form of a power law $KM^a$ where the proportionality factors $K$ and the exponents $a$ depend on the polymer, solvent, and temperature as well being slowly varying functions of $M$. It is the goal of a comprehensive quantitative theory of polymers [3] in dilute solution to explain the variation of $K$ and $a$ with the polymer-solvent system and the temperature over the full experimentally accessible range.

Because large scale polymer properties like $R_G$, $A_2$, $D$ and ($\eta$) are measures of large scale or long wave length polymer properties, the theoretical description of these polymer properties does not require a detailed treatment of the short range microscopic details of the polymer such as the specific bond lengths, bond angles and hindered rotation potentials. Rather, *it suffices to employ apparently simple models to capture the essential large length scale characteristics of long chain molecules* [3-7].

The above noted popular random flight model of polymers treats the chain as having a set of effective monomer units sequentially number 0, 1, ..., at the spatial positions $r_0$, $r_1$, ..., $r_n$. Because the properties of interest involve large distance scales and implicitly large $n$, the central limit theorem allows us to take the individual bonds to have an effective Gaussian length distribution [1,2] with an rms value of $l$. The excluded volume interaction is modeled by introducing a short range repulsive contribution to the energy when a pair of segments occupies the same position in

space. This model then describes the dimensionless (free) energy associated with a chain configuration $\{r_k\}$ of the form [1-4]

$$H/kT = (d/2l^2) \sum_{j=1}^{n} |r_j - r_{j-1}|^2$$

$$+ (\beta_0/2) \sum_{i \neq j = 0}^{n} \delta(r_i - r_j) \qquad (1)$$

where $\beta_0$ is the volume excluded by an effective segment due to the presence of another. The spatial dimensionality is $d$, and it is convenient to consider the description of polymers as a function of the dimensionality of space where normal solutions involve $d = 3$ while the case $d = 2$ is associated with polymers at a surface or an interface.

The probability distribution function for the chain configuration $\{r_k\}$ is governed by the Boltzmann factor $\exp(-H/kT)$. Each of the monomer units in (1) describes the center of mass position of a collection of several actual monomers in the real polymer. This coarse graining is permissible because we are interested in long wavelength polymer properties. The first term on the right hand side in (1), therefore, represents the entropy of the polymer chain configuration associated with the many internal degrees of freedom in these coarse grain effective units [1-3]. Retention of only this term yields the simple Gaussian chain model of polymers at the theta temperature, a simple model for which the long wavelength polymer properties are easily evaluated. The entropic or elastic energy accounts in the model for the connectivity of the polymer chain.

The second term on the right hand side of (1) contains a pairwise sum over all the effective monomers and thereby converts the simple Gaussian chain model into a true many-body problem. The complexity of the treatment of excluded volume is readily seen by expanding $\exp(-H/kT)$ of (1) in powers of $\beta_0$ and evaluating polymer properties as a formal power series in $\beta_0$. This excluded volume perturbation theory is readily shown [2-7] to be an expansion in the dimensionless quantity $\beta_0 l^{-d} n^{\epsilon/2}$ with $\epsilon = 4 - d$. Hence, for high molecular weight polymers where $n$ is large, the expansion is in a large parameter as long as $\beta_0 l^{-d} < n^{-\epsilon/2}$ and $d < 4$. Consequently, the perturbation theory alone is of little use except very near the theta temperature where the empirical $\beta_0$ vanishes. The power law dependence of polymer properties on the molecular weight also indicates the difficulty of using these perturbation expansions in powers of $\beta_0$ since it is hard to see how a few terms in such an expansion, in a large parameter, can simply be resumed in order to provide the empirically observed nonanalytic power law dependence with fractional and often continuously varying exponents as a function of temperature.

Renormalization group methods are designed specifically to effect the resummation of such asymptotic expansions in a large expansion parameter [3-10]. The theories work by analytically continuing the mathematical description to continuous dimensionality d in order to exhibit the singularities of the perturbation theory as a function of dimensionality [4-10]. As in analytic function theory, the singularities govern the dominant properties of the functional dependence on excluded volume and chain length and thereby enable resummation of the asymptotic perturbation expansion.

The process begins by noting that the perturbation expansion becomes a controllable one by expanding in powers of $\beta_0$ and of $\epsilon$. This $\epsilon$ expansion method yields for instance

$$n^{\epsilon/2} = 1 + (\epsilon/2)\ln(n) + (\epsilon^2/8)\ln^2(n) + \ldots, \qquad (2)$$

so that when $n$ is large, a small $\epsilon$ can be chosen such that $\epsilon\ln(n)$ is a small expansion parameter. This procedure makes the perturbation expansion a mathematically well behaved method of computation, but it remains to show how perturbation expansions can be usefully applied to real $d=3$ or $d=2$ polymers. The renormalization group method accomplishes the analytic continuation to $d=3$ by focusing on the nature of the singularities of the perturbation expansion in powers of excluded volume for large scale polymer properties. The condition that the observable properties of polymers mathematically exist for $d=4$ leads to a resummation approach which is embodied in the renormalization group equation that summarizes the general analytic dependence of large length scale polymer properties on the molecular weight and excluded volume interaction.

This renormalization group equation has implicit within it a coarse graining length scale $L$ that plays several important roles. First, $L$ is a phenomenological parameter which is used to average out short length scale properties of the theoretical model which are irrelevant in the description of large length scale polymer properties. Comparison with experiment shows $L$ to characterize a correlation range along the chain for excluded volume interactions [4]. Roughly speaking, a portion of the chain of length $L$ interacts with the remainder of the polymer as if this portion were effectively a hard sphere. When the excluded volume is weak, $L$ is comparable to the size of the polymer, so there is effectively no excluded volume interaction. At the other extreme of strong excluded volume, $L$ approaches an asymptotic limit [11] which can be taken as a useful empirical definition of the step length $l$.

The theoretical analysis uses a continuum limit of the energy (1) in which the Gaussian chain model alone would lead to Wiener integrals [1] and where the excluded volume term must be appended with a short distance cut off to remove counting of self-excluded volume interactions [3-9]. The double perturbation expansions in $\beta_0$ and $\epsilon$ are fairly straightforward albeit extremely tedious for the polymer

properties of interest. The renormalized theory involves polymer properties on a coarse grain length scale $L$ in a region far from that for which perturbation expansions are valid. The physical properties are described in terms of the renormalized excluded volume $\beta$ and renormalized chain lengths $N$ replacing $\beta_0$ and the chain length $nl$ of the perturbation theory which is strictly valid only for $\beta_0$ very small. The details of the renormalization group method are too lengthy [3-10] to discuss here, so we turn to a description of some of the major results of the theory.

For instance, consider properties $Q$ for linear, ring, star and comb polymers like the radius of gyration or the hydrodynamic radius, etc., that naively scale as the $p$th power of the polymer radius. The Gaussian chain value for this property is written as $Q_0 = G_Q <S^2>_0^{p/2}$ with $<S^2>_0 = Nl/6$ the Gaussian chain square radius of gyration $R_G^2$. The property $Q_0$ is assumed to be known since it is relatively easily evaluated using the Gaussian chain model. An approximation to the second order renormalization group calculation in $d=3$ yields [7,11]

$$Q = \begin{cases} Q_0 (1 + 32\bar{z}/3)^{p/8}[1 + a_Q(32\bar{z}/3)/(1 + 32\bar{z}/3)], & \bar{z} < 0.75 \\ Q_0(6.441\bar{z})^{p(2\nu-1)}(1 + a_Q), & \bar{z} > 0.75 \ , \end{cases} \qquad (3)$$

where $\nu = 0.592$ to order $\epsilon^2$ and $a_Q$ is a pure number that depends only on the property $Q$ and emerges only from a first order calculation in $\epsilon$. The variable $\bar{z}$ is an empirical parameter which is often observed to depend on molecular weight and temperature in the form $\bar{z} = AM^{1/2}[1 - (\theta/T)]$ with $A$ a polymer and solvent dependent quantity and $\theta$ the theta temperature where $A_2$ vanishes.

Perturbation expansions in excluded volume are effectively expansions in a parameter like $\bar{z}$. Expansion of (3) in $\bar{z}$ shows the expansion coefficients to be growing rapidly with the power of $\bar{z}$. Hence, eq (3) represents the results of a rather sophisticated resummation of the perturbation expansion in powers of $\bar{z}$ based in fact on only the first two terms in the series and an analysis of the singularities of this perturbation expansion as a function of the dimensionality of space. A slightly different form of the prediction emerges for quantities like second virial coefficient which vanish at the theta point, and the reader is referred to our previous papers for these simple analytic formulas [7,11].

Figure 1 provides an example of the comparison between theory and experiment [12] for the interpenetration function $\psi(\bar{z})$ defined by

$$A_2 = (4\pi <S^2>)^{3/2}(N_A/2M^2)\psi(\bar{z}) \qquad (4)$$

as a function of the radius of gyration expansion factor defined by

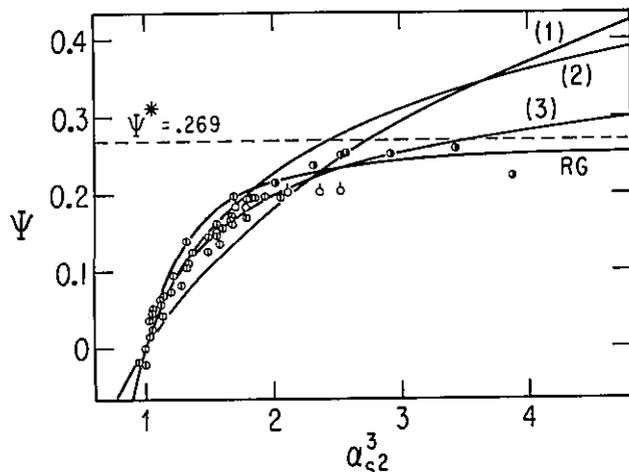$$\alpha_{S^2}^2 = <S^2>/<S^2>_0 \ . \qquad (5)$$

505

Figure 1—Comparison of theories [2] and experiment for variation of $\psi$ with $\alpha_{s2}^3$ [3]. The figure is reproduced from Yamakawa [2] with the parameter free RG predictions added. Data points from Norisuye et al. [14] are $\bullet$ for polychloroprene (PC) in $CCl_4$ at 25°C, $\circ$ for PC in $n$-butyl acetate at 25°C, and $\odot$ for PC in transdecalin at various temperatures. Similar data [15] and agreement with theory is available for polyisobutylene, polystryrene and poly-$p$-methylstyrene in various solvents. The curves (1)-(3) represent older theories as reviewed by Yamakawa [2].

The renormalization group prediction is given by the solid line marked RG, whereas the other lines represent older theories [2] which lack systematic mathematical guidance and therefore which are unable to predict the physically obvious fact that $\psi$ must lead to a universal value in the limit of large chain expansion or equivalently large excluded volume interaction $\bar{z}$.

An important feature of figure 1 is that the theoretical curve is obtained with absolutely no *adjustable parameters*, and that it is also derived using purely analytical methods. Hence, its derivation is in the true spirit of *analytical chemistry*, using the methods of *mathematical analysis* to provide *quantitative* descriptions of the properties of *chemical systems*. This type of work merges the disciplines of analytical chemistry and applied mathematics and is therefore at the heart of the goals of chemometrics.

Renormalization group calculations have been performed for a wide range of polymer properties [12] in dilute solutions for linear, ring, star, and block copolymers. The agreement between theory and experiment is generally as good or

better than that presented in figure 1. The situation is somewhat more complicated for the dynamical properties $R_H$ and $[\eta]$ where our theoretical calculations show that the effective exponents $a$ in good solvents often depend on an additional parameter, called the draining parameter [13]. We believe that the success of this renormalization group description of the excluded volume dependence of dilute solution properties will enable us to describe polymer properties in a variety of mathematically more complicated and physically very interesting situations such as the properties of polymers in interaction with a surface or interface, the properties of polyelectrolyte solutions where polymers have a distribution of charges and there are small counter ions in solution, and the properties of polymer mixtures in solution.

## References

[1] A more mathematically oriented introduction to polymers as interacting random walks is given in K.F. Freed, Ann. Prob. **9**, 537 (1981), a paper which predates our current work with the more powerful renormalization group methods.

[2] See, for instance, H. Yamakawa, *Modern Theory of Polymer Solutions* (Harper and Row, New York, 1971) and references therein.

[3] Freed, K.F., Accts. Chem. Res. **18**, 38 (1985).

[4] Oono, Y., T. Ohta and K.F. Freed, J. Chem. Phys. **74**, 6458 (1971).

[5] Oono, Y., and K.F. Freed, J. Phys. A **15**, 1931 (1982).

[6] Kholodenko A.L., and K.F. Freed, J. Chem. Phys. **78**, 7390 (1983).

[7] Douglas, J.F., and K.F. Freed, Macromolecules **17**, 1854 (1984).

[8] Amit, D.J., *Field Theory, The Renormalization Group and Critical Phenomena* (McGraw-Hill, New York, 1978).

[9] Ramond, *Field Theory, A Modern Primer* (Benjamin/Cummings, Reading, Mass., 1981).

[10] Freed, K.F., *Renormalization Group Theory of Macromolecules* (Wiley, to be published).

[11] Douglas, J.F., and K.F Freed, Macromolecules **15**, 1800 (1983).

[12] Douglas, J.F., and K.F. Freed, Macromolecules **17**, 2344 (1984); **18**, 201 (1985).

[13] Douglas, J.F., and K.F. Freed, Macromolecules **17**, 2354 (1984).

[14] Norisuye, T., K. Kawahara, A. Teramoto and H. Fujita, J. Chem. Phys. **49**, 4330 (1968); K. Kawahara, T. Norisuye and H. Fujita, J. Chem. Phys. **49**, 4339 (1968).

[15] Berry, G.C., J. Chem. Phys. **44**, 4550 (1966); G. Tanaka, S. Imai and H. Yamakawa, J. Chem. Phys. **52**, 2639 (1970); T. Matsumoto, N. Nishioka and H. Fujita, J. Polym. Sci. **10**, 23 (1972).

# Fourier Representations
# of Pdf's Arising in Crystallography

## George H. Weiss

### National Institutes of Health, Bethesda, MD 20205

and

## Uri Shmueli

### Tel-Aviv University, Ramat-Aviv 69978, Tel-Aviv, Israel

A survey is given of some recent calculations of univariate and multivariate probability density functions (pdf's) of structure factors used to interpret crystallographic data. We have found that in the presence of sufficient atomic heterogeneity the frequently used approximations derived from the central limit theorem in the form of Edgeworth or Gram-Charlier series can be quite unreliable, and in these cases the more exact, but lengthier, Fourier calculations must be made.

Key words: characteristic functions; direct methods of phase determination; Fourier series; intensity statistics.

Few scientific disciplines depend so heavily on techniques based on the central limit theorem and associated expansions in orthogonal polynomials as does crystallography. Ever since the pioneering work of Wilson [1,2],[1] and Karle and Hauptman [3-5], the central limit theorem has played a vital role in translating crystallographic scattering data into structural information and, indeed, it is built into many computer routines for this purpose. As we will show, when the central limit theorem is applied to data from unit cells with a considerable variation in the atomic weights of the constituent atoms it can lead to serious qualitative errors. That this

---

**About the Authors:** George H. Weiss is an applied mathematician and Uri Shmueli is a chemist.

---

[1] Figures in brackets indicate literature references.

is true is well known to crystallographers who have made heavy use of Edgeworth and related expansions to correct zeroth order approximations based on the central limit theorem [5,6]. It is not generally appreciated, however, that serious errors can persist even with these correction terms, provided that atomic heterogeneity is sufficiently great. This suggests the value that may be attached to exact results when these are available and are readily computed. This paper reports on recent efforts we and several collaborators [7-11] have made in this direction.

Two general classes of probabilistic methods are used to deduce structural information from radiation intensity diffracted from crystals, the so-called intensity statistics and direct methods of phase determination. In order to make this exposition self-contained, we will sketch how such information can be derived from data

on intensities in a particularly simple case, and refer the interested reader to two monographs that give detailed accounts of these subjects in more general cases [12,13]. The arrangement of atoms in a unit cell of a crystal is most often restricted by the space group to which the crystal belongs [14], and in the general case, only the arrangement within the asymmetric part of the cell needs to be determined. The intensity of the diffracted radiation can be represented in terms of structure factors $F(\mathbf{h})$, where the vector $\mathbf{h}$ and its components $(h,k,l)$, the orders of diffraction, specify the geometric relation between incident and scattered beams and their relative orientation to the basis vectors of the lattice of the diffracting crystal [14]. The structure factors are Fourier coefficients of the (periodic) density function of the scattering matter, and both their magnitude and phase are required in order to reconstruct the density—i.e., the actual atomic arrangement. Thus, $F(\mathbf{h})$ is in general a complex quantity, which we write as $F(\mathbf{h})=A(\mathbf{h})+iB(\mathbf{h})$. The function $F(\mathbf{h})$ can be expressed as a sum of contributions from individual atoms in the unit cell as

$$F(\mathbf{h})=\sum_j f_j \exp(2\pi i \mathbf{h}_j \mathbf{r}_j)=\sum_j f_j \exp(i\theta_j) \qquad (1)$$

where $\mathbf{r}_j$ is the location of atom $j$, the $f_j$ are so-called scattering or form factors which can be approximated, in the normalized-structure-factor representation (see below), by the atomic numbers of the corresponding atoms, and $\theta_j=2\pi\mathbf{h}\cdot\mathbf{r}_j$. The space of $\mathbf{h}$ is surveyed by varying the orientation of the crystal with respect to the incident beam.

Since $F(\mathbf{h})$ is a complex quantity, it can be represented as a vector in a plane which is the sum of $n$ vectors, the $j$'th being $f_j \exp(i\theta_j)$. The fundamental difficulty faced by crystallographers is that only the magnitude $|F|$ is measurable (although some recent work may change this situation [15]), and the phase of $F(\mathbf{h})$ must be inferred indirectly. To do so, one can establish a correspondence between the vector $F$ and a random walk first studied by Pearson [16]. Using theorem of Weyl [17], one can show that if the components of $r$ are rationally independent, i.e., there exists no vector of integers $m$ such that $\mathbf{m}\cdot\mathbf{r}=$ integer, then the set of angles, $\{\theta_j\}$, can be regarded as consisting of independent random variables, each of which is uniformly distributed over the interval $(0,1)$[17]. Thus the properties of the $F(\mathbf{h})$ can be determined using probabilistic methods, as was first pointed out by Wilson [1,2].

For a typical and important case in which probabilistic techniques allow one to derive structural information, consider how one can distinguish between centrosymmetric and noncentrosymmetric (space groups $P\bar{1}$ and $P1$, respectively) unit cells on the basis of intensity statistics alone. A centrosymmetric unit cell is one in which, for every atom located at $\mathbf{r}_j$, there is an identical one at $-\mathbf{r}_j$. Consequently if we write $F=A+iB$ where

$$A=\sum_j f_j \cos\theta_j, \quad B=\sum_j f_j \sin\theta_j \qquad (2)$$

it follows that $B\equiv 0$ by symmetry in the presence of centrosymmetry. When the unit cell is noncentrosymmetric $B$ is not necessarily equal to 0. Hence the value of $F$ can be represented as a one-dimensional random walk in $P\bar{1}$ and by a two-dimensional random walk in $P1$. In what follows we will use the physics notation that "$<\ >$" denotes the average of the variable contained in brackets. It will also prove convenient to work with the normalized structure factor $E=F/<|F|^2>^{\frac{1}{2}}$ which, since $<F>=0$, has the property that $<|E^2|>=1$. Wilson's argument uses the central limit theorem to deduce the pdf of scattered intensities. In $P\bar{1}$, for which $B\equiv 0$, the form of the pdf of $E$ that follows from the central limit theorem is

$$p(|E|)=(2/\pi)^{\frac{1}{2}} \exp(-E^2/2). \qquad (3)$$

The corresponding pdf for the two dimensional case for unit cells without a crystallographic center of symmetry is

$$p(|E|)=2|E|\exp(-|E|^2). \qquad (4)$$

The qualitative difference between eqs (3) and (4) thus allows the experimental distinction to be drawn purely on a comparison of intensity data with the two forms for the pdf.

Notice, however, that the use of the central limit theorem presupposes the validity of certain assumptions, the major one of which is the presence of a large number of atoms in the unit cell and the second of which is that the $f_j$ appearing in eq (1) should not exhibit too great a heterogeneity. The first of these assumptions holds for most crystalline materials of interest, but the second may be violated particularly when there are a small number of atoms that are considerably heavier than the majority of atoms comprising the molecule. When this is the case it is customary to replace, e.g., eq (3) by the Edgeworth series

$$p(|E|)=(2/\pi)^{\frac{1}{2}} \exp(-E^2/2)\{1+\sum_n a_n H_n(|E|/\sqrt{2})\} \quad (5)$$

where the $n$'th coefficient, $a_n$ is expressible as a linear combination of the moments of $A$ in eq (2) and $H_n(x)$ is the $n$'th Hermite polynomial. These are readily calculated for the simpler space groups [18], and all space-group results are available for fourth, sixth, and eighth

moments[19,20]. Furthermore, the Edgeworth expansion may also not be too useful in the presence of extreme heterogeneity. This is illustrated in figure 1 in which the asymmetric unit of a cell in $P\bar{1}$ consists of 14 carbon atoms and one uranium atom, with a ratio of $f$'s approximately equal to $15\frac{1}{3}$. With 0 or 2 moments the Edgeworth series fails to reproduce the maximum and the 4 and 8 moment approximation locates the maximum quite far from its actual position. It is therefore desirable to have an exact easily computable representation for the pdf which is robust with respect to changes in atomic heterogeneity.

Just such a representation was first suggested by Barakat in a study of the freely jointed chain as a model for polymer configurations [21] and of laser speckle [22]. Let us write $g_j = f_j / (\Sigma_j f_j^2)^{\frac{1}{2}}$ so that
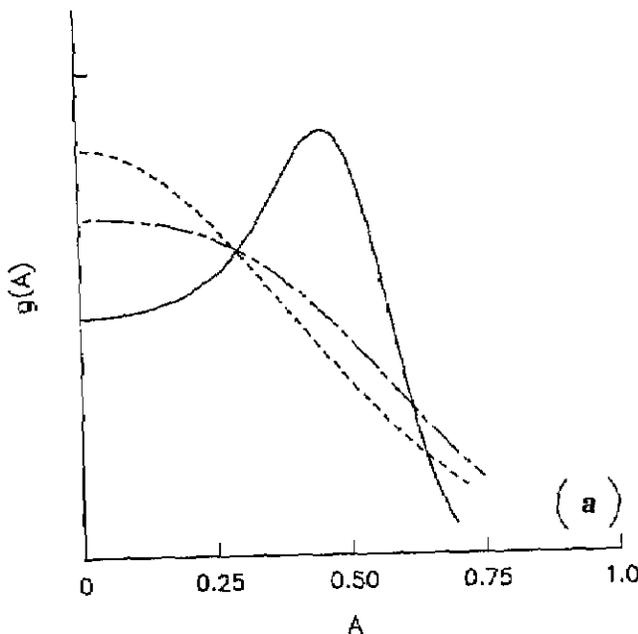
$$E = \Sigma_j g_j \exp(i\theta_j) = A + iB \tag{6}$$

and let us set

$$S = \Sigma_j g_j \tag{7}$$

so that $-S \leqslant A, B \leqslant S$. As an example we consider the case of a centrosymmetric unit cell for which $B = 0$. The pdf of $A$, $g(A)$, has the property that it can differ from zero only in the interval $S^2 \geqslant A^2$. Within this interval we will expand $g(A)$ in a Fourier series:

$$g(A) = \frac{1}{2S} \left\{ 1 + 2 \sum_{n=1}^{\infty} a_n \cos\left(\frac{\pi m A}{S}\right) \right\} \tag{8}$$

where

$$a_n = \int_S^S g(A) \cos(\frac{\pi m A}{S}) \alpha A = \int_{-\infty}^{\infty} g(A) \cos(\frac{\pi m A}{S}) \alpha A$$

$$= C(\frac{\pi m}{S}) \tag{9}$$

where $C(\omega)$ is the characteristic function generated by $g(A)$. The Fourier series in eq (8) corresponds to a sampling theorem [23] for pdf's with a compact support. When the unit cell is noncentrosymmetric so that $B = 0$ in general, it is more convenient to expand the pdf of $|E| = (A^2 + B^2)^{\frac{1}{2}}$ in a Fourier-Bessel function series

$$p(|E|) = \frac{2|E|}{S^2} \sum_{i=0}^{\infty} D_j J_0(\gamma_j |E|/S) \tag{10}$$

where the $\gamma_j$ are successive roots of $J_0(\gamma) = 0$ and the coefficients, $D_j$, are

$$D_j = C(\gamma_j/S)/J_1^2(\gamma_j) \tag{11}$$

again written in terms of the characteristic function.

Two questions that require an answer relate to the advantage of representations such as those in eqs (8) or (10) and the feasibility of numerical evaluation of the series. In the absence of atomic heterogeneity, or when there is a very large number of atoms in a unit cell, the central limit results are perfectly adequate for crys-
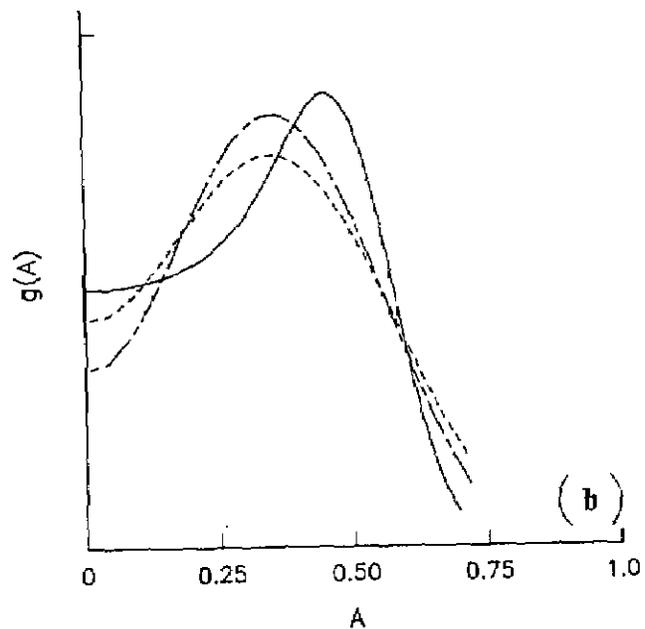


Figure 1a-(a) Approximations to the exact pdf $g(A)$ (devoted by the solid line) for a unit cell in spacea group $P\bar{1}$, consisting of 14 carbon atoms and a single uranium atom (atomic weight ratio 15 $\pm$1) in the asymmetric unit. For convenience $A_{max}$ has been set equal to 1. Note that the pdf is symmetric around $A = 0$. The approximations are a Gaussian (---) and the Gaussian corrected by two moments (-----). (b) Approximations to the same pdf as in figure 1a by an Edgeworth series using 4 moments (---) and 8 moments (-----).

tallographic applications. However, when there are fewer than about 40 atoms in the unit cell combined with *one or two outstandingly heavy atoms*, one has to convolve the Gaussian with the appropriate pdf for the heavy atoms [12]. In principle, using the series of eqs (8) and (10) finesses this difficulty, provided that the convergence properties are not overwhelming. In practice, *in the case of intensity statistics we have found no problem in evaluating the Fourier or Fourier-Bessel function series*, requiring no more than about 40 terms for the most extreme amounts of heterogeneity, and many fewer terms in the absence of heterogeneity. The evaluation *of the analogous series for direct methods can* present much tougher numerical problems, as we will see. Finally, a problem not so far discussed is the ease with which expressions for the characteristic function can be calculated. We have found that it is not too difficult to evaluate the characteristic function for all but a handful of space groups whose structure factor is found in the International Tables [24]. As an example, let us write the structure factor for $P\bar{1}$ as

$$A = 2 \sum_{j=1}^{n/2} g_j \cos\theta_j \qquad (12)$$

where the $g_j$ are known and the $\theta_j$ are uniformly distributed in $(0,2\pi)$. The characteristic function is

$$C(\omega) = \langle \exp(2i\omega \sum_{j=1}^{n/2} \cos\theta_j) \rangle \qquad (13)$$

$$= \prod_{j=1}^{n/2} C_j(\omega)$$

where

$$C_j(\omega) = \langle \exp(2i\omega g_j \cos\theta) \rangle$$

$$= \frac{1}{2\pi} \int_{\pi}^{\pi} \exp(2i\omega g_j \cos\theta) a\theta = J_0(2\omega g_j). \qquad (14)$$

Other examples merely test one's ability to evaluate integrals. For example, we have recently examined the Fourier representation of the pdf of the intensity for a unit cell in $P\bar{1}$ in which there is an auxiliary or noncrystallographic center of symmetry located at **d**, so that a single atom located at $r_j$ generates one at $-r_j$ and $-r_j \pm 2\mathbf{d}$ [9,10]. In this case one can show that

$$A(\mathbf{h}) = 4 \sum_{j=1}^{n/4} g_j \cos(2\pi\mathbf{h}\cdot\mathbf{d})\cos[2\pi\mathbf{h}\cdot(r_j - \mathbf{d})] \qquad (15)$$

and the corresponding characteristic function is

$$C(\omega) = \frac{2}{\pi} \int_0^{\pi/2} (\prod_{j=1}^{n/4} J_0(4\omega g_j \cos\theta)) d\theta. \qquad (16)$$

Again the numerical problems associated with this representation were not severe and allowed us to generalize the theory first presented by Rogers and Wilson for the equal-atom case [9,25]. It is possible though algebraically messy to generate the orthogonal polynomials corresponding to the Rogers-Wilson pdf, but the Fourier series representation is relatively straightforward. One can also analyze *partially bicentric structures using the same techniques* [10].

Our present development of Fourier representations of crystallographic pdf's has led us into the examination of direct methods in which one is interested in the joint pdf *of several, usually correlated, structure factors. One* of the simplest examples of these is the so-called $\Sigma_1$ relationship [13,26], in which one uses the joint pdf of $E(\mathbf{h})$ and $E(2\mathbf{h})$ to determine the probability that the phase of $E(2\mathbf{h})$ is positive given a knowledge of $E(\mathbf{h})|$ and $|E(2\mathbf{h})|$. *For simplicity we consider structures in* $P\bar{1}$ letting $E(\mathbf{h}) = E$ and $E(2\mathbf{h}) = G$. Then

$$p(E,G) = \frac{1}{4S^2} \sum_{r=-\infty}^{\infty} \sum_{s=-\infty}^{\infty} C\left(\frac{\pi r}{S}, \frac{\pi s}{S}\right) \\ \cos\left(\frac{\pi r E}{S}\right)\cos\left(\frac{\pi s G}{S}\right) \qquad (17)$$

where

$$C(\omega_1,\omega_2) = \langle \exp(i\omega_1 E + i\omega_2 G) \rangle = \prod_{j=1}^{n/2} C_j(\omega_1,\omega_2) \qquad (18)$$

in which since $$G = 2 \sum_{j=1}^{n/2} g_j \cos(2\theta_j)$$

$$C_j(\omega_1,\omega_2) = \frac{1}{2\pi} \int_{-\pi}^{\pi} \exp[2ig_j(\omega_1\cos\theta + \omega_2\cos2\theta)] d\theta \\ \equiv R_j + iI_j \qquad (19)$$

where $R_j$ and $I_j$ can be expanded in terms of Bessel functions as

$$R_j(\omega_1,\omega_2) = J_0(2g_j\omega_1)J_0(2g_j\omega_2) + 2 \sum_{m=1}^{\infty} \\ (-1)^m J_{4m}(2g_j\omega_1)J_{2m}(2g_j\omega_2) \qquad (20)$$

$$I_j(\omega_1,\omega_2) = 2 \sum_{m=0}^{\infty} (-1)^{m+1} J_{4m+2}(2g_j\omega_1)J_{2m+1}(2g_j\omega_2).$$

From eq (18) it follows that $C(\omega_1,\omega_2) = R(\omega_1,\omega_2) + iI(\omega_1,\omega_2)$ where $R$ and $I$ can be computed from the $R_j$ and $I_j$. The probability that $G$ is positive given $E$ can now be written exactly as

510

$$p_+(2\mathbf{h}1\mathbf{h})=\tfrac{1}{2}(1+\tfrac{\Omega}{\Gamma})\tag{21}$$

where

$$\Omega=\frac{1}{S^2}\sum_{s=1}^{\infty}\sum_{t=1}^{\infty}I\left(\frac{\pi s}{S},\frac{\pi t}{S}\right)\cos\left(\frac{\pi sE}{S}\right)\sin\left(\frac{\pi t\,|G|}{S}\right)$$

$$\Gamma=\frac{1}{4S^2}\left\{1+2\sum_{S=1}^{\infty}R\left(\frac{\pi s}{S},0\right)\left[\cos\left(\frac{\pi sE}{S}\right)+\cos\left(\frac{\pi tG}{S}\right)\right]\right.$$

$$\left.+4\sum_{s=1}^{\infty}\sum_{t=1}^{\infty}R\left(\frac{\pi s}{S},\frac{\pi t}{S}\right)\cos\left(\frac{\pi sE}{S}\right)\cos\left(\frac{\pi tG}{S}\right)\right\}.\tag{22}$$

The exact eq (21) should be compared to the much simpler approximation furnished by the use of the central limit theorem [26],

$$p_+(2\mathbf{h}1\mathbf{h})\sim\tfrac{1}{2}\left[1+\tanh\left(\frac{\sigma_2}{2\sigma_2^{\frac{3}{2}}}\,|G|(E^2-1)\right)\right]\tag{23}$$

where

$$\sigma_m=2\sum_{j=1}^{n/2}g_j^m\,.\tag{24}$$

Although eq (23) and generalizations of it are much used in the crystallographic literature, there has been no real test of its accuracy in the presence of heterogeneity, since until now there has been no attempt to calculate the exact pdf. A comparison of the result of evaluating eq (21) with that obtained from eq (23) for an assumed composition $C_{30}Kr_2$ in the half unit of a $P\bar{1}$ structure is shown in figure 2 for $G=1.75$ [11]. A substantial difference between the two predictions is immediately evident. Further evidence of the inaccuracy of eq (23) in the presence of atomic heterogeneity is provided in figure 3 where we examine the effects of the variation in atomic weights for a unit cell in which the half unit is $C_{30}X_2$, where $X$ varies. In the absence of heterogeneity eq (23) provides perfectly satisfactory results, but its utility decreases considerably with an increase in the atomic weight of the X atom.

We are presently examining the analogous properties of the $\Sigma_2$ relationship, in which one determines probabilistic relations between phases from properties of the joint pdf of $E(\mathbf{h})$, $E(\mathbf{k})$ and $E(-\mathbf{h}-\mathbf{k})$, which requires the evaluation of higher order Fourier series by the same basic techniques. While this investigation is very similar both in spirit and results to those for $\Sigma_1$ discussed in the last paragraphs, it appears to be much more difficult to evaluate the series for the pdf of the three-phase
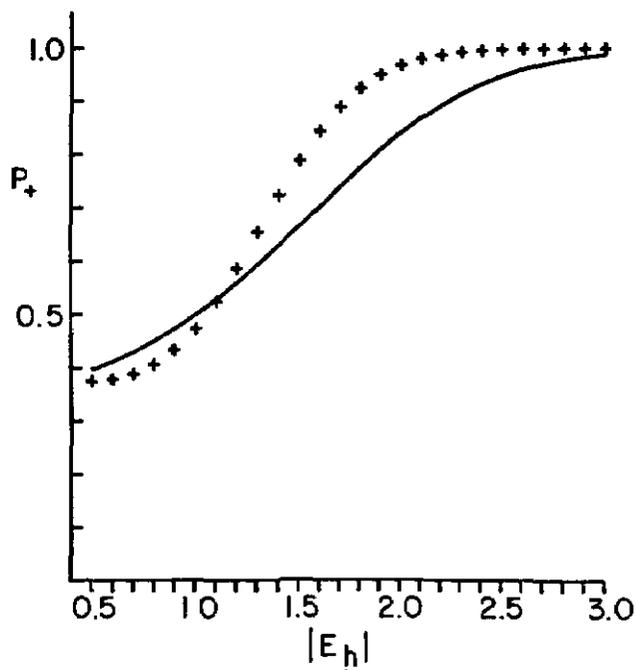


Figure 2–A graph of the exact expression for $p_+(2\mathbf{h}\,|\mathbf{h})$, the probability that the phase of $E(2\mathbf{h})$ is positive, as a function of $E\,(+++)$ compared to the approximation provided by eq (23) (the solid curve) for a molecule with the assumed composition $C_{30}Kr_2$ in the asymmetric unit. The magnitude, $|E(2\mathbf{h})|$, was chosen equal to 1.75 for this example.

invariant, $\Phi$, defined in terms of the phases of the triplet of structure factors $E(\mathbf{h})$, $E(\mathbf{k})$, $E(-\mathbf{h}-\mathbf{k})$, by

$$\Phi=\phi(\mathbf{h})+\phi(\mathbf{k})+\phi(-\mathbf{h}-\mathbf{k}).\tag{25}$$

To convey some notion of the difficulties we point out that the characteristic function to be evaluated is

$$C_j(\mathbf{w})=\langle\exp\{ig_j(\omega_1A_1+\omega_2B_1+\omega_3A_2+\omega_4B_2+w_5A_3$$
$$+\omega_6B_3)\}\rangle\tag{26}$$

where $E(\mathbf{h})=A_1+iB_1$, $E(\mathbf{k})=A_2+iB_2$, $E(-\mathbf{h}-\mathbf{k})=A_3+iB_3$. A detailed evaluation of $C_j(\mathbf{w})$ results in the expression

$$R_j(\mathbf{w})=\prod_{k=1}^{6}J_0(f_j\omega_k)+2\sum_{m=1}^{\infty}(-1)^m\prod_{k=1}^{6}J_{2m}(f_j\omega_k)$$

$$I_j(\mathbf{w})=2\sum_{m=0}^{\infty}(-1)^{m+1}\prod_{k=1}^{6}J_{2m+1}(f_j\omega_k)$$

$$C(\mathbf{w})=R(\mathbf{w})+iI(\mathbf{w})=\prod_j C_j(\mathbf{w})=\prod_j(R_j+iI_j)\,.\tag{27}$$

The resulting expression for the pdf of $\Phi$ is in terms of sevenfold Fourier series, each coefficient of which is an
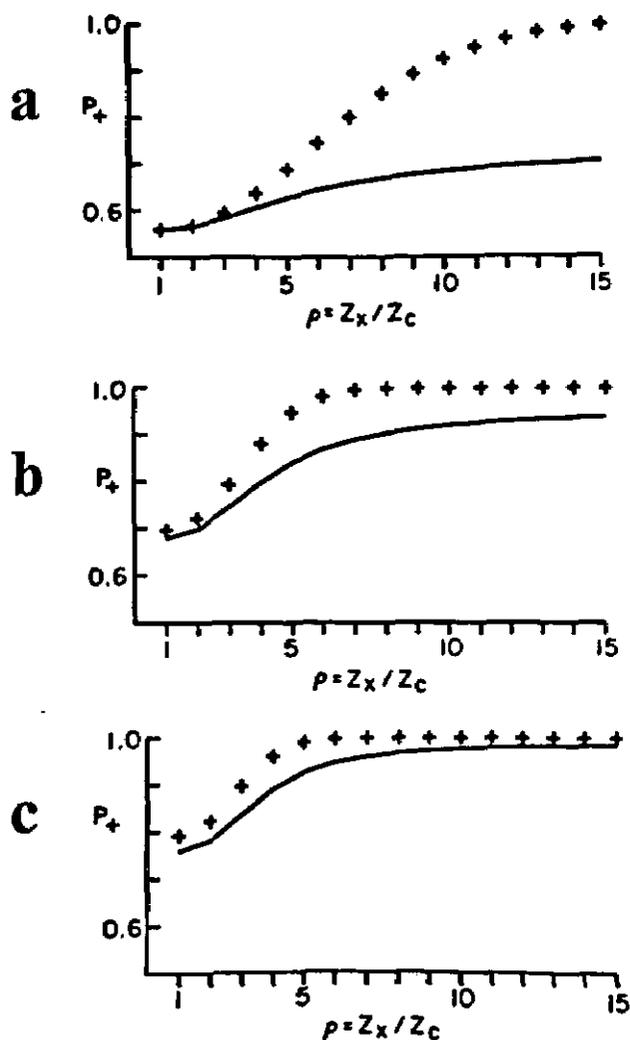
511

Figure 3–This figure shows the effects of heterogeneity on $p + (2h \mid h)$, as a function of the ratio of atomic numbers, $p = Z_x/Z_c$ for a molecule with the composition $C_{30}X_2$. The $+++$'s are the exact results and the solid lines are the approximation of eq (23). The values chosen are (a) $|E(h)| = |E(2h)| = 1.5$, (b) $|E(h)| = |E(2h)| = 2.0$, (c) $|E(h)| = |E(2h)| = 2.25$. Note that the approximation is always on the conservative side. It is not known whether this is always true.

infinite series of the form shown in the last equation. Whether the resulting calculations can be made in a reasonable amount of time remains to be seen, but the difficulties to be overcome are exemplified by this problem.

A final word is in order about the philosophy behind the series of projects that we have undertaken. It would hardly be sensible to want to eliminate methods based on the central limit theorem that have served crystallographers so well in the past. However, it is useful to establish the limitations of these methods by having more exact representations available. Indeed we have explored such limitations in the case of tests based on the

$\Sigma_1$ relationship as indicated earlier and are presently considering more complicated crystallographic techniques. Furthermore, as the processing of crystallographic data becomes more and more automated it becomes increasingly attractive to have exact, rather than approximate formulae in the computer. We hope, in the coming years, to explore the feasibility of doing this for a variety of techniques, as well as contributing to the development of further ones based on the availability of exact representations.

## References

[1] Wilson, A. J. C. Determination of absolute from relative x-ray intensity data, Nature 150, 151–152 (1942).

[2] Wilson, A. J. C. The probability distribution of x-ray intensities, Acta Crystallogr. 2, 318–321 (1949).

[3] Hauptman, H., and J. Karle, Crystal structure determination by means of a statistical distribution of interatomic vectors, Acta Crystallogr. 5, 48–54 (1952).

[4] Hauptman, H., and J. Karle, The probability distribution of the magnitude of a structure factor, I. The centrosymmetric crystal, Acta Crystallogr. 6, 131–135 (1953); II. The noncentro-symmetric crystal, ibid 6, 136–141.

[5] Hauptman, H., and J. Karle, Solution of the Phase Problem 1. The Centrosymmetric Crystal, ACA Monograph No. 3 (Polycrystal Book Service, New York) 1953.

[6] Klug, A., Joint probability distributions of structure factors and the phase problem, Acta Crystallogr. 11, 515–543 (1958).

[7] Weiss, G. H., and J. E. Kiefer, The Pearson random walk with unequal step sizes, J. Phys. A16, 489–495 (1983).

[8] Shmueli, U.; G. H. Weiss, J. E. Kiefer, and A. J. C. Wilson, Exact random walk models in crystallographic statistics. I. Space groups P1̄ and P1. Acta Crystallogr. A40, 651–660 (1984).

[9] Shmueli, U.; G. H. Weiss and J. E. Kiefer, Exact random walk models in crystallographic statistics, II. The bicentric distribution for the space group P1̄, Acta Crystallogr. A41, 55–59 (1985).

[10] Shmueli, U., and G. H. Weiss, Centric, bicentric, and partially bicentric intensity statistics in Structure and Statistics in Crystallography (Adenine Press, Quilderland New York) to appear.

[11] Shmueli, U., and G. H. Weiss, Exact joint probability distribution for centrosymmetric structure factors. Derivation and application to the $\Sigma_1$ relationship in the space group P1̄, Acta Crystallogr. (to appear).

[12] Srinivisan, R. and Parthasarathy, S., Some Statistical Applications in x-ray Crystallography, Pergamon Press, Oxford (1976).

[13] Giacovazzo, C., Direct Methods in Crystallography, Academic Press, New York (1980).

[14] Woolfson, M. M., An Introduction to x-ray Crystallography, Cambridge University Press, Cambridge (1970).

[15] Post B.; J. Nicolosi and J. Ladell, Experimental procedures for the determination of invariant phases of centrosymmetric crystals, Acta Crystallogr. A40, 684–688 (1984).

[16] Pearson, K. The problem of random walk, Nature 72, 294 (1905).

[17] Hardy, G. H., and E. M. Wright, Introduction to the Theory of Numbers, Oxford University Press, Oxford, 5th ed. (1979).

512

[18] Foster, F., and A. Hargreaves, The use of moments of x-ray intensity in space group determination I. Derivation of theoretical moments, Acta Crystallogr. 16, 1124–1133 (1963); II. Practical application, *ibid* 16, 1133–1138 (1963).

[19] Shmueli, U., and U. Kaldo, Calculation of even moments of the trigonometric structure factor. Methods and results, Acta Crystallogr. A37, 76–80 (1981).

[20] Shmueli, U., and U. Kaldor, Moments of the trigonometric structure factor, Acta Crystallogr. A39, 619–621 (1983).

[21] Barakat, R., Isotropic random flights, J. Phys. A6, 796–804 (1973).

[22] Barakat, R., First-order statistics of combined random si-

nusoidal waves with applications to laser speckle patterns, Optica Acta 21, 921 (1974).

[23] Hamming, R., *Numerical Methods for Scientists and Engineers,* (McGraw-Hill, New York, 2d ed. (1973).

[24] *International Tables for x-ray Crystallography* ed. N.F.M. Henry, K. Lonsdale, Kynoch Press, Birmingham, 2d ed. (1965).

[25] Rogers, D., and A. J. C. Wilson, The probability distribution of x-ray intensities. V. A note on some hypersymmetric distributions, Acta Crystallogr. 6, 439–449 (1953).

[26] Cochran, W., and M. M. Woolfson, The theory of sign relations between structure factors. Acta Crystallogr. 8, 1–12 (1955).

# DISCUSSION

of the Weiss-Shmueli paper,
Fourier Representations of Pdf's
Arising in Crystallography

## E. Prince

Institute for Materials Science and Engineering
National Bureau of Standards

This interesting paper by Drs. Weiss and Shmueli represents a substantially exact solution of a problem that has concerned crystallographers for more than 35 years, the analysis in terms of atomic structure of x-ray diffraction data. (Similar information can be obtained from the diffraction of electrons and neutrons, but, for reasons that are both experimental and theoretical, this information is mainly used to supplement that obtained from x-ray diffraction, which remains the basic tool of the structural crystallographer.) The observed intensity in x-ray diffraction is given by

$$I = SL \, |F(\mathbf{h})|^2,$$

where $S$ is a scale factor, $L$ is a geometrical factor, and $F(\mathbf{h})$, commonly called the structure factor, is the Fou-

rier transform of the electron density in a crystal. It may be written in the form

$$F(\mathbf{h}) = \int \rho(\mathbf{r}) \exp(2\pi i \mathbf{h} \cdot \mathbf{r}) d\mathbf{r}.$$

The density function, $\rho(\mathbf{r})$, in a crystal is periodic in three dimensions, so that it can be represented as a convolution of a function consisting of $\delta$ functions located at the nodes of a space lattice and a density defined in a small region known as a unit cell. Because of the periodicity the Fourier transform has appreciable values only at the nodes of a lattice in transform space, called by crystallographers the reciprocal lattice. Because it is a physical quantity, $\rho(\mathbf{r})$ is non-negative, and, furthermore, because a crystal is composed of atoms, it can be

513

represented as a sum of functions that are, at least to a good approximation, spherically symmetric about a finite set of nuclear positions. If we designate by $f_j$ the Fourier transforms of these atomic functions, the Fourier transform of the crystal can be written

$$F(\mathbf{h}) = \sum_{j=1}^{n} f_j \exp(2\pi i \mathbf{h} \cdot \mathbf{r}_j),$$

which is the conventional structure factor formula.

If it were possible to measure the values of the structure factor throughout transform space, it would be possible to compute the inverse transform and determine the density function, $\rho(\mathbf{r})$, directly. However, $F(\mathbf{h})$ is, in general, a complex quantity, and, because it appears in the intensity formula only as $|F(\mathbf{h})|^2$, only its amplitude can be measured, and that only within a finite region of transform space. In the early days of structural crystallography this "phase problem" was treated by using chemical intuition to devise a trial model for which $F$ could be calculated and then, using the calculated phases along with the observed amplitudes, to compute a density map. If the crystallographer was lucky, this map would show a sufficiently clear picture of the structure to suggest adjustments to the model, and several iterations of the process would converge to a structure that made chemical sense. In the days before high-speed, digital computers these computations were done on mechanical desk calculators using tabulated sines and cosines written on strips of cardboard known, after the two British crystallographers who introduced them, as Beevers-Lipson strips. The process was very laborious, and a single structure analysis could consume many months.

In studies, beginning in the 1940s, of the structures of boron hydrides chemical information to suggest a reasonable starting model was often not available. Even if the crystal possessed a center of symmetry, so that the imaginary parts of the contributions to $F$ from pairs of atoms would cancel, thereby constraining $F$ to have real values, the number of possible combinations of signs in the density summation could be enormous. It was realized, however, by Harker and Kasper [1][1] that many of the sign combinations would result in violations of the non-negativity condition on the electron density, and they were able to derive a number of inequality conditions that must be satisfied by certain combinations of $F$ values in order to keep the density positive. The use of the Harker-Kasper inequalities in the solution of the structure of decaborane, $B_{10}H_{14}$, in 1950 by Kasper, Lucht, and Harker [2] was the first successful application of direct methods to the determination of a crystal structure using diffraction data alone.

The Harker-Kasper inequalities were applicable only to centrosymmetric crystals, with their resultant real values of the structure factors. In view of the amount of labor involved in the computation of a density map it was certainly important to be able to determine whether a crystal did in fact have a center of symmetry. In the late 1940s this problem was attacked by Wilson and his coworkers [3]. In the limit of a large number of identical atoms distributed at random in the unit cell, the contributions of the individual atoms to $F$ are random walk steps, and the central limit theorem can be invoked to show that the distribution of $F$ is approximately normal with zero mean. $|F|^2$ is then distributed as $\chi^2$ with one degree of freedom if the crystal has a center of symmetry and with two degrees of freedom otherwise. The presence of other symmetry operations, such as rotation axes or mirror planes, constrains certain subsets of the structure factors to be real, so that statistical tests on the observed intensities can be an aid to determining the proper symmetry group for a crystal.

The Harker-Kasper inequalities may be viewed as a limiting case of a more general problem, which may be stated as follows: Given the magnitudes of a set of structure factors and the phases of a subset of them, (It is always possible to assign the phases of three structure factors arbitrarily. This merely defines the origin.) what are the probability density functions for the phases of others? Harker and Kasper identified particular cases where a discrete phase could be assigned with unit probability. The more general problem was attacked by Hauptman and Karle[2] in a long series of papers, beginning in the early 1950s [4], in which they have developed increasingly powerful methods for defining narrow ranges within which phases are likely to lie with high probability.

Most of the statistical methods that have been developed for determining a structural model are based on assumptions similar to those used by Wilson (1949), namely, that the crystal was composed of a large number of nearly identical atoms located at random within the unit cell. There is, however, another limit in which the solution of the phase problem is well known. This is the case (such as a simple metal) where the unit cell contains only one atom. In this limit all structure factors are identical in both magnitude and phase. Structure studies of very large molecules, such as proteins, have depended heavily on the preparation of crystals in which the unit cell contains one heavy atom, or a few at most, along with the very large number of atoms of carbon, nitrogen, and oxygen. The distribution of in-

[1] Figures in brackets indicate literature references.

tensities in this case cannot be similar to that in either limiting case, but must rather represent some sort of intermediate situation. Crystallographers have used various approximate methods to treat these real situations, and the results are strongly dependent on the validity of the approximations. The results given here by Weiss and Shmueli provide an accurate solution to which the approximate methods may be compared, and for this reason they are of tremendous interest to the crystallographic community.

## References

[1]  Harker, D., and J. S. Kasper, Acta Cryst. 1, 70–75 (1948).
[2]  Kasper, J. S.; C. M. Lucht and D. Harker, Acta Cryst. 3, 436–455 (1950).
[3]  Wilson, A. J. C., Acta Cryst. 2, 318–321 (1949).
[4]  Hauptman, H., and J. Karle, *Solution of the Phase Problem 1. The Centrosymmetric Crystal*, ACA Monograph No. 3, New York, Polycrystal Book Service (1953).

515

# Aggregated Markov Processes and Channel Gating Kinetics

## Donald R. Fredkin and John A. Rice

### University of California, San Diego; La Jolla, CA 92093

A finite state Markov process is aggregated into several groups. Rather than observing the underlying Markov process, one is only able to observe the aggregated process. What can be learned about the underlying process from the aggregated one? Such questions arise in the study of gating mechanisms in ion channels in muscle and nerve cell membranes. We discuss some recent results and their implications.

Key words: aggregated Markov process; channel gating kinetics.

## 1. Introduction

We are concerned with the following mathematical problem: a finite state Markov process in continuous time is aggregated, by which we mean that its states are grouped into a smaller number of aggregates. The Markov process is assumed to be in equilibrium. One is not able to observe the process itself, but only what aggregate the process is in as time goes along. From the aggregated process we wish to draw inferences about the underlying Markov process. How much can we learn? This is a rather general question and more sharply defined questions can be asked. For example, to what extent is the graph structure of the Markov process identifiable? Are some aspects of it identifiable and some not? (By the graph we mean a diagram showing the states and the interconnections between them, but not the numerical values of the transition rates.) Can some graph structures be ruled out as incompatible with the aggregated process? If

About the Authors: Donald R. Fredkin is with the Department of Physics and John A. Rice with the Department of Mathematics at the University of California, San Diego : La Jolla.

the graph structure is known or hypothesized *a priori*, what functionals of the various rate constants can be identified?

These are questions of *identifiability*. There are also problems related to the efficient statistical use of observation of the aggregated process over a finite time interval.

Problems such as these arise in modeling and data analysis for biophysical studies of gating mechanisms in ion channels in muscle and nerve cell membranes. The next section of this paper briefly describes the biophysical context. In the third section we summarize the mathematical results that we have been able to obtain and discuss their implications and possible applications. In the fourth section we make some concluding remarks and mention several open questions.

## 2. Biophysical Background

Ion channels are transmembrane proteins with the ability to open a pore through which ions can flow with high conductance; in the absence of such pores, the lipid bilayer membranes of cells are virtually impermeable to most charged particles. Most channels are either voltage "gated" (controlled), such as the sodium channel in nerve axons,

which is the fundamental non-linear circuit element involved in the propagation of nerve impulses, or chemically gated, such as the post-synaptic acetylcholine receptor. We have already reviewed the biophysical context [1][1] and a broad review of ion channels is now available [2]. We shall therefore confine ourselves here to an illustrative example.

The present work started with the desire to extract information from experiments on the chemically gated acetylcholine receptor. Chemically extracted protein is incorporated in an artificial lipid bilayer membrane separating two compartments containing electrolyte solution. The voltage difference between the two sides of the membrane is fixed, and the current through the membrane is measured. In the presence of an agonist (chemical stimulant, such as acetylcholine), the current is found to fluctuate randomly between two levels, reflecting the open or closed state of the channel. (Great pain is taken to arrange to have only one active channel, as shown by the absence of time intervals with current a multiple of the minimum quantum.) It is known that the agonist must bind to the channel to permit opening, and the time scale for channel opening and closing is much shorter than the time scale for agonist binding, as shown by chemical kinetics studies, so the simplest kinetics would be a Markov process with three states: $C_1$ has no bound agonist and the channel is closed, $C_2$ has bound agonist and the channel is still closed, and $O$ has bound agonist and the channel is open, and transitions $C_1 \leftrightarrow C_2 \leftrightarrow O$. The transition $C_2 \leftrightarrow O$ is visible in the experiment as a jump in electric current through the membrane. The transition $C_1 \leftrightarrow C_2$, which involves the binding or dissociation of agonist, is invisible because it does not involve a change in channel conductance. Under these circumstances, we have the simplest example of an aggregated Markov process, with aggregates $\{C_1, C_2\}$ and $\{O\}$. We would like to see if the scheme $C_1 \leftrightarrow C_2 \leftrightarrow O$ is consistent with the data, and, if it is, we would like to estimate the four transition rates in the scheme from the data. In this particular case, it is easy to see, using the results described in the next section, that the transition rates can be estimated from the data, and in fact it is sufficient to use the one dimensional densities for this purpose.

The model we have just described is radically oversimplified and inconsistent with experiment. One feature which is accessible via these experiments but not via agonist binding studies is that the one dimensional density for the channel open "state" (actually, aggregate) is the sum of at least two exponentials [3]. According to the next section, this demonstrates the existence of at least two open states $O_1$ and $O_2$. The aggregates are now $\{C_1, C_2\}$ and $\{O_1, O_2\}$. We would like to accept or reject schemes like $C_1 \leftrightarrow C_2 \leftrightarrow O_1 \leftrightarrow O_2$, $\{C_1 \leftrightarrow C_2 \leftrightarrow O_1, C_2 \leftrightarrow O_2\}$, and $C_1 \leftrightarrow C_2 \leftrightarrow O_1 \leftrightarrow O_2 \leftrightarrow C_2$. The theorems of the next section imply immediately that the one dimensional densities contain all the information available if

---

[1]Numbers in brackets indicate literature references.

any of these schemes is correct, that the third scheme, which contains a cycle, is not identifiable, and that any of these schemes can be excluded if correlation if observed between two consecutive durations of channel opening, which it is [4,5].

## 3. Results

We first introduce some notation. $P(t)$ will denote the transition matrix of the Markov process. We will assume throughout that the process is in equilibrium. As is standard, we let

$$Q = \lim_{t \to 0} \frac{P(t) - I}{t} \text{ , where } I \text{ is the identity matrix.}$$

The aggregates will be indexed by lower case Greek letters; $n_\alpha$ is the number of states in aggregate $\alpha$. We order the states so that states in the same aggregate are contiguous, and we partition the matrix $Q$ into sub-matrices $Q_{\alpha\beta}$. We will assume throughout that the submatrices $Q_{\alpha\alpha}$ are diagonalizable, which holds if the law of detailed balance is valid for the system.

Supposing that the process enters aggregate $\alpha$ at time $t = 0$, we denote the probability density of the length of time $T$ spent in that aggregate before leaving by $f_\alpha(t)$. It is shown in [6] and [1] that

$$f_\alpha(t) = \pi_\alpha e^{Q_{\alpha\alpha} t} q_\alpha$$

where $q_\alpha = \Sigma_{\beta \neq \alpha} Q_{\alpha\beta} u_\beta$, $u_\beta$ is a column vector of $n_\beta$ ones, and $\pi_\alpha$ is a row vector giving the probabilities that aggregate $\alpha$ is entered via each of its states.

Under the assumption that $Q_{\alpha\alpha}$ is diagonalizable this one dimensional density can be re-expressed as:

$$f_\alpha(t) = \sum_{i=1}^{n_\alpha} a_\alpha^i e^{-\lambda_\alpha^i t} \text{ .}$$

An implication of this result is that a lower bound on $n_\alpha$ can be obtained by counting the number of exponential components. This result has been widely used in channel gating studies by fitting experimental results to sums of exponentials and judging how many exponentials to include by a chi-squared statistic.

Two dimensional densities can also be used to obtain information about the Markov process. It can be shown that the density of spending a length of time $s$ in aggregate $\alpha$ and then a length of time $t$ in aggregate $\beta$ is

$$f_{\alpha\beta}(s,t) = \sum_{i=1}^{n_\alpha} \sum_{j=1}^{n_\alpha} a_{\alpha\beta}^{ij} e^{-\lambda_\alpha^i s - \lambda_\beta^j t} \text{ .}$$

(Here we use $i$ and $j$ as indices, not powers.) It is noteworthy

that the same exponential parameters occur in both the one and two dimensional densities. This might be used to judge the plausibility of an underlying Markov process model.

The matrix of linear coefficients with $ij$ entry $a_{\alpha\beta}^{ij}$ yields information about how the aggregates are interconnected:

**Theorem A.** [7] *Let $A_{\alpha\beta}=[a_{\alpha\beta}^{ij}]$ be the matrix of coefficients of a two dimensional density above. Then the rank of $A_{\alpha\beta}$ is less than or equal to the rank of $Q_{\alpha\beta}=p_{\alpha\beta}$, say, and $A_{\alpha\beta}$ depends on at most $p_{\alpha\beta}(n_\alpha+n_\beta-p_{\alpha\beta})$ parameters. The rank of $Q_{\alpha\beta}$ is less than or equal to the smaller of: the number of states in $\alpha$ which are linked directly to states in $\beta$, and the number of states in $\beta$ which are linked directly to states in $\alpha$. It may be possible to empirically obtain a lower bound on the complexity of interconnection by fitting two dimensional densities and using chi-square tests. We hope to implement and test such a procedure in the future.*

Higher dimensional densities can be considered also. The following theorem shows that under certain conditions no more information about the process can be obtained from these densities, however.

**Theorem B.** [7] *If, for each aggregate $\alpha$, the eigenvalues $\lambda_\alpha^i$ are distinct and $\alpha_\alpha^i\neq0$ for all $i$, the higher dimensional densities $f_{\alpha_0\ldots\alpha_r}(t_0,\ldots,t_r)$, $r>1$ are completely determined by the two-dimensional densities.*

Thus, in principle, all the available information about the underlying process can be extracted from the two dimensional densities if the hypotheses of the theorem are satisfied. By counting the number of independent parameters involved in the two dimensional densities, we have the following theorem:

**Theorem C.** [7] *Under the assumptions of Theorem B, the finite dimensional distributions depend on at most $\sum_{\alpha\neq\beta} p_{\alpha\beta}(n_\alpha+n_\beta-p_{\alpha\beta})$ parameters.*

Thus for example, if there are two aggregates, there is information on at most $2p(n_\alpha+n_\beta-p)$ parameters. If a model depends on more than this many parameters, its parameters are not uniquely identifiable.

Theorem A above suggests one way to study the complexity of interconnection between two aggregates. Labarca et. al [5] have also used certain correlation functions for this purpose. For a particular aggregate $\alpha$, say, the sequence of dwell times in that aggregate, $T_1,T_2,\ldots$ is a stationary process, and it can be shown that the covariance function of that process is of the form given in the following theorem:

**Theorem D.** [1] *The covariance function is of the form*

$$\Gamma_\alpha(k)=\sum_{i=1}^{M-1} u_i\kappa_i^{|k|}$$

*where $0\leq\kappa_i<1$ and $k\neq0$ and where $M$ is the rank of the matrix $Q_\alpha\cdot$ which is composed of the off-diagonal blocks corresponding to aggregate $\alpha$ in the matrix $Q$. If $M=1$, $\Gamma(k)=0$ for $k\neq0$.*

The rank of $Q_\alpha\cdot$ is less than or equal to the smaller of: the number of states in alpha which are linked directly to states in any other aggregate, and the number of states in other aggregates which are linked directly to states in $\alpha$.

Correlation functions can thus be used to obtain information similar to that in the two dimensional densities as in Theorem A. In the case that there are more than two aggregates, the two dimensional densities contain finer information, however, since a lower bound on the rank of each of the matrices $Q_{\alpha\beta}$ which constitute $Q_\alpha$ can be obtained.

It is interesting that the covariance function of Theorem D is of the form of the covariance function of a moving average–autoregressive process, although the stationary process is distinctly non-Gaussian. It may well be that techniques developed for order estimation in the time series literature can be used to estimate $M$. Labarca et al. [5] were primarily interested in testing whether $M$ was greater than 1, which is relatively simple since the large sample distribution of the empirical correlation coefficients in the case $M=1$ can be used.

## 4. Further Considerations

We believe that although the results summarized in the previous section are useful, there are still many open and interesting questions.

We have developed some necessary conditions for identifiability which can sometimes be used to conclude that hypothetical models are unidentifiable. It would be useful to have checkable sufficient conditions as well.

Our analysis applies to stationary Markov processes. The nonstationary case is of both theoretical and practical interest, and we hope to consider that situation in the future. Records from the sodium channel are typically nonstationary because of the presence of an absorbing inactivated state.

We have only begun to explore practical consequences of our results for the analysis of experimental data [3]. It is tempting to consider estimating the two dimensional distributions and basing further analysis on them. Horn and Lange [8] have proposed likelihood analysis of sample paths and Horn and Vandenburg [9] have applied these techniques to data from the sodium channel. Advantages of their approach are that it is applicable to nonstationary data and multi-channel data. The likelihood method is computationally intensive, however; Horn uses an array processor on a VAX 11/730 and reports that days of computer time are necessary. An analysis based on the two dimensional distributions would be much faster; it is not clear what loss of statistical efficiency would be incurred.

Finally, it would be desirable to develop more data-analytic and model-free methods for analyzing experimental records to give qualitative insights that might suggest phys-

ical mechanisms. Labarca et al. [3] have used box-plots to advantage in analyzing data from the chloride channel.

## 5. References

[1] Fredkin, D.R.; M. Montal and J. Rice, Identification of aggregated Markovian models: application to the nicotinic acetylcholine receptor. *Proceedings of the Berkeley Conference in Honor of Jerzy Neyman and Jack Kiefer*, Ed. L. Le Cam and R. Olshen, Vol. I, pp 269-290, Wadsworth Publishing Co., Belmont, CA (1985).

[2] Hille, B., *Ionic Channels of Excitable Membranes*, Sinauer Associates, Inc. (Sunderland, MA) (1984).

[3] Labarca, P.; J. Rice, D.R. Fredkin and M. Montal, Kinetic analysis of channel gating: application to the cholinergic receptor channel and the chloride channel from *Torpedo californica. Biophys. J.* **47** 469-478. (1985).

[4] Jackson, M.B.; B.S. Wong, C.E. Morris, H. Lecar and C.N. Christian. Successive openings of the same acetylcholine receptor channel are correlated in their open times. *Biophys. J.* **42** 109-114. (1983).

[5] Labarca, P.; J. Lindstrom and M. Montal, The acetylcholine receptor channel from *Torpedo californica* has two open states. *J. Neurosci.* **4** 473-496. (1984).

[6] Colquhoun, D., and A.G. Hawkes, On the stochastic properties of single ion channels. *Proc. R. Soc. London B* **211** 205-235 (1981).

[7] Fredkin, D.R., and J. Rice, On aggregated Markov processes. To appear in *J. Appl. Prob.* (1984).

[8] Horn, R. and K. Lang, Estimating kinetic constants from single channel data. *Biophys. J.* **43** 207-223 (1983).

[9] Horn, R. and Vandenberg, C., Statistical properties of single sodium channels. *J. Gen., Physiol.* **84** 505-533 (1984).

# Automated Pattern Recognition:
# Self- Generating Expert Systems for the Future

## Thomas L. Isenhour
### Utah State University, Logan, UT 84322

Chemometrics and pattern recognition had their start in chemistry in the late 1960's. The most recent review of the area by Michael DeLaney listed 438 journal articles and books. The three most important areas of future development will be Expert Systems, Relational Data Bases, and Robotics. It should now be possible to combine existing robotics and artificial intelligence software to create a system which will generate its own expert systems using relational data bases. The data will be in the chemical domain and the system I describe we are calling the Analytical Director. The Analytical Director will be an artificial intelligence/robotic expert system for the analytical laboratory. The Analytical Director will develop, test, implement and interpret chemical analysis procedures. It will learn from its own experience, the experience of others and communicate what it has learned to others. The Analytical Director will be a self-generating Expert System. I believe that such systems will, in the future, provide all the advantages of pattern recognition, expert systems and relational data bases in experimental settings. Problems will continue to be defined by human beings, but more and more, the laboratory will design, execute and evaluate its own experiments.

Key words: artificial intelligence; chemical analysis; expert systems; pattern recognition; relational data bases; robotics.

Chemometrics and pattern recognition had their start in chemistry in the late 1960's. Two areas of application appeared almost simultaneously, those being learning machines and project dendral, both applied to spectroscopic interpretation. The former originated in my research group at the University of Washington and have been carried on by us, and my two students, Peter Jurs and Bruce Kowalski, and the latter was developed at the Stanford Artificial Institute by a consortium from Chemistry and Computer Science. Over the past 15 years a variety of applications have occurred which include the following list and probably others: statistics, modeling and parameter estimation, resolution, calibration, signal processing, image analysis, factor analysis, pattern recognition, optimization, artificial intelligence, graph theory and structure handling, and library searching.

The most recent review of the area by Michael DeLaney listed 438 journal articles and books. This was an *Analytical Chemistry* article which covered the last two years of activity. Clearly, pattern recognition and its applications now have an established place in chemistry.

The topic of this presentation is, however, not what has happened up to the point. Rather it is what will happen in the forseeable future. I believe that the three most important areas of future development will be expert systems, relational data bases, and robotics. We will talk a little about each one of these and then go on to deal in detail with what may soon become one of the most, if not *the* most, sophisticated tool that the experimental chemist has ever acquired.

An expert system is, simply stated, a piece of software that behaves like an expert. The origin of expert systems were instruction manuals that told you what to do based upon what you encountered in a step by step fashion.

**About the Author:** Thomas L. Isenhour is with Utah State University's Department of Chemistry and Biochemistry.

Almost everyone is familiar with manuals on "How to Fix Your Chevy Station Wagon," etc. These are non-mysterious, and sometimes non-useful, written recipies designed to lead you by the hand through a repair, construction, gourmet meal preparation, etc., that an expert could do but you could not, at least on your own. Of course, the quality of such systems depends on the knowledge of the expert, and on the successful transfer of that knowledge from that expert to the author and from the author to you.

The modern expert system is usually a computer program that attempts the same thing, with the same limits of success. That is, the quality is dependent on the knowledge of the expert, the successful transfer of that knowledge from that expert to the author and from the author to the user. It is a little more, however, because the machine can use your input to do computations, etc., and while all of this could be done with a programmed manual, it certainly is faster by computer and potentially more accurate. There is less chance of you failing to follow the directions for correct use than in a written expert system.

Relational data bases are collections of data that interrelate along traditional and non-traditional lines. In a sense, the relational data base is the social scientist's dream. Given a set of descriptors for each entry, the relational data base allows very easy cross correlations such as, how many children in Miami who had braces before the age of 12 also have maternal grandparents alive in Manhattan. Again, a relational data base was always possible with pencil and paper, or even better with three-by-five cards, but it can be greatly facilitated in a computer. However, the use is still defined by the selection of the descriptors and the quality of the queries.

A robot, according to one handy dictionary, is "a machine devised to function in place of a human agent." This is quite a broad definition and could refer to an automatic ticket dispenser at a parking garage, an autopilot operating from an inertial guidance system in a jet airplane, or a computer interfaced with an autoanalyser at a clinical laboratory.

It is my contention that it is now possible to combine existing robotics and artificial intelligence software to create a system that will generate its own expert systems using relational data bases. The data will be in the chemical domain and the system I describe we are calling the Analytical Director. Simply stated, the Analytical Director will be an artificial intelligence/robotic expert system for the analytical laboratory. The Analytical Director will develop, test, implement, and interpret chemical analysis procedures. It will learn from its own experience, and that of others, and it will communicate what it has learned to others.

The subject of this research is neither automation nor the robot's roll in automation, but rather EXPERT laboratory management working through a combination of robotics and artificial intelligence. We propose to combine robotics and artificial intelligence into an expert system for the analytical chemistry laboratory. We propose to demonstrate that an Analytical Director can develop, test, implement, and interpret chemical analysis procedures. It is a misconception that the best use of robots will be exhaustive testing of possible solutions to problems. While computationally exhaustive methods are often quite successful, they are rarely useful to analytical chemistry. Artificial intelligence is required for a real breakthrough in automated laboratory methodology.

Consider the following analytical problem which might arise from a relatively simple problem.

Given:
10 possible components to a mixture
10 reagents
10 possible temperatures
10 pH values

If each reaction combination were chemically independent, that is, if the results of any combination could be learned by a linear addition of the separate tests, then 10,000 procedures could be carried out to determine the entire system. This might be feasible if, for example, each test could be completed in one minute. (This would require just about one week of continuous work assuming the robot suffered no maintenance problems or other delays.)

However, chemical reactions are not usually independent. For example, if one of the components were Fe(III) and two of the reagents were $CSN^-$ and citrate ion, there would clearly be complex equilibria interactions. If we redo the calculation considering from 1 to 10 possible components, from 1 to 10 possible reagents and any combination of 10 temperatures and 10 pH values, it requires $1.63 \times 10^{15}$ tests. Again carried out at one minute intervals, assuming the robot could work through the entire set of procedures without interruption, $3.10 \times 10^9$ years would be needed. Experimental design methods might achieve a few orders of magnitude improvement but nothing like the 10 orders of magnitude necessary to make this approach feasible.

Furthermore, in real analytical situations the number of variables and dimensions is often much greater. It is clear that analytical chemistry cannot be done by exhaustive trial and error. Therefore, if robotics is to have any real effect upon the field an intelligent robot must be created that can choose meaningful experiments and profit from its experience, as well as the experience of others.

We propose to construct such a system and test it initially on a very limited analytical problem to prove that artificial intelligence can be used to seek efficient paths to complex analytical problems without resorting to exhaustive trial and error. To do so our first model system will be very simple, involving 3 ions, 5 reagents, 2 temperatures and 3 pH's. This system can be exhaustively tested with 4320 procedures. Assuming 1 minute tests, 3 days would be required. Exhaustive testing of this model system will produce a set of observations that will facilitate the development and test-

ing of generalized data structures and optimization routines to be used by the Analytical Director for more complex problems. For complex problems, the Analytical Director must develop artificial intelligence methods that circumvent the testing approach.

We have selected a developmental domain that is a closed system of simple analytical chemical problems. The domain will be wet and photometric analysis of simple cations. The set of manipulative skills required is purposely limited to the abilities of the Zymark system, and the chemical reactions and spectrophotometric measurements possible limited

to those that can be performed with the available equipment. This way, we plan to be able to test thoroughly the creative capabilities of the artificial intelligence programs to be developed for the Analytical Director. We have further selected the domain of water analysis for an advanced test of the Analytical Director. Only by using an unbounded problem will we be able to demonstrate the true capability of the Analytical Director.

Given a set of standards, reagents, and manipulative skills, the Analytical Director will develop its own set of tests for each individual cation. These data will be stored in a relational data base keyed on ions, reagents, conditions and results of spectroscopic measurements. It will be assumed initially that no chemistry is known for these elements or reagents. After developing possible individual tests, these tests will be cross compared to identify likely interfering reactions. Tests will be characterized by their quality as defined by time, expense and reproducibility. As the best compromise of these components is a value judgment, an adjustable value coefficient will be developed. Then possible mixture methods will be systematically tested and compared for success. As this proceeds the relational data base will continue to expand. Finally, new unknowns will be introduced into the system to test the ability of the Analytical Director to adapt to new circumstances. At this point literature information will be added to the data base. The Analytical Director will thereby "learn" from the experience of others. Further, the Analytical Director will report developed procedures for possible use by others.

In summary, the Analytical Director will be a self-generating expert system. I believe that such systems will, in the future, provide all the advantages of pattern recognition, expert systems, and relational data bases in experimental settings. Problems will continue to be defined by human beings, but more and more the laboratory will design, execute, and evaluate its own experiments.

# Regression Analysis of Compartmental Models

## T. L. Lai

### Columbia University, New York, NY 10027

Herein we study the problem of assessing, on the basis of noisy and incomplete observations, how much information there is in the data for model identification in compartmental systems. The underlying concept is that of an "information distance" between competing models, and estimation of this distance on the basis of the given data is discussed. Useful reduction of the dimensionality of the corresponding least squares problem is accomplished by regarding the decay rate constants as primary parameters of interest and the other parameters of the model as nuisance parameters. Estimation of the decay rate function is also discussed.

Key words: compartmental models; decay rate constants; separable least squares problems; system identification.

## 1. Introduction

Compartmental models are widely used in many fields of science and engineering—pharmacokinetics, biochemistry, physiology, radioactive isotopes and tracers, to name a few. The basic equations of a typical (linear) compartmental model with $k$ compartments are

$$dx_i/dt = c_{i0} + \sum_{j=1, j \neq i}^{k} (c_{ij} x_j - c_{ji} x_i) - c_{0i} x_i , \qquad (1)$$

$0 \leqslant t < \infty, i = 1, ..., k$; where $c_{ij}$ are nonnegative constants (called the "turnover rate constants"), compartment 0 denotes the environment, and $x_i = x_i(t)$ is the amount of material in compartment $i$ at time $t$. Letting $f_{i0} = c_{i0}$ and $f_{ij} = c_{ij} x_j$ for $j \neq 0$, $f_{ij}$ represents the mass flow rate to compartment $i$ from compartment $j$. Under certain assumptions (cf. [3], [12])[1], integration of differential eq (1) leads to

$$x_i(t) = \beta_i + \sum_{j=1}^{k} \alpha_{ij} e^{-\lambda_j t}, \qquad (2)$$

---

---

[1] Figures in brackets indicate literature references.

$i = 1, ..., k$; where $\lambda_1, ..., \lambda_k$ are positive constants (called the "decay rate constants") depending only on the turnover rate constants $c_{ij}$, and $\beta_i$ and $\alpha_{ij}$ are nonnegative constants.

Noting that it is often not possible to have accurate and complete measurements from all the compartments, Berman and Schoenfeld [2] considered the degrees of freedom in choosing a compartmental model compatible with the data when measurements are incomplete. In particular, when one can only measure aggregate input-output characteristics so that one observes only $x(t) = \sum_{i=1}^{k} x_i(t)$, then one can only identify the rate constants $\lambda_1, ..., \lambda_k$ of the model (2) from the equation

$$x(t) = B + \sum_{j=1}^{k} A_j e^{-\lambda_j t}, \qquad (3)$$

where $B = \sum_{i=1}^{k} \beta_i$ and $A_j = \sum_{i=1}^{k} \alpha_{ij}$. In pharmacokinetic applications, Wagner [10] has shown that many basic quantities of interest can be expressed in terms of the parameters of the reduced model (3), the use of which he recommends in lieu of the full model (2) whose specification usually leads to ambiguities in these applications because of noisy and incomplete measurements.

The difficulties in model formulation and identification are compounded when the measurements are not only incomplete but also are subject to error. This leads to the question of assessing, on the basis of noisy

and incomplete observations, how much information there is in the data for model identification. If the amount of information is found to be inadequate for meaningful specification of the full model, then it may be more useful to work with a reduced model or directly with certain scientifically important characteristics (functions of the parameters) of the model.

Since both eqs (2) (for the individual compartments) and (3) (for the whole system) express the response as a polyexponential function of time $t$, the statistical problems in the present context are basically those of parameter estimation for polyexponential regression models. We discuss herein 1) the information content in the data from these regression models, and 2) estimation of model parameters and certain functions thereof.

## 2. Noisy Data and Information Content

Consider the polyexponential regression model

$$y_i = \beta + \sum_{j=1}^{k} \alpha_j \exp(-\lambda_j t_i) + \epsilon_i, \ i = 1, \ldots, n, \qquad (4)$$

where the $\epsilon_i$ represent random errors. The errors $\epsilon_i$ are usually assumed to be independent random variables with zero means. Letting $\theta = (\lambda_1, \ldots, \lambda_k; \alpha_1, \ldots, \alpha_k, \beta)$ and

$$f_\theta(t) = \beta + \alpha_1 e^{-\lambda_1 t} + \ldots + \alpha_k e^{-\lambda_k t},$$

common models for $\mathrm{Var}(\epsilon_i)$ are:

(i) $\mathrm{Var}(\epsilon_i) = \sigma^2$ (constant variance error model),
(ii) $\mathrm{Var}(\epsilon_i) = f_\theta(t_i)\sigma^2$ (constant coefficient of variation model),
(iii) $\mathrm{Var}(\epsilon_i) = f_\theta(t_i)\sigma^2$ (Poisson-type error model).

Statistical methods for estimating the unknown parameters of the regression function try to "average out" the random errors in various ways. Assuming that observations are taken at equally spaced times $t_i = (i-1)\Delta$ and that $n = (2k+1)m$, the method of Lanczos [6, p. 273] and Cornell [4] uses the sample means

$$\bar{y}_r = m^{-1} \sum_{i=(r-1)m}^{rm-1} y_i \ (r = 1, \ldots, 2k+1)$$

to estimate the moving averages $\mu_r = m^{-1} \sum_{i=(r-1)m}^{rm-1} f_\theta(t_i)$ of the regression function. Replacing $\mu_r - \mu_{r-1}$ by $\bar{y}_r - \bar{y}_{r-1}$ in Prony's algebraic equation defining $\zeta_i = e^{-\lambda_i m \Delta}$, solution of the algebraic equation then gives the estimate of $\lambda_i$. This method therefore tries to average out the random errors by introducing the new parameters $\mu_r$ and using the sample means $\bar{y}_r$ to estimate $\mu_r$.

Since the $\bar{y}_r$ are strongly consistent estimates of $\mu_r$, it follows that the Cornell-Lanczos method is consistent, as was established by Cornell [4]. However, Lanczos [6] gave an example to demonstrate "surprising numerical snags which may develop on account of the exceedingly nonorthogonal behavior of exponential functions." The true regression function in Lanczos' example is

$$f_\theta(t) = 0.0951 \ e^{-t} + 0.8607 \ e^{-3t} + 1.5576 \ e^{-5t}. \qquad (5a)$$

On the basis of 24 successive decay observations in the time interval $(0,1.2]$ of the form 2.51, 2.04, ..., 0.07, 0.06, which are accurate up to 1/2 unit of the second decimal, and assumming knowledge of $\beta = 0$ and $k = 3$, the preceding method gives the fitted model

$$f_\theta(t) = 0 + 0.305 \ e^{-1.58t} + 2.202 \ e^{-4.45t}. \qquad (5b)$$

Although the true parameters are disappointingly far from their estimates, the fitted function (5b) is remarkably close to the true function (5a), and one cannot distinguish between the two models within the errors of the measurements.

An alternative method to estimate the unknown parameters is that of least squares. The estimate $\hat{\theta}$ of $\theta$ is the value of $\gamma$ that minimizes

$$S(\gamma) = \sum_{i=1}^{n} w_i [y_i - f_\gamma(t_i)]^2, \qquad (6)$$

where the $w_i$ are suitably chosen weights. Note that

$$S(\gamma) - S(\theta) = \sum_{1}^{n} w_i [f_\theta(t_i) - f_\gamma(t_i)]^2 + 2 \sum_{1}^{n} w_i [f_\theta(t_i) - f_\gamma(t_i)] \epsilon_i. \qquad (7)$$

Moreover, if the weights $w_i$ are so chosen that $w_i \mathrm{Var}(\epsilon_i)$ are bounded, then

$$\sum_{1}^{n} w_i [f_\theta(t_i) - f_\gamma(t_i)] \epsilon_i =$$
$$o\left(\sum_{1}^{n} w_i [f_\theta(t_i) - f_\gamma(t_i)]^2\right) \ w.p. \ 1 \qquad (8)$$

(with probability 1) as

$$d(\theta, \gamma) = \sum_{1}^{n} w_i [f_\gamma(t_i) - f_\theta(t_i)]^2 \to \infty. \qquad (9)$$

The quantity $d(\theta, \gamma)(= E\{S(\gamma) - S(\theta)\})$ defined in (9) is a measure of the separation (distance) between $\gamma$ and the parameter vector $\theta$ reflected by the data. When the $\epsilon_i$ are normal $N(0, 1/w_i)$ random variables, $1/2 \ d(\theta, \gamma)$ is the Kullback-Leibler information number and the least squares estimate coincides with the maximum

likelihood estimate. Thus, the least squares method averages out the errors $\epsilon_i$ in the weighted sums $\sum_1^n w_i [f_\theta(t_i) - f_\gamma(t_i)] \epsilon_i / d(\theta, \gamma)$ for all choices of $\gamma$ under consideration. Consistency of this method under assumption (9) is therefore an immediate consequence of (7) and (8) if there are only finitely many such choices among which one is the true parameter vector, and its extension to infinite parameter spaces involves additional compactness arguments (cf. [7], [11]).

Since $f_\gamma(t)$ is linear in the parameters $\beta$, $\alpha_1$, ..., $\alpha_k$, least squares estimates of these "linear" parameters are given by the standard formulas in multiple linear regression theory for every fixed $\lambda = (\lambda_1, ..., \lambda_k)$. For fixed $\lambda$, define

$$S^*(\lambda) = \min_{\beta, \alpha_1, ..., \alpha_k} \sum_{i=1}^n w_i [y_i -$$
$$(\beta + \alpha_1 e^{-\lambda_1 t_i} + ... + \alpha_k e^{-\lambda_k t_i})]^2, \qquad (10)$$

and the original least squares problem is reduced to that of minimizing $S^*(\lambda)$ involving only the decay rate constants $\lambda_1$, ..., $\lambda_k$. This approach, suggested by Golub and Pereyra [5] and Osborne [8], has the great advantage of reducing the dimensionality of the parameter vector from $2k + 1$ to $k$. Using this approach to fit a two-compartment model to the Lanczos data, Varah [9] has recently shown that $S^*(\lambda_1, \lambda_2)$ (with equal weights) has ill-conditioned Hessian matrices and is very flat over a broad region containing the minimum.

We now show how a global analysis of the least squares function $S^*(\lambda)$ enables us to assess how much information there is in the data to specify the model. Not only does such analysis provide a relatively stable numerical algorithm for finding the least squares estimates of the model parameters, but it also sheds light on the range of models that are compatible with the data.

To fix the ideas, consider the Poisson-type error model $\text{Var}(\epsilon_i) = f_\theta(t_i)\sigma^2$, in which $f_\theta(t)$ is large, at least during the initial portion of the sampling times, as is often the case in tracer measurements. For large $f_\theta(t_i)$,

$$\log y_i = \log f_\theta(t_i) + \log\{1 + \epsilon_i / f_\theta(t_i)\}$$
$$\approx \log f_\theta(t_i) + \eta_i / (f_\theta(t_i))^{\frac{1}{2}}, \qquad (11)$$

where $\eta_i = \epsilon_i / (f_\theta(t_i))^{\frac{1}{2}}$ has mean 0 and constant variance $\sigma^2$. This suggests that $f_\theta(t_i)$ can be estimted with small relative error when $f_\theta(t_i)$ is large. Therefore we introduce "ideal" weights of the form $w_i^* = 1/\max\{f_\theta(t_i), C\}$, where $C$ is some large constant, and define the actual weights

$$w_i = 1/\max\{\bar{f}_\theta(t_i), C\}, \qquad (12)$$

where $\bar{f}_\theta(t)$ denotes some initial estimate of $f_\theta(t)$ such

that $\bar{f}_\theta(t)$ is proportionally close to $f_\theta(t)$ if $f_\theta(t)$ is large.

With this choice of weights $w_i$, we define the least squares function $S^*(\lambda)$ by (10) and study its global and local properties by using both discretization and gradient methods. The idea is to partition the $k$-dimensional parameter space $\Lambda$ (set of all possible values for $\lambda$) into a finite number of subregions. The minimum $S^*(\lambda_D)$ in each subregion $D$ is found by standard gradient-type (such as the Marquardt or Fletcher) algorithms. The minimum of $S^*(\lambda_D)$ over all subregions $D$ then gives the least squares estimate $\hat{\lambda}$ $(S^*(\hat{\lambda}) = \min_{\lambda \in \Lambda} S^*(\lambda))$. Moreover, those values of $S^*(\lambda_D)$ that are proportionally close to $S^*(\lambda)$ also give a range of models compatible with the data.

Figure 1 illustrates the results of this analysis in the regression model

$$y_i = 100 e^{-t_i} + 1000 e^{-5t_i} + \epsilon_i, \, i = 1, ..., n, \qquad (13)$$

where $n = 50$, $t_i = (0.01)i$, and the $\epsilon_i$ are independent normal random variables with zero means and $\text{Var}(\epsilon_i) = Ey_i$. Here $\lambda_1 = 1$, $\lambda_2 = 5$, $\alpha_1 = 100$ and $\alpha_2 = 1000$ are the unknown parameters, and $\beta = 0$ is assumed known. The initial estimates $\bar{f}_\theta(t_i)$ in the weights (12), where we set $C = 30$, are obtained by using the Cornell-Lanczos estimate of $\theta$. Prior knowledge of the inequality constraints $0 < \lambda_1 \leqslant 4$ and $4 \leqslant \lambda_2 \leqslant 10$ is assumed, and we divide this parameter space into 24 unit squares. In 100 simulation experiments performed, we obtained results similar to those in figure 1. In figure 1 reporting one such simulation, $S^*(\lambda_D)$ is shown inside each unit square $D$, near the minimizing point $\lambda_D$ which is represented by a small triangle. The solid triangles denote those $\lambda_D$ whose $S^*(\lambda_D)$ values lie within 10% of the minimum value, which is underlined in the figure. At the true parameter vector $\lambda = (1,5)$, $S^*(\lambda) = 51.5$, which differs from the minimum value of 47.1 by about 9%. The curves represent the contour $S^*(\lambda) = 52$, so that $S^*$ lies within 10% of its minimum value in the region between the curves.

The wide range of models compatible with the data in figure 1 is in sharp contrast to figure 2, where in addition to the data during the time interval $[0, 0.5]$ of figure 1 we took 75 additional observations (generated by the model (13)) at equally spaced times in the subsequent period $(0.5, 1.25)$. In figure 2 there is a relatively small range of parameter vectors $\lambda$ whose $S^*(\lambda)$ values are near the underlined minimum value at the least squares estimate $\hat{\lambda}$, which is remarkably close to true parameter $\lambda = (1,5)$.

Regarding $S(\gamma) - S(\theta)$ as an estimate of the "information distance" $d(\theta, \gamma) = \sum_1^n w_i [f_\theta(t_i) - f_\gamma(t_i)]^2$, we can use it to assess the compatibility of the model $f_\gamma$ with the data. To illustrate this idea, consider the 14 models repre-
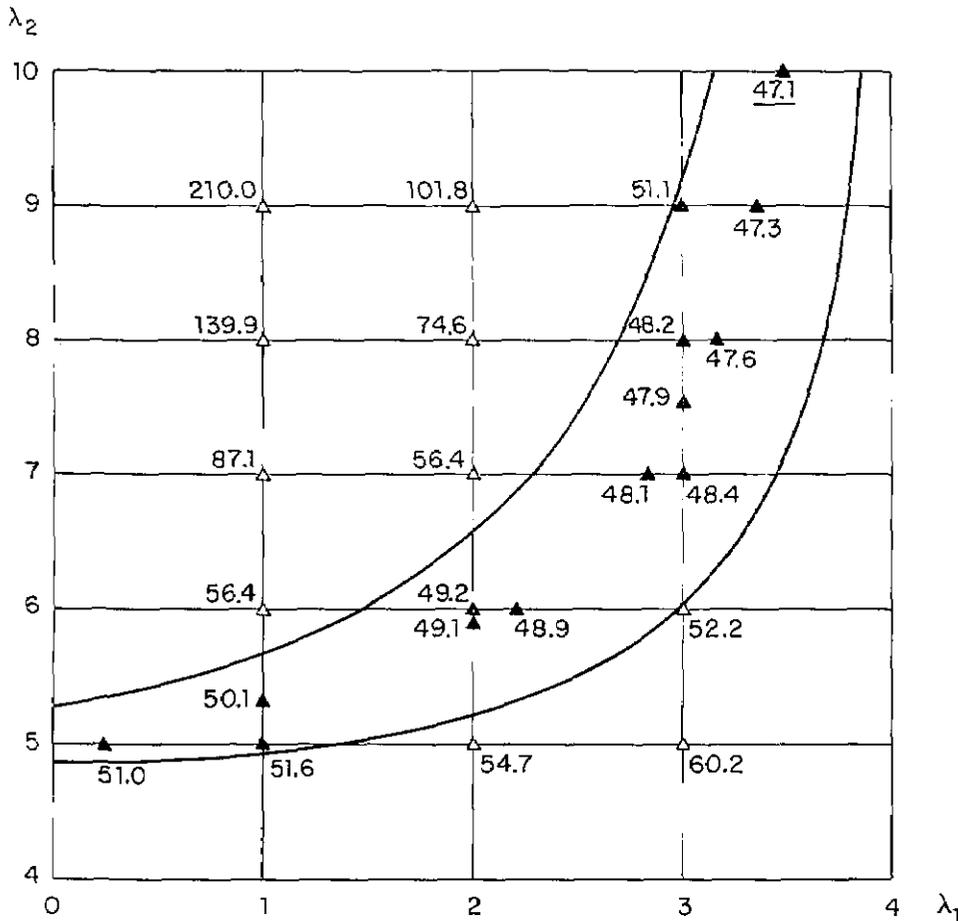
sented by solid triangles in figure 1, whose $S^*$ values lie 10% of the minimum value of 47.1. We tabulate below $S(\gamma)-S(\hat{\theta})$ for each of these 14 values of $\gamma$. In addition, the corresponding information distances $d(\theta,\gamma)$ and $d^*(\theta,\gamma)=\sum_1^n[f_\theta(t_i)-f_\gamma(t_i)]^2/f_\theta(t_i)$ are also tabulated for comparison. Note that $d^*(\theta,\gamma)$ is remarkably close to $d(\theta,\gamma)$; moreover, $d^*(\theta,\gamma)$ is so small (in good agreement with $S(\gamma)-S(\hat{\theta})$) that $f_\gamma(t)$ is within 2% of $f_\theta(t)$ over the observed time range $0.01\leqslant t\leqslant 0.5$ in each of the 14 cases.

| $S(\gamma)-S(\hat{\theta})$ | $d(\theta,\gamma)$ | $d^*(\theta,\gamma)$ | $S(\gamma)-S(\hat{\theta})$ | $d(\theta,\gamma)$ | $d^*(\theta,\gamma)$ |
|---|---|---|---|---|---|
| 0 | 0.63 | 0.63 | 1.0 | 2.20 | 2.24 |
| 0.2 | 3.18 | 3.25 | 1.8 | 1.58 | 1.61 |
| 4.0 | 1.67 | 1.69 | 2.1 | 2.80 | 2.84 |
| 0.5 | 2.72 | 2.78 | 2.0 | 2.11 | 2.14 |
| 1.1 | 5.21 | 5.28 | 3.0 | 1.67 | 1.69 |
| 0.8 | 3.01 | 3.06 | 4.5 | 0.06 | 0.06 |
| 1.3 | 1.20 | 1.23 | 3.9 | 0.67 | 0.63 |

In the case of normal $N(0,1/w_i)$ random errors $\epsilon_i$, $\exp\{-1/2\,(S(\gamma)-S(\hat{\theta}))\}$ is the generalized likelihood ratio for testing $H_0{:}\theta=\gamma$. We can also construct confidence regions for $\theta$ by using contours of the function

$S(\gamma)$, as will be shown elsewhere. In this connection, Bates and Watts [1] recently proposed another useful method involving parameter transformations to improve the standard asymptotic approximations for constructing confidence regions.

The separation of the full parameter vector $\theta$ into its linear and nonlinear components is not only of computational interest, but it also has basic statistical implications. Analogous to the preceding paragraph, in the case of independent $N(0,1/w_i)$ random errors, $\exp\{-1/2\,(S^*(\lambda)-S^*(\hat{\lambda}))\}$ is the generalized likelihood ratio for testing whether $\lambda$ is the true vector of decay rate constants, with $\beta$, $\alpha_1$, ..., $\alpha_k$ as the nuisance parameters. This idea can be easily extended to simultaneous equations (2) defining $k$-compartment models, where we have

$$y_v(t)=\beta_v+\sum_{j=1}^{k}\alpha_{vj}e^{-\lambda_j t}+\epsilon_v(t),\quad v=1,\ ...,\ k.$$

We can similarly define

$$S^*(\lambda)=\min_{(\beta_v,\alpha_{vj})_{1<v,j<k}}\sum_t\sum_{v=1}^{k}w_v(t)[y_v(t)-\beta_v-\sum_{j=1}^{k}\alpha_{vj}e^{-\lambda_j t}]^2,$$

$\lambda_2$

Figure 2.

281.5

1584.7△    348.8    212.4    287.6

924.7    242.9    188.5    304.2

538.3    169.8    164.2    335.3

△150.5

249.4    141.8    153.6    385.9

130.6△△130.5

130.6 | 130.7    228.4    461.0
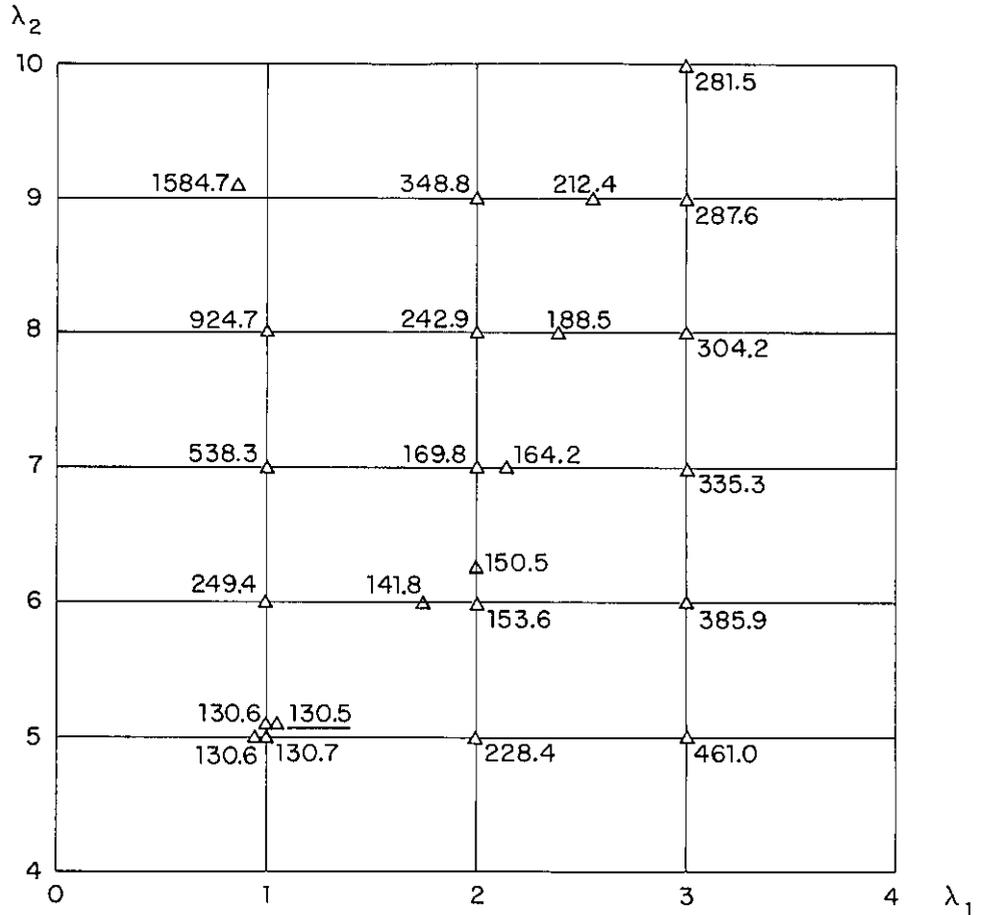
0        1        2        3        4    $\lambda_1$

and interpret $\exp\{-1/2\,(S^*(\lambda)-S^*(\hat{\lambda}))\}$ as a generalized likelihood ratio for testing whether $\lambda$ is the true vector of decay rate constants on the basis of data from all $k$ compartments.

## 3. Decay Rates and Their Estimation

Consider the polyexponential regression model (4). As illustrated in figure 1, one often encounters a wide range of rate parameter vectors $(\lambda_1, ..., \lambda_k)$ that are compatible with one's data. In such circumstances, since there is not enough information to estimate the individual rate constants, it is more meaningful to consider them in a combined characteristic, such as the fractional rate of decay

$$r(t)=\lim_{\delta\to 0}\ \delta^{-1}\{1-f_\theta(t+\delta)/f_\theta(t)\}=-(d/dt)\,\log f_\theta(t)$$

at different time points $t$ of interest. Letting $\lambda_0=0$ and $\alpha_0=\beta$, the logarithmic derivative of $f_\theta(t)=\sum_{j=0}^{k}\alpha_j\,e^{-\lambda_j t}$ is given by $-1$ times

$$r(t)=\sum_{i=0}^{k}p_i(t)\lambda_i,\ \text{where}\ p_i(t)=\alpha_i\,e^{-\lambda_i t}/\sum_{j=0}^{k}\alpha_j\,e^{-\lambda_j t}. \quad (14)$$

Thus, $r(t)$ is a convex combination of the rate constants $\lambda_i$, with the proportional size of the $i^{\text{th}}$ exponential term as the natural weighting factor for the decay rate $\lambda_i$.

To estimate $r(\tau)$ at a particular time point $\tau$ within the range of sampling times, we propose to balance "global information" from all sampling times (leading to weighted least squares estimates described in sec. 2) with "local information" from only the sampling times near $\tau$. We start by using the method of section 2 on all the data to find a region $\Lambda_0$ of rate parameters $(\lambda_1, ..., \lambda_k)$ that are compatible with the data. To choose a vector in $\Lambda_0$ that will provide the best estimate of $r(\tau)$, we note that the logarithmic derivative $-r(\tau)$ is a "local" quantity involving only sampling times near $\tau$, and it is therefore reasonable to weight the observations not only by their variability but also by how far their sampling times are from $\tau$, putting more weight on sampling times near $\tau$. With this new set of weights $w_i(\tau)$, we calculate the (constrained) least squares estimate $\hat{\theta}(\tau)$ of the parameter vector $\theta$, under the constraints $\lambda\epsilon\Lambda_0$ and $\alpha_i\geqslant 0(i=0, 1, ..., k)$. Substituting the unknown parameters in eq (14)

529

by these least squares estimates, we obtain the estimate $\hat{r}(\tau)$ of $r(\tau)$.

A detailed discussion of the procedure sketched above, together with a comparative study of this approach and the popular curve-peeling methods for compartmental analysis (cf. [12]), will be presented elsewhere.

## References

[1] Bates, D. M., and D. G. Watts, Parameter transformations for improved approximate confidence regions in nonlinear least squares Ann. Statist. 9, 1151 (1981).

[2] Berman, M., and R. Schoenfeld, Invariants in experimental data on linear kinetics and the formulation of models. J. Appl. Phys. 27, 1361 (1956).

[3] Coddington, E. A., and N. Levinson, Theory of Ordinary Differential Equations. McGraw-Hill, New York (1955).

[4] Cornell, R. G. A method for fitting linear combinations of exponentials. Biometrics 18, 104 (1962).

[5] Golub, G. H., and V. Pereyra, The differentiation of pseudo-inverses and nonlinear least squares problems whose variables separate. SIAM J. Numer. Anal. 10, 413 (1973).

[6] Lanczos, C., Applied Analysis. Prentice Hall, Englewood Cliffs (1956).

[7] Ljung, L., Consistency of the least squares identification method. IEEE Trans. Autom. Contr. AC-21, 779 (1976).

[8] Osborne, M. R., Some special nonlinear least squares problems. SIAM J. Numer. Anal. 12, 571 (1975).

[9] Varah, J. M., On fitting exponentials by nonlinear least squares. SIAM J. Sci. Stat. Comput. 6, 30 (1985).

[10] Wagner, J. G., Linear pharmacokinetic equations allowing direct calculation of many needed pharmacokinetic parameters from the coefficients and exponents of polyexponential equations which have been fitted to data. J. Pharmacokin. Biopharm. 4, 443 (1976).

[11] Wu, C. F., Asymptotic theory of nonlinear least squares estimation. Ann. Statist. 9, 501 (1981).

[12] Zeirler, K., A critique of compartmental analysis. Ann. Rev. Biophys. Bioeng. 10, 531 (1981).

# DISCUSSION

of the T. L. Lai paper,
Regression Analysis of Compartmental Models

## T.-H. Peng

Oak Ridge National Laboratory

One of the major tasks of marine geochemists is determining the uptake by the sea of $CO_2$ derived from the combustion of fossil fuels. Until valid models of the general circulation of the ocean are constructed, this task will have to be done with box models calibrated through use of the distribution of natural radioisotopes and transient tracers.

We need to explore how sensitive the uptake of fossil fuel $CO_2$ is to the basic design of these models and how the design can be improved by simultaneously fitting the distributions of several tracers. Five different 11-box thermocline circulation models of the temperate North Atlantic were constructed for this purpose.* Anthropogenic tritium, $^3$He, and radiocarbon are used as tracers to calibrate these models. The temporal input functions of these tracers differ considerably from one another. So also do the geographic patterns of their inputs and their geochemistries in the sea.

Using the basic equation of the box model [e.g., eq (1) of T.L. Lai's presentation at this conference] and the finite difference method of computation for mass balance in each box, these thermocline ventilation models with differing circulation patterns were calibrated to yield a tritium distribution similar to that observed during the Geochemical Ocean Section Studies (GEOSECS) survey in 1973. These models were then run for $^3$He and bomb-produced $^{14}$C. While the models differ significantly in their ability to match the observed $^3$He and $^{14}$C distributions, these differences are not large enough to clearly single out one model as superior. This insensitivity of tracer to tracer ratio to model design is reflected by the nearly identical uptake of $CO_2$ by the various models. This result also suggests that the uptake of $CO_2$ by the sea is limited more by the rates of physical mixing within the sea than by the rate of gas exchange across the sea surface.

# Measurement and Control
# of Information Content
# in Electrochemical Experiments

## Sam P. Perone and Cheryl L. Ham

### Lawrence Livermore National Laboratory, Livermore, CA 94550

One of the most important problems in chemical analysis is the interpretation of analytical data. The difficulty of this task has been further compounded by the data explosion. Chemical information relevant to the particular analysis problem is hidden within excessive amounts of data. This problem could be alleviated through knowledge and control of the information content of the data. Information theory provides a means for the definition, evaluation, and manipulation of quantitative information content measurements. This paper provides a general review of some of the basic concepts in information theory, including history, terminology, entropy, and other information content measures. The application of information theory to chemical problems requires some modifications. The analyst is usually only interested in a subset of the information (data) which has been collected. Also, this relevant chemical information is dependent upon not only the informational goals of the problem, but the completely specified procedure as well. This paper reviews chemical applications of information theory which have been reported in the literature including applications to qualitative analysis, quantitative analysis, structural analysis, and analytical techniques. Measures of information and information content and figures of merit for performance evaluations are discussed. The paper concludes with a detailed discussion of the application of information theory to electrochemical experiments and the empirical determination of the information content of electroanalytical data.

## Introduction

Data interpretation is one of the most challenging problems of chemical analysis. Both the data-rich and the data-limited cases stress the need for efficient methods to extract chemical information from the available data. Data-rich analyses result from the ability of modern chemical instrumentation to generate enormous amounts of data in short periods of time. The current trend towards more exotic

hybrid instruments buries the chemical information even deeper within the data. Alternatively, data-limited analyses often result from limitations in appropriate sensors, accessible techniques, time, and manpower. The need for efficient methods to extract chemical information is superseded only by the need to acquire information-rich data.

Improving accessibility of chemical information emphasizes the importance of good experimental design and requires a re-evaluation of the traditional approach to chemical analysis. The typical approach involves using foreknowledge about the samples to choose an anlytical procedure. The analytical procedure, which may involve more than one analytical technique, is used to produce as much data as possible which is collected for later analysis. Data interpretation is performed by the analyst using as much intuition and background knowledge as possible. In

the data-rich case, emphasis is often on data reduction. The data-limited case emphasizes information extraction.

A more efficient and desirable approach to chemical analysis would be to maximize the amount of information obtained relevant to the current problem, while minimizing the amount of analytical effort, time, and collected data. The evaluation, selection, and optimization of analytical procedures need to be investigated further. In order to study this problem, a method for the quantification of chemical information must be defined. This paper reviews some relevant information theory concepts and describes applications to chemical analysis. Illustrations have been taken from the literature, as well as from our own recent work, to demonstrate the value of applied information theory for optimized chemical analysis methods.

## Information Theory Concepts

Information theory [1-4][1] is concerned with the study of information and its transmission. The earliest information theorists studied the encoding and decoding of secret codes. Modern work on information theory can be traced back to the 1920's beginning with Carson's study of frequency modulation. Nyquist determined the minimum bandwidth required for the transmission of a definite quantity of information. Hartley established a definite bandwidth-time product required for the tranmission of a definite quantity of information. However, most of the present work in information theory is based upon probablistic models of communication developed by C.E. Shannon in 1948 [5]. Simply stated, the basic principle is that a message with a high probability of random occurrence conveys little information. The most information is conveyed by the message that is least likely to spontaneously occur.

This principle is formalized by the concept of entropy which equates information and uncertainty. Entropy is a quantitative measure of the amount of information supplied by a probabilistic experiment. It is based upon classical Boltzmann entropy from statistical physics. Shannon's formula, eq (1) defines the entropy or the average information

$$H_{av} = I_{av} = - \sum_{i}^{N} p(x_i) \log_2 p(x_i) , \qquad (1)$$

to be equal to the weighted average of the specific information for each event in the system under consideration. Specific information [6] is information conveyed by the occurrence of a particular event and is quantified by the $-\log_2$ of the probability of the event ($p(x_i)$). Entropy is limited by a maximum of $\log_2 N$ (Hartley's equation), where

$N$ is the number of events, when considering a set of mutually exclusive but equally probable events [1,6,7]. For example, consider the measurement of three distinguishable intensity levels. If the probabilities of measuring these levels are 0.25, 0.30, and 0.45, respectively, the averge entropy would equal 1.5 bits. The amount of specific information conveyed by the measurement of each level would equal 2.0, 1.7, and 1.1 bits, respectively. Notice that the least likely level to be measured does indeed convey the most information. The maximum entropy is equal to 1.6 bits. A thorough treatment of the mathematical basis of entropy and its properties is given in a book by Mathai and Rathie [2].

Redundancy [6] is the difference between the maximum information and the average information eq (2). Relative

$$I_{red} = \log_2 N - \sum_{i}^{N} p(x_i) \log_2 p(x_i) \qquad (2)$$

redundance is the ratio of the redundance to the maximum information [6]. Relative information content is the ratio of the actual average information to the maximum information [1]. Redundancy can then be expressed as the remaining fraction not due to relative information [1]. In the above example there is 0.1 bit of redundancy and 0.062 relative redundance. If the actual average information is equal to 1.4 bits, the relative information content is equal to 0.88 and the redundancy equals 0.12.

## Types of Information

The concept of information as used in information theory refers to the choice or uncertainty of outcomes when regarding distinguishable elements with respect to some random mechanism. Information is a system based upon elements that have been agreed upon for the representation of information (characters) and the relationships between them (codes). It is not a measure of meaning as used in the usual sense, which implies a subjective evaluation. The concept of information as used in chemical analysis encompasses the uncertainty regarding the quantity, identity, chemical structure, or properties of the analyte of interest.

Preinformation or foreknowledge [1] is the prior knowledge concerning the occurrence of events. The information conveyed by the occurrence of more than one independent event is simply the sum of the information conveyed by each event individually. However, if the occurrence of a second event is dependent upon the occurrence of the first event, the foreknowledge of the first event reduces the amount of information conveyed by the second event. Therefore, the amount of information conveyed by a series of events within the system under consideration is always less than or equal

to the sum of the information conveyed by each of the events separately.

Preinformation is probably the most commonly used information theory concept in chemical analysis. Chemical preinformation [8] is information that is known prior to performing the current analysis. It may result from experience, preliminary analyses, etc. It is used to reduce the effort required to solve the analytical problem. Preinformation may be quantified through the use of entropy-based measures [8]. The uncertainty before the analysis for a discrete variable, such as chemical identity, is quantified by the use of the a priori probabilities for identification in Shannon's equation (eq 1). The uncertainty for a continuous variable such as concentration or signal intensity is expressed by integrating the a priori probability density function over the range of interest eq. (3).

$$H(X) = -\int_{x_1}^{x_2} p(x) \log_2 p(x) \, dx \qquad (3)$$

Joint information [1] is information that is provided by more than one event. It can be quantified by substituting the joint probability of occurrence for the events into Shannon's equation (eq 1). In the case of independent events, the joint probability of occurrence is simply the product of the a priori probabilities of occurrence. For nonindependent events, it is the product of the a priori probability for the first event with the conditional probabilities for the other events.

Mutual information describes the amount of information in one event that determines the state of another event. It may also be thought of as the average amount of information required to distinguish members of different classes. Isenhour et al. [9] investigated the relationship between mutual information and classification in the determination of chemical functionality for 200 compounds based upon binary encoded (peak/no peak) infrared spectra. Mutual information was calculated as the difference between the total average entropy and the average conditional entropy. The total average entropy is the average amount of information required to distinguish between the 200 spectra under consideration. It is calculated as a weighted average of the probability of occurrence of a peak maximum for each spectral interval using Shannon's formula (eq 1). The average conditional entropy is the average amount of information required to distinguish between members of the same class. Average conditional entropy is calculated as the sum of the class conditional entropies weighted by class size. In the case of two separable, equally probable classes, the independent mutual information is equal to one bit. A value for mutual information greater than one bit implies the inclusion of redundant information in the data. The square root of mutual

information was shown to be linearly related to the maximum likelihood classification ability.

## Figures of Merit

The application of information theory concepts to chemistry is most familiar in the evaluation of analytical methods. Figures of merit such as accuracy, precision, and detection limit have long been used to evaluate the attainment of informational goals such as concentration, resolution, and sensitivity, respectively. Figures of merit are measures of goal achievement for completely specified procedures that can be used for evaluating, selecting, and comparing analytical procedures. Other quantifiable factors that can be used to determine the applicability of analytical procedures to a particular problem include sensitivity, selectivity, speed of analysis, personnel requirements, and cost of the analytical procedure.

Grys [10] described five new functional concepts: accuracy, limit of detection, firmness, efficiency, and cost, which result in judgements of acceptibility of analytical methods. Accuracy is expressed in terms of recovery and reproducibility. The contribution due to recovery can be calculated by summing the percentage of different losses throughout the whole procedure. Reproducibility is expressed by the ratio of the full range to 100 times the ideal range.

Limit of detection is the concentration of a sample that gives a reading that is equal to twice the confidence half-interval for a series of ten determinations of the blank test value determined to 99% certainty. It is measured in mg per kg, or ppm, and is given by multiplying the standard deviation by 2.17, a factor that is determined from the $t$ test for $t_{0.01}$ and $n = 10$.

Firmness is an index of the effects of different factors upon the results. It is equal to the total deviation from expected values caused by the presence of equimolar amounts of interfering substances or connected with 5% changes in optimum reaction conditions such as acidity or reagent concentrations.

Efficiency provides information about the time consumption during the course of the whole procedure. It is expressed as the time of effective labor for one sample in minutes divided by 100.

Cost is a measure of the expenditure of materials and equipment used for the analysis of one sample by a new method in relation to the least expensive method. The cost of any desired method by which the analysis can be performed may be substituted for that of the least expensive method. It is given by dividing the ratio of the cost of the new method to the old method by 1000.

Eckschlager [11] discussed two informational variables, time-information performance and specific price of informa-

tion, which can be utilized in the evaluation and optimization of analytical methods. Time-information performance [12] can be rewritten as the ratio of information content to the time required for the analysis, including the analysis itself plus the time required to prepare the equipment for the analysis of the next sample. The time required for analysis can be partitioned into two segments, the basis time and the time required for the performance of $N$ parallel determinations. The specific price of the information [11] is defined as the ratio of the cost of the analysis to the amount of information obtained through simultaneous determination of $N$ components.

Danzer and Eckschlager [13] defined a general measure of information efficiency as the product of efficiency coefficients that are based upon the ratio of the value of the variable characterizing the properties required for the solution of particular analytical assignments to the actual value that the method provides. A ratio greater than one implies that more of the property is required than is provided by the method and the efficiency coefficient is assigned the value of zero. Otherwise, the efficiency coefficient is assigned the value of the ratio. They also defined a measure of information profitability as the ratio of the information efficiency to the specific price of the information. Results for the determination of manganese in low-alloy steels were tabulated for seven analytical methods: titrimetry, potentiometric titration, photometry, atomic absorption spectrophotometry, optical emission spectroscopy, optical emission spectrography, and optical emission spectrometry. The results demonstrated that information efficiency and information profitability are not always correlated. For example, although potentiometric titration has almost five times the information efficiency of titrimetry, both methods have the same information profitability when the duration of the analysis should be less than one day.

# Informing Power

Informing power is a measure of the amount of information available in a given analytical procedure. The concept was developed by Kaiser [14] with respect to spectrochemical methods of analysis eq (4) as a function of the resolving

$$P_{inf} = \int_{\nu_a}^{\nu_b} R(\nu) \, \log_2 S(\nu) \, \frac{d\nu}{\nu} \qquad (4)$$

power, $R(\nu)$, the maximal number of discernable steps for the amplitude, $S(\nu)$, and the spectral range, $\nu_a$ to $\nu_b$. If the resolving power and the maximal number of steps are fairly constant over the spectral range under consideration, informing power reduces to eq (5). For example, a grating

$$P_{inf} = R_{av} \, \log_2 S_{av} \, \ln(\nu_b/\nu_a) \qquad (5)$$

spectrograph system with a resolving power of $2 \times 10^5$, a spectral range from 2000 to 8000 Å, and 100 discernable steps in measured intensity levels at each wavelength would have an informing power of $2 \times 10^6$ bits. It is obvious that here the resolving power is the most important factor in maximizing informing power. In the case of a nondispersive, monochromatic method, resolving power between peaks at different wavelengths is not applicable and informing power is simply the $\log_2$ of the number of discernable amplitude steps at that wavelength. For example, 100 discernable intensity steps yields an informing power of 7 bits. The informing power for the corresponding polychromatic method is that for the monochromatic method multiplied by the number of frequencies. If the number of steps is different for each of the different frequencies, then the informing power is the same as for a collection of monochromatic methods, and the $\log_2$ of the number of steps is summed over each of the different frequencies.

Fitzgerald and Winefordner [15] extended the application of informing power to time-resolved spectrometric systems with the addition of a second resolving parameter, $R_t$. If both resolving powers and number of discernable steps are nearly constant over the range, then informing power reduces to eq (5) multiplied by $R_t \ln(t_2/t_1)$. For example, an atomic fluorescence spectrometer with an average resolving power of 3000 over a spectral range from 200 to 500 nm with an averge of 200 discernable intensity steps has an informing power of $5.7 \times 10^4$ bits. The informing power is increased to $9.7 \times 10^6$ bits for a range of $10^{-9}$ ($t_1$) to $10^{-6}$ sec with a measurement time limited by the lifetimes of the excited species of $10^{-7}$ sec ($t_2$) and a $\delta t$ of $10^{-9}$ sec ($R_t$ equals $(t_2 - t_1)/\delta t$). Comparisons of the informing power for a single beam molecular absorption spectrophotometer, normal molecular absorption phosphorimetry and time-resolved phosphorimetry showed an increase in the informing power by a factor of two for the normal phosphorimeter over the spectrophotometer. The addition of a time resolution element to phosphorimetry increased informing power by a factor of 450. The addition of a time resolution element to atomic fluorescence spectrometry increased the informing power by a factor of 170. Informing power was also used to compare analytical methods as well as to compare analytical instruments. A general photon counter was shown to have an informing power three times larger than that for an analog synchronous detection system.

Yost extended the application of informing power to tandem mass spectrometry [16,17], a method capable of generating enormous amounts of data. In the case of a quadrupole mass filter, the minimum resolution element, $\delta x$ is constant rather than the resolving power, $R(x)$. Informing power can then be calculated as shown in eq (6). A quad-

534

$$P_{\text{inf}} = \frac{1}{\delta x} \log_2 S(x) [x_b - x_a] \qquad (6)$$

rupole mass spectrometer with unit mass resolution, a mass range of 1000, and an ion intensity range of $2^{12}$ bits would have an informing power of $1.2 \times 10^4$ bits. The addition of another resolution element produces a double integral in the informing power equation that is equivalent to the the product of the informing power of the two elements. The addition of a capillary gas chromatograph with a nearly constant $10^5$ theoretical plate resolution over a one hour analysis time, results in $6.6 \times 10^6$ bits of informing power. The addition of a second quadrupole mass spectrometer with the same characteristics results in $1.2 \times 10^7$ bits of informing power. The combination of a capillary gc/ms/ms system results in an informing power of $6.6 \times 10^9$ bits, an increase by a factor of $5.5 \times 10^5$ over the original quadrupole mass spectrometer. The effect of experimental parameters on informing power was also demonstrated by Yost [16,17]. The variables associated with the collisionally activated dissociation process are potential resolution elements. Energy- and pressure- resolved ms/ms has an informing power of $3.6 \times 10^9$ bits.

The informing power metric can also be applied to electrochemistry. Using a current range of $-20$ to $+20$ $\mu$A that can be measured to within .005 $\mu$A yields $4 \times 10^3$ discernable steps. Eq (6) can be used to calculate the informing power for a cyclic staircase voltammetry (CSCV) experiment in which each current pulse is sampled and analyzed. An experiment with a staircase step of 13.5 mV ($\delta x$) and a potential range scanned from 0.0 to $-1.73$ V yields $3.1 \times 10^3$ bits of information.

Boudreau and Perone [18] demonstrated quantitative resolution in programmed potential step voltammetry for overlapped peaks with 30 mV separation between half wave potentials. If only resolved peaks are analyzed and the smallest resolution is 30 mV, the informing power is $1.4 \times 10^3$ bits. The addition of a time resolution element increases the informing power for electrochemical methods. Taking 45 equally spaced current measurements on each step at a sweep rate of 1.00 V/sec for the CSCV experiment increases the amount of information to $3.5 \times 10^7$ bits. The amount of information obtained from CSCV experiments can be easily manipulated by changing or adding resolution parameters.

Informing power can be used as a figure of merit for a completely specified method or system. Although the informing power of instrumental techniques may seem excessive when compared to the maximum information as calculated by Hartley's formula, it must be remembered that informing power is simply a measure of the maximal number of bits of information available in the procedure, not necessarily the useable or necessary amount of information.

Limitations in informing power arise from differences between practical and calculated resolving power. The lack of independence between the bits of information, noise, and interference result in the reduction of the useful informing power. Informing power can be partitioned into the amount of information required for the solution of the problem and the amount of redundant information required to provide a given level of confidence.

## Information Content

One of the most important concepts in information theory is that of informational gain or information content [19]. This is equal to the change in entropy due to the experiment and is quantified by the difference between the entropy using *a priori* probabilities and the entropy using *a posteriori* probabilities eq (7). The use of Shannon's

$$I(X|Y) = H(X) - H(X|Y) \qquad (7)$$

formula eq (1) to calculate the entropy does not guarantee a non-negative information content. However, another informational measure eq (8) always results in non-negative

$$I(X|Y) = -\sum_i^N p(x_i|Y) \log_2 [p(x_i|Y)/p(x_i)] \qquad (8)$$

values. For equal *a priori* probabilities, information content as calculated by Eqs (7,8) are equivalent. Since information content as discussed above can only be established after the analysis, these measures cannot be used as a quality criterion for selecting an analytical procedure. However, they can be used to evaluate the performance of a procedure.

Measures of information content has been applied to information theory models of structural analysis, qualitative analysis, quantitative analysis, trace analysis and instrumental analysis. Eckschlager and Stepanek have published a book [7] and a review article [8] on the application of information theory to analytical chemistry.

## Structural Analysis

One of the most difficult analytical tasks is the unambiguous determination of chemical structure. However, application of information theory to structural analysis is based upon a relatively simple entropy model [8] and an informational measure introduced by Brillouin [20]. The input consists of a finite number, $n_0$, of equally probable identities such as functional groups or conformational arrangements.

The output is a portion of a signal that corresponds to the identity, such as an IR band, NMR peak, or MS m/z peak, encoded only as to its presence or absence. The number of possible, but as yet undistinguished, structural arrangements is $n$, where $1 \leq n \leq n_0$ and $n = 1$ for an unambiguous determination of the structure for the analyte.

The uncertainty prior to analysis can be expressed by substituting the *a priori* probabilities into Shannon's equation eq (1). Since the *a priori* probabilities are qual, that is $1/n_0$, then the situation reduces to the case of maximum information (Hartley's equation) and the uncertainty is equal to $\log_2 n_0$. The uncertainty after analysis is equal to $\log_2 n$. The decrease in uncertainty due to the analysis corresponds to the informational gain eq (7) and is equal to $\log_2 (n_0/n)$. It assumes its maximum value in the case of the unambiguous determination of the structure and is equal to $\log_2 n_0$.

## Qualitative Analysis

The input for the model of qualitative analysis consists of a set of discrete identities, $X_i$, where $i = 1,2,. . .n_0$. If the output consists of a number of discrete, equally likely identities, the limiting case of Shannon's equation (Hartley's equation) can be used to calculate the entropy, and the information gained can be expressed by a ratio of the number of possible components before and after the analyses eq (9). [7]. For example, consider the case of an addition of

$$I(p,p_0) = \log_2 (n_0/n) \qquad (9)$$

HCl to a solution that contains only one of a list of 25 possible cations, three of which can be precipitated by HCl. The information gained as evidenced by a precipitate would be equal to $\log_2 (25/3)$ or 3 bits. The lack of a precipitate would imply an informational gain of only 0.2 bits. To consider various combinations of components, the total number of possible combinations is given by eq 10, where

$$n_0 = \sum_m^M n_m(M) = \sum_m^M [M!/(m!(M-m)!)] \qquad (10)$$

$M$ is the total number of components, and is divided into groups of $m$ components. In the case of a solution that contains from one to six cations, the total number of combinations is equal to 53. If two of them can be precipitated by HCl, there are 15 combinations in which neither cation is present and the information gained is 1.8 bits. For the 38 remaining cases in which either one or both cation is present, as evidenced by the appearance of a precipitate, an informational gain of 0.5 bit results.

In the case of instrumental or chromatographic qualitative or identification analyses, the output is a set of discrete signals in positions $Y_j$, where $j = 1,2,. . .m$ and $m \geq n_0$ [8]. The entropy can be expressed by Shannon's formula eq (1) and reaches a maximum when all of the possible identities, $X_i$, are equally likely. It is equal to zero when one identity is confirmed and the others are excluded, as would be the case for the *a posteriori* entropy for an unambiguous identification. The interpretation of these signals leads to an input-output relationship for the system that is represented by a set of a posteriori conditional probabilities, $p(x_i|y_j)$. The interpretation of these signals is also dependent upon preinformation represented by the *a priori* probabilities that may be calculated by Bayes theorem eq (11) [19].

$$p(x_i|Y) = [p(x_i)p(Y|x_i)]/\left[\sum_i p(x_i)p(Y|x_i)\right] \qquad (11)$$

The information content of an analytical signal is defined as the decrease in uncertainty eq (7). In the case of unambiguous determinations, $H(X|Y) = 0$ and $I(X|Y) = H(X)$, and $H(X)$ is also considered as the information required for unambiguous determination. However, most qualitative or identification procedures are chosen so as to minimize the uncertainty in identification for every possible signal. This is quantified by the informational measure of equivocation. Equivocation [8,19] is a measure of the expected or average value of the uncertainty after analysis eq (12). Equivocation

$$E = H(X|Y) = \sum_j p(y_j) H(X|y_j) \qquad (12)$$

and information content are complementary quantities, their sum equaling the entropy of the identification procedure. For an "ideal" procedure or an unambiguous determination, equivocation equals zero and information equals entropy.

Cleij and Dijkstra [21] demonstrated the use of information content and equivocation in the evaluation of thin-layer chromatographic procedures. Information content and equivocation were calculated for the identification of DDT and 12 related compounds for 33 different TLC systems. Calculations of the equivocation for the combinations of two TLC systems showed that the best combinations are not produced by combinations of the best individual TLC systems. This reflects the correlation between the best individual TLC systems.

Another method for quantifying information content is from the perspective of the possible signals rather than the possible identities [19]. Signal entropy, $H(Y)$, is the uncertainty in the identity of the unknown signal and can be quantified by substituting the probabilities of measuring the signals into Shannon's equation, eq (1). The conditional

536

entropy, $H(Y|x_i)$, is the uncertainty in the signal if the compound is known to be $x_i$. It can be considered a measure of noise and is given by substituting the conditional probabilities into Shannon's equation eq (1). Expected values for entropy and information content can be expressed in a manner analogous to that shown above. Also, since entropy is strong-additive [2,19], information content can be expressed in terms of the signal entropies.

Dupuis and coworkers [22,23] applied these methods to gas-liquid chromatography. The information content for 10 stationary phases used in gas-liquid chromatography was calculated on the basis of compound identification by retrieval of retention indices from a compiled library for a set of 248 compounds [22]. The information content per column ranged from 6.5 to 7.0 bits. The information content for combinations of columns is dependent upon both the number of columns and the sequence of columns. Ten sequences of the 10 columns yielded an information content of 43.3 bits. The study was expanded to include 16 gas-liquid chromatography stationary phases [23]. The complete data set of 248 compounds, a subset of 48 aliphatic alcohols, a subset of 35 aldehydes/ketones, and a subset of 60 esters were explored. For all four sets of compounds, combinations of stationary phases that yielded the most information consisted of one nonpolar phase plus one or more polar phases.

Van Marlen and Dijkstra [24] calculated the information content for the identification of binary coded mass spectra by retrieval and determined the optimal sequence of masses which contained the most information. A set of approximately 10,000 low resolution mass spectra were binary encoded using a threshold intensity level of 1% of the intensity of the base peak. Masses greater than 300 did not yield any additional information. Individually, masses of 300 or less contained from zero (m/z=3,4,5,6,7) to one (m/z=29,40,51,53,57,69,77) bit of information. The optimal mass sequence of 120 masses contained 40.9 bits of information, demonstrating the obvious redundancy in the binary coded spectra. The optimal mass sequence is dependent upon the distribution of the peaks. A set of 200 binary coded alkane spectra yielded 9 bits of information for 25 selected masses.

## Quantitative Analysis

The model for quantitative analysis is a two-stage model [8]. The input is a continuous distribution that produces a continuously variable signal. In the second stage, the signal specified by both position and intensity is decoded into results. The distribution of the results for parallel determinations is generally normal. Preinformation indicates that content of the component lies within a specified range of $x_0$ to $x_1$ so the a priori probability density is that of a rectangular

or uniform distribution. The information content is considered a divergence information measure that represents the error term in the measurement of the inaccuracy of the preinformation [7,8]. If the results confirm the a priori assumptions for the component, the information content is given by eq (13). The effect of a systematic error,

$$I(p,p_0) = \log_2 [(x_1 - x_0)/(\sigma\sqrt{2\pi e})] \qquad (13)$$

$\delta$, reduces the information content by factor of $\delta^2/2\sigma^2$ [7]. The use of parallel determinations, $n_p$, and the estimate of $\sigma$, s, results in eq (14), which utilizes the student's $t$ test at

$$I(p,p_0) = \log_2 [(x_1 - x_0)n_p/(2st(v)] \qquad (14)$$

the significance level of 0.038794. This level of significance is chosen so that twice the $t$ value at infinity equals $\sqrt{2\pi e}$ as the number of degrees of freedom approaches infinity. The information content as measured by eq (14) is the practical form of eq (13) since usually only the estimate of the standard deviation is known.

Poisson distributed results can be approximated by a normal distribution with the population mean, $\mu$, equal to the constant representing the average number of random points per unit time, $\lambda$, and the population standard deviation, $\sigma$, equal to $\sqrt{2}$. This changes the equation for information content to that shown in eq (15) [25]. However, this

$$I(p,p_0) = \log_2 [(x_1 - x_0)/(\sqrt{2\pi\mu e})] \qquad (15)$$

approximation is less valid for small values of lambda.

## Trace Analysis

The model for trace analysis [7,8] is essentially the same as for quantitative analysis except that the output signal is often barely distinguishable from the background noise. In the first case, the information content of the component to be determined is less than or equal to the detection limit of the analytical method and the only conclusion is that the content is somewhere between zero and the detection limit. The a posteriori probability distribution is equal to the inverse of the detection limit of the method. The information content is given as the $\log_2$ of the ratio of the highest estimated content of the component to the detection limit for the component. In the second case where the content of the component to be determined is greater than the detection limit, the content can be determined quantitatively. The a posteriori probability distribution is a shifted log-normal distribution. This information content differs from the information content of the first case by the addition of $\log_2$ $[\sqrt{n_p}/k\sigma 2\pi e]$ term, where $n_p$ is the number of parallel

determinations and $k$ is the asymmetry parameter for the shifted log-normal distribution. If the mean value for the determination is close to the detection limit, a truncated Gaussian distribution is used to describe the *a posteriori* distribution. The information content can be calculated as a function of the highest estimated value, the mean value, the standard deviation, the frequency, and the distribution function. The information content for both the log-normal and the truncated Gaussian distributions converge to $\log_2 [x_1/\sigma\sqrt{2\pi e}]$.

## Electrochemistry

Perone and coworkers have examined the effects of various experimental parameters on the qualitative information content of electrochemical data [26–31]. Because pattern recognition methods were used to obtain structural classification information empirically, an empirical measure of information gain was used to assess the effects of experimental parameters. This involved defining an appropriate figure of merit with which to measure the extent of informational goal achievement. Changes in information content are then determined by observing changes in attainment of the informational goal.

Byers, Freiser, and Perone [28,29] analyzed 45 compounds using cyclic staircase voltammetry. The data set consisted of 19 nitrobenzenes, 9 nitrodiphenyl ethers, and 17 ortho-hydroxy azo compounds. Of the nitrodiphenyl ethers, 4 were strong herbicides and 5 were either weak or nonherbicides. The informational goals for the problem were the classification of the 45 compounds by their structural type and the classification of the 9 nitrodiphenyl ethers according to herbicidal activity. Seven experimental variables, percent ethanol in the solvent, pH, surfactant concentration, number of cycles, scan rate, mercury drop hang time, and sampling time were varied in a fractional factorial design to generate a complete data base of collected cyclic staircase voltammograms and cyclic differential capacity curves [28] for subsequent analysis [29]. Faradaic and capacitive variable effect curves were calculated from the data. The average entropy for the three classes was 1.5 bits. The maximum entropy was 1.6 bits. The figure of merit for the informational goals was the percent correct classification as achieved by $k$-nearest neighbor analysis. For the structural characterization studies, the best overall percent classification ranged from 76% using capacitive variable effect features to 93% for both capacitive variable effect features and faradaic variable effect features. Overall accuracy for structural classification using the faradaic variable effect curve features ranged from 67% for percent ethanol to 93% for number of cycles. For the herbicidal prediction using variable effect curve features, the percent correct classification ranged from 78% for pH, faradaic sampling time, ca-

pacitive scan rate, and drop hang time to 100% for % ethanol, surfactant, faradaic number of cycles, and scan rate.

Barnes and Perone [30] studied the enhancement of chemical process information through experimental design. A simple model of controlled potential electrochemical processes based upon the Cottrell equation [31] was developed and implemented. The informational goal was to determine the effects of input voltage sequence, data collection fraction to analyze, and preprocessing scheme upon the determination of the diffusion coefficient. The figure of merit for goal achievement was based upon the least squares criterion function. Three input voltages, 180 mV step, pseudorandom binary sequence with peak voltage, $(E-E_0)$, of 180 mV, and white gaussian noise with mean of 90 mV and variance of 3 mV, were presented to the model. Comparisons of the identification results for the three inputs showed that the diffusion current model is fairly insensitive to the input sequence. Closer inspection of the model reveals that the anodic current is overwhelmed by the charging current. Therefore, the best input sequence for the model is that which is most easily generated.

In order to investigate the effects of timing, 3,000 data points corresponding to three milliseconds in time were generated using a step input. The least squares identifier was applied to 1 msec intervals of data both with and without the removal of charging current effects. When the charging current was present, the best results were obtained for the analysis of data taken after 10 time constants of the charging network. With the charging current removed, the best results were obtained with data taken within the first 1 msec interval when the amplitude of the anodic current and the signal to noise ratio are maximized. The best preprocessing scheme included filtering of unneccessary measurement components from the signal of interest, such as the removal of the charging current and signal averaging.

The application of information theory concepts to analytical chemistry can illuminate methods to increase the efficiency of chemical analysis. Early work shows encouraging promise for these types of applications. Optimum conditions have been established for obtaining structural, herbicidal activity, and diffusion coefficient information from voltammetric data. It hs been demonstrated that the informational goal(s) will dictate the most favorable choice of experimental conditions. The use of objective systematic information enhancement methods can highlight experimental parameters that are often traditionally overlooked.

## References

[1] Young, John F., Information Theory, Butterworth: London, (1971).
[2] Mathai, A.M., and P.N. Rathie, Basic Concepts in Information Theory and Statistics, John Wiley & Sons: New York, (1975).
[3] Guiasu, Silviu, Information Theory with Applications, McGraw-Hill

International Book Company: New York, (1977).

[4] Pierce, John R., An Introduction to Information Theory: Symbols, Signals, and Noise, Second Edition, Dover Publications, Inc.: New York, (1980).

[5] Shannon, C.E., Bell System Technical Journal, 27, 379 and 623 (1948).

[6] Malissa, Conveners H., and J. Rendl, Z. Anal. Chem., 272, 1 (1974) English version: I.L. Marr, Talanta, 22, 597 (1975).

[7] Eckschlager, Karel, and Vladimir Stepanek, "Information Theory as Applied to Chemical Analysis", John Wiley & Sons: New York, (1979).

[8] Eckschlager, Karel, and Vladimir Stepanek, Anal. Chem., 54(11), 1115 (1982).

[9] Ritter, S.L.; S.R. Lowery, H.B. Woodruff and T.L. Isenhour, Anal. Chem., 48(7), 1027 (1976).

[10] Grys, Stanislaw, Z. Anal. Chem., 273, 177 (1975).

[11] Eckschlager, Karel, Anal. Chem., 49(8), 1265 (1977).

[12] Danzer, K., Z. Chem., 15, 326 (1975).

[13] Danzer, Klaus, and Karel Eckschlager, Talanta, 25, 725 (1978).

[14] Kaiser, H., Anal. Chem., 42(2), 24A (1970).

[15] Fitzgerald, J.J., and J.D. Winefordner, Rev. Anal. Chem., 2(4), 229 (1975).

[16] Yost, R.A., Spectra, 9(4), 3 (1983).

[17] Fetterolf, D.D., and R.A. Yost, Int. J. Mass Spectrom. Ion Processes, 62, 33 (1984).

[18] Boudreau, P.A., and S.P. Perone, Anal. Chem., 51(7), 811 (1979).

[19] Cleij, P., and A. Dijkstra, Fresenius Z. Anal. Chem, 298, 97 (1979).

[20] Brillouin, L., "Science and Information Theory", Academic Press: New York, (1962).

[21] Cleij, P., and A. Dijkstra, Fresenius Z. Anal. Chem, 294, 361 (1979).

[22] Dupuis, Foppe, and Auke Dijkstra, Anal. Chem., 47(3), 379 (1975).

[23] Eskes Arie, Foppe Dupuis, Auke Dijkstra, Henri De Clercq, and Desire L. Massart, Anal. Chem., 47(13), 2168 (1975).

[24] van Marlen, Geert, and Auke Dijkstra, Anal. Chem., 48(3), 595 (1976).

[25] Eckschlager, K., Coll. Czech. Chem. Commun., 41, 2527 (1976).

[26] Perone, S.P., ACS Symposium Series, 265 (Comput. Lab.), 99 (1984).

[27] Burgard, D., and S.P. Perone, Anal. Chem., 50(9), 1366 (1978).

[28] Byers, W. Arthur, and S.P. Perone, Anal. Chem., 55(4), 615 (1983).

[29] Byers, W. Arthur, B.S. Freiser, and S.P. Perone, Anal. Chem., 55(4), 620 (1983).

[30] Barnes, Freddie, and S.P. Perone, unpublished results, personal communication from Freddie Barnes, September 1984.

[31] Bard, Allen J., and Larry R. Faulkner, "Electrochemical Methods, Fundamentals and Applications", John Wiley & Sons: New York, (1980).

# DISCUSSION

of the Perone-Ham paper, Measurement and Control of Information Content in Electrochemical Experiments

## Herman Chernoff

Statistics Center
Massachusetts Institute of Technology

The Shannon theory of information has had a profound impact in science and technology. Shannon defined information in terms of the reduction of uncertainty which, in turn, was measured by entropy. He was concerned mainly with the use of information to measure the ability to transmit data through noisy channels, i.e., channel capacity.

Statisticians have developed other, somewhat related, notions of information. In statistical theory, the major emphasis has been on how well experimental data help to achieve the goals in the classical statistical problems of estimation and hypothesis testing. These measures serve two useful functions. They serve to set a standard for methods of data analysis, methods whose efficiencies are measured in terms of the proportion of the available information that is effectively used. They also serve to design efficient experiments.

For the problem of estimation, Fisher introduced the Fisher Information which we now define. Suppose that it is desired to estimate a parameter $\theta$ using the result of an experiment which yields the data $X$ with the density $f(x|\theta)$. The *Fisher Information* for $\theta$ corresponding to $X$ is given by the matrix

$$J = I_X(\theta) = E_\theta(Y\ Y^T) \tag{1}$$

where $Y$ is the *score function* defined by

$$Y = Y(X, \theta) = \frac{\partial[\log f_X(x|\theta)]}{\partial\theta}\bigg|_{x=X} \tag{2}$$

If $\theta$ is a multidimensional vector, $J$ is a nonnegative definite

539

symmetric matrix with the additive property

$$I_{X,Z}(\theta)=I_X(\theta)+I_Z(\theta) \qquad (3)$$

if $X$ and $Z$ are independent. As a consequence

$$I_{X,X,\ldots,X}(\theta)=nI_X(\theta)=nJ \qquad (4)$$

if the left subscript refers to the independent replicaton of $X$, $n$ times. For such an experiment, it has been shown, under mild regularity conditions, that the Maximum Likelihood Estimate (MLE) $\hat{\theta}$ will be approximately normally distributed with mean $\theta$ and covariance matrix $J^{-1}/n$ for large $n$. Moreover, the Cramér-Rao theorem states that one cannot expect to find a reasonable estimate that does better.

Some implications of the above paragraph are illustrated by three simple examples below.

### Example 1. *Mean of a Normal Distribution.*

Let $X$ be normally distributed with mean $\theta$ and known variance $\sigma^2$, and let $X_1,X_2,\ldots,X_n$ be a sample of $n$ independent observations on $X$. It is easy to show that $J_X=\sigma^{-2}$ and that the MLE $T_1=\bar{X}=n^{-1}(X_1+\ldots+X_n)$ is normally distributed with mean $\theta$ and variance $\sigma^2/n$. However, the statistician, who fears for outliers and may wish to use a more robust estimator than the sample mean, may prefer to use $T_2$, the sample median. It can be shown that $T_2$ is approximately normally distributed with mean $\theta$ and variance $\pi\sigma^2/2n$ for large $n$. The equation

$$\frac{\sigma^2}{n_1}=\frac{\sigma^2}{n_2}\frac{\pi}{2} \qquad (5)$$

implies $n_1/n_2=2/\pi=.64$ which is a natural measure of the efficiency of $T_2$, indicating that, with $T_1$, we need only 64% of the data to achieve the same accuracy as with $T_2$. If the effective waste of 36% of the data seems excessive, the statistician can improve on efficiency with little sacrifice of robustness, e.g., by using the upper and lower quartiles as well as the median, or by using trimmed means.

### Example 2. *Experiments With Information Matrices.*

Let $\theta=(\theta_1,\theta_2)^T$, and let $X$ and $Z$ be two experiments with information matrices

$$J_X=\begin{Vmatrix} 4 & 3 \\ 3 & 4 \end{Vmatrix} \text{ and } J_Z=\begin{Vmatrix} 4 & -3 \\ -3 & 4 \end{Vmatrix}.$$

It is desired to estimate $\theta_1$ using replications of either $(X,X)$ or $(Z,Z)$ or $(X,Z)$. Let $J^{11}$ be the upper left member of $J^{-1}$ which measures the (asymptotic) variance of $\hat{\theta}_1$, the MLE of $\theta_1$. Then

$$J_{XX}=\begin{Vmatrix} 8 & 6 \\ 6 & 8 \end{Vmatrix}, \qquad J_{XX}^{11}=0.286,$$

$$J_{ZZ}=\begin{Vmatrix} 8 & -6 \\ -6 & 8 \end{Vmatrix}, \qquad J_{ZZ}^{11}=0.286,$$

and

$$J_{XZ}=\begin{Vmatrix} 8 & 0 \\ 0 & 8 \end{Vmatrix}, \qquad J_{XZ}^{11}=0.125.$$

This clearly indicates, that in the presence of *nuisance parameters* such as $\theta_2$, one may squeeze much more useful information out of a combination of two equally informative experiments than by repeating one of these two or, in this case, even four times.

### Example 3. *Estimate Safe Dose Level in Probit Model.*

For an experiment at does level $d$, the probit model attributes the probability of a response to be

$$p(d,\theta)=\Phi[(d-\mu)/\sigma] \qquad (6)$$

where $\theta=(\mu,\sigma)^T$, $\Phi$ is the standard normal cumulative distribution and the "safe" dose level to be estimated is defined as $\mu-2.87\sigma$. If one is permitted to select a sequence of $n$ dose level, $d_1,d_2,\ldots,d_n$ with which to challenge $n$ subjects, the optimal choice or *design*, for estimating $\mu-2.87\sigma$ can be shown to assign about 23% of the doses at level $d=\mu+1.57\sigma$ and the remaining 77% of the doses at level $d=\mu-1.57\sigma$.

This optimal design illustrates several points.

1. This design is *locally optimal*, i.e., it requires a knowledge of $\theta$ to provide the best estimate of a function of $\theta$. Superficially, it seems silly, for if we knew $\theta$, we would not need to estimate it. In fact, it indicates that as data cumulates, one knows more about $\theta$ and can sequentially use that information to provide improved experiments.

2. In this experiment, the repeated use of one dose level $d_o$ would provide only an estimate of the function $p(d_o,\theta)$ and would yield no other useful information about $\theta$ or $\mu-2.87\sigma$. At least two dose levels are required. What is .somewhat surprising is that no more than two dose levels are required for an optimal design. A more general theorem states that if it is desired to estimate $r$ functions of $k$ parameters upon which the distribution of the data depend, then an optimal design can be constructed using at most $k+(k-1)+\ldots+(k-r+1)$ of the available (elementary) experiments.

3. The optimal design is not necessarily a practical one. Most investigators would be interested in using a variety of dose levels as a means of checking the basic model. Theory permits us to measure the loss of information inherent in the use of practical, but suboptimal designs, so that one can

decide on whether the loss is so extravagant that other alternatives should be considered.

We mention briefly that for testing hypotheses, there are several measures of information which are of potential use, depending on the type of problem. Perhaps the most useful measure is the *Kullback-Leibler Information* (KL)

$$I_X^*(\theta,\phi) = \int f(x|\theta) \, \log \left[\frac{f(x|\theta)}{f(x|\phi)}\right] dx \qquad (7)$$

which measures two important aspects of the ability to use a sample of $n$ observations on $X$ with distribution $f(x)$, to discriminate between the hypotheses $H_1:f(x)=f(x|\theta)$ and $H_2:f(x)=f(x|\phi)$.

The Kullback-Leibler Information is additive as is the Fisher Information but it is not symmetric since, $I_X^*(\theta,\phi)$ is not generally equal to $I_X^*(\phi,\theta)$. For large samples, it is possible to find tests which, for fixed type 1 error probability $\alpha=P(\text{reject } H_1|H_1)$, have the type 2 error probability $\beta=P(\text{accept } H_1|H_2)$ approach 0 at a rate determined by $I^*$. We have, roughly

$$\beta^* \sim e^{-nI_X^*(\theta,\phi)} \qquad (8)$$

Another property of KL is that for optimal sequential testing as the cost $c$, per observation, approaches zero, the expected costs $R(\theta)$ and $R(\phi)$ associated with the sequential procedure when $H_1$ and $H_2$ are true, satisfies

$$R(\theta) \approx -c \, \log \, c/I_X^*(\theta,\phi)$$

$$R(\phi) \approx -c \, \log \, c/I_X^*(\phi,\theta) . \qquad (9)$$

This implies that if we suspect $H_1$ is true, we should select the experiment which maximizes $I^*(\theta,\phi)$ and if we suspect that $H_2$ is true, we should maximize $I^*(\phi,\theta)$. Here again, as in the estimation problem, we are in a position to improve the experimental design as information cumulates, and our belief in $H_1$ or $H_2$ increases.

To return to chemical experimentation, one should point out that an experimental set up which yields vast amounts of bits of information is not very useful if the analysis of the data does not make efficient use of the data. To discriminate between two alternatives requires only one bit of *effective* information in the Shannon sense. The choice between experiments which yield 1,000 and 10,000 bits must involve how much effective information is readily available from the analysis.

Some bibliography on the uses of information in statistics is contained in Chernoff (1972).

[1] Chernoff, H., *Sequential Analysis and Optimal Design* SIAM monograph 8, SIAM, Philadelphia (1972).

# Pattern Recognition Studies of Complex Chromatographic Data Sets

## P. C. Jurs, B. K. Lavine, and T. R. Stouch

### The Pennsylvania State University, University Park, PA 16802

Chromatographic fingerprinting of complex biological samples is an active research area with a large and growing literature. Multivariate statistical and pattern recognition techniques can be effective methods for the analyisis of such complex data. However, the classification of complex samples on the basis of their chromatographic profiles is complicated by two factors: 1) confounding of the desired group information by experimental variables or other systematic variations, and 2) random or chance classification effects with linear discriminants. We will treat several current projects involving these effects and methods for dealing with the effects.

Complex chromatographic data sets often contain information dependent on experimental variables as well as information which differentiates between classes. The existence of these types of complicating relationships is an innate part of fingerprint-type data. ADAPT, an interactive computer software system, has the clustering, mapping, and statistical tools necessary to identify and study these effects in realistically large data sets.

In one study, pattern recognition analysis of 144 pyrochromatograms (PyGCs) from cultured skin fibroblasts was used to differentiate cystic fibrosis carriers from presumed normal donors. Several experimental variables (donor gender, chromatographic column number, etc.) were involved in relationships that had to be separated from the sought relationships. Notwithstanding these effects, discriminants were developed from the chromatographic peaks that assigned a given PyGC to its respective class (CF carrier vs normal) largely on the basis of the desired pathological difference. In another study, gas chromatographic profiles of cuticular hydrocarbon extracts obtained from 179 fire ants were analyzed using pattern recognition methods to seek relations with social caste and colony. Confounding relationships were studied by logistic regression. The data analysis techniques used in these two example studies will be presented.

Previously, Monte Carlo simulation studies were carried out to assess the probability of chance classification for nonparametric and parametric linear discriminants. The level of expected chance classification as a function of the number of observations, the dimensionality, and the class membership distributions were examined. These simulation studies established limits on the approaches that can be taken with real data sets so that chance classifications are improbable.

Key words: classification effects; multicomponent spectra; pattern recognition.

Profiling of complex biological materials with high performance chromatographic methods is an active research area with a large and growing literature, e.g., [1–10][1]. Such chromatographic experiments often yield chemical profiles containing hundreds of constituents. These chromatograms can be viewed as chemical fingerprints of the complex samples. Objective analysis of the profiles depends upon the use of multivariate statistical methods. In this regard pattern recognition techniques have been found to be of utility.

Pattern recognition methods have been used to distinguish between individuals in a particular diseased state and normal individuals [7–10]. These methods attempt to classify a sample according to a specific property (e.g., diabetic vs normal) by using measurements that are indirectly related to that property. Mea-

[1] Bracketed numbers indicate literature references.

surements related to the property in question are made. An empirical relationship is then derived from a set of data for which the property of interest and the measurements are known (a training set). Such a relationship or classification rule may be used to infer the presence or absence of this property in objects that are not part of the original training set.

For pattern recognition analysis, each chromatogram is represented as a point, $X = (x_1, x_2, x_3, ..., x_d)$ where component $x_j$ is the area of the $j$th peak. A set of chromatograms is represented by a set of points in a $d$-dimensional Euclidean space. The expectation is that the points representing chromatograms from one class will cluster in one limited region of the space separate from the points corresponding to the other class. Pattern recognition is a set of methods for investigating data represented in this manner in order to assess the degree of clustering and general structure of the data space. The four main subdivisions of pattern recognition methodology are mapping and display, discriminant development, clustering, and modelling [11-14]. The ADAPT computer software system [15] has routines in all these areas, and many were used in the two example studies below.

An assumption in pattern recognition is that the ability to categorize the data into the proper classes is meaningful. Successful classification is thought to imply that a relationship between the measurements or features and the property of interest exists. However, classification based on random or chance separation can be a serious problem. For example, the probability of fortuitously obtaining 100% correct classification for a two class problem using a nonparametric linear discriminant can be calculated from the following equation

$$P = 2 \sum_{i=0}^{d} C_i^{n-1} / 2^n \qquad (1)$$

where $C_i^{n-1} = (n-1)!/[(n-1-i)!i!]$, $n$ is the number of objects in the data set, and $d$ is the dimensionality or number of descriptors per object [16,17]. Figure 1 shows a plot of $P$ versus the ratio of the number of objects to the number of descriptors per object ($n/d$) for $n = 50$. The only assumption made concerning the data is that it be in general position, that is, none of the $d+1$ data points should be contained in a $(d-1)$-dimensional hyperplane. When $n/d$ is large, the probability of achieving complete separation due to chance is small. As the number of descriptors approaches the number of objects used in the study, the probability of such an occurence increases. When $n/d = 2$, the probability of complete separation is one-half. Such classifications arise due to chance and are not due to any relationship between the objects in the data set. A linear discriminant
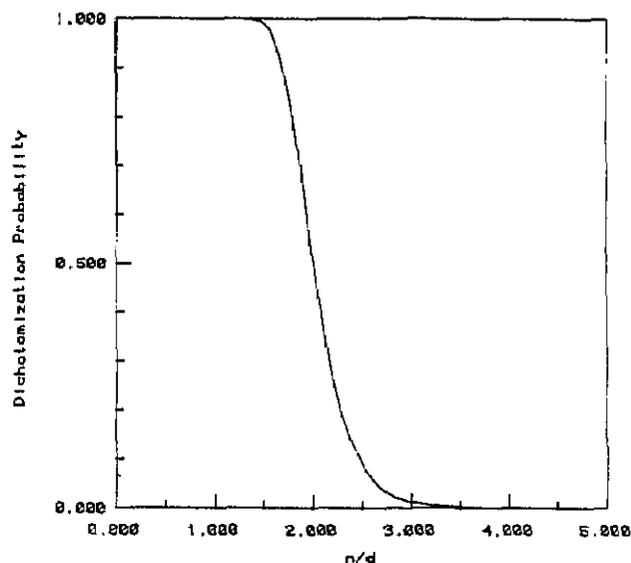


Figure 1-The probability of complete separation into classes by a nonparametric linear discriminant function versus the ratio of the number of objects to the number of descriptors per object.

function developed with an inappropriately small $n/d$ will probably have no predictive ability beyond random guessing.

If $n/d > 3$, the probability of complete separation due to chance is small [18, 19]. However, classification rules using linear discriminants are often developed using training sets that are not completely linearly separable. Recently, Stouch and Jurs have reported Monte Carlo simulation studies [20] assessing the degree of fortuitous classification for such situations. Figure 2 is a plot of results obtained in hundreds of Monte Carlo experi-
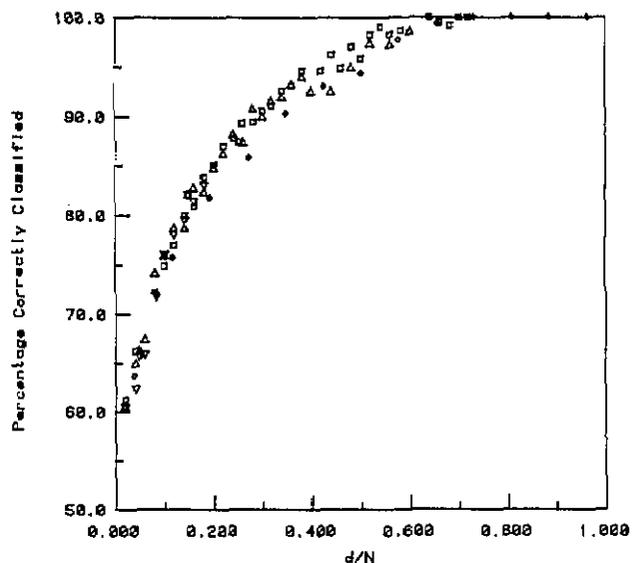


Figure 2-Plot of the percentage of correctly classified patterns versus the ratio of the number of descriptors per pattern to the number of patterns. Each plotting character represents the mean of a number of Monte Carlo experiments.

544

ments. It shows the percentage of objects correctly classified versus the $d/n$ ratio. The patterns used to develop this curve were random, and equal class sizes were used. The percentage correctly classified for a given $d/n$ value can only be due to chance. Although the probability of obtaining 100% correct classification for $n/d > 3$ is small, chance classification success rates range between 85% and 95%. The influence of the class membership distribution upon chance classification was also investigated, and unequal class sizes lead to even higher success rates due to chance. Figure 3 shows the cumulative probability of achieving any degree of separation due to chance for evenly-divided classes for three values of $n/d$. At $n/d = 5$, the probability is 50% that 77% of the objects will be correctly classified due to chance. Chance classifications can be a serious problem in linear discriminant analysis of chromatographic fingerprint data. Hence, the results obtained with real data sets must be compared to the results achievable by chance in order to assure that meaningful relations have been discovered.

A second complicating aspect of the classification of complex samples on the basis of their chromatographic profiles is the confounding of the desired group information by experimental variables or other systematic variations. If the basis of classification for patterns in the training set is other than the desired group difference, unfavorable classification results for the prediction set will be obtained despite a linearly separable training set. The existence of these types of complicating relationships is an inherent part of fingerprint-type data. We will discuss several current projects involving these effects and methods for dealing with them.
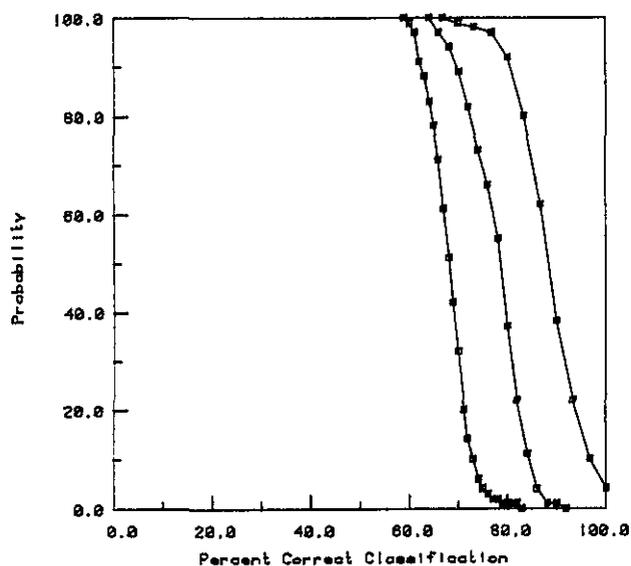


Figure 3–Plot of the cumulative probability of achieving any degree of separation due to chance versus that degree of separation. Three values of $n/d$ are shown. Evenly divided training sets were used.

## Cystic Fibrosis Heterozygotes vs Normal Subjects

The first study involves the application of pyrolysis gas chromatography (PyGC) and pattern recognition methods to the problem of identifying carriers of the cystic fibrosis (CF) defect [21]. The biological samples used in this experiment were cultured skin fibroblasts grown from 24 samples obtained from parents of children with CF and from 24 presumed normal donors. A typical CF heterozygote pyrochromatogram is shown in figure 4. The pyrolysed fibroblasts were analyzed on fused silica capillary columns with temperature programming. For each subject, triplicate pyrochromatograms were taken.

The 144 pyrochromatograms were standardized using an interactive computer program [22]. Each pyrochromatogram was divided into 12 intervals defined by 13 peaks that were always present. The retention times of the peaks within the intervals were scaled linearly for best fit with respect to a reference pyrochromatogram. This peak matching procedure yielded 214 standardized retention time windows. Each pyrochromatogram was also normalized using the total area of the 214 peaks. This set of chromatographic data—144 PyGCs of 214 peaks each—was autoscaled so that each PyGC peak had a mean of zero and a standard deviation of one within the entire set of pyrochromatograms.

To apply pattern recognition methods to this overdetermined data set, the necessary first step was feature selection. The number of peaks per chromatogram must be reduced to at least one-third the number of independent PyGCs in the data set, so at most 16 peaks could be analyzed at one time. For the final results of the analysis to be meaningful, this feature selection must be done objectively, that is, without using any class membership information.

For experiments of the type that we are considering here it is inevitable that there will be relationships between sets of conditions used in generating the data and patterns that result. One must realize this in advance when approaching the task of analyzing such data. One must isolate the information pertinent to the pathological alteration characteristic of CF heterozygotes from the large amount of qualitative and quantitative data due to experimental conditions that is also contained in the complex capillary pyrochromatograms.

We have observed that experimental variables (cell culture, batch number, passage number, donor gender, and column identity) can contribute to the overall classification process. For example, a decision function or classification rule was developed from the 12 peaks
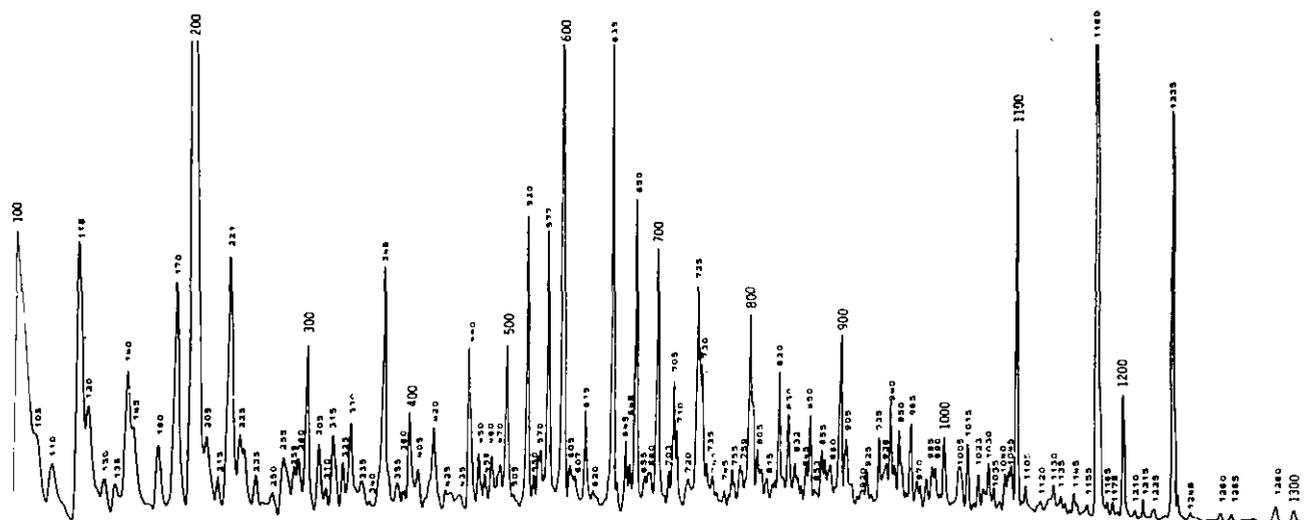
545

Figure 4–A representative pyrochromatogram from the CF study. The peak identities are those assigned using the peak-matching software. The major peaks are those with assignments that are multiples of 100.

comprising interval three. The CF PyGCs were linearly separable from the PyGCs of the presumed normal donors. However, when the points from this 12-dimensional space were mapped onto a plane that best represents the pattern space (the plane defined by the two largest principal components), groupings related to column identity were observed. Furthermore, classifiers could be developed from these 12 peaks that yielded favorable classification results for many of the experimental variables.

Notwithstanding the effects of the experimental variables described above, a discriminant or decision function has been developed from the PyGC peaks that separates the pyrochromatograms of CF heterozygotes from those of presumed normal subjects, by and large, on the basis of valid chemical differences. The development of such discriminant is described in detail below.

The 65 peaks that were present in at least 90% of the PyGCs were used as a starting point for the analysis. We assessed the ability of each of these 65 peaks alone to discriminate between PyGCs with respect to gender, passage number, and column identity. Twelve peaks that had larger classification success rates for the CF vs normal than for any other dichotomy were selected for further analysis. This procedure identifies those peaks that contain the most information about CF vs normal as opposed to the experimental variables. We were attempting to simultaneously minimize both the probability of chance separation and that of confounding with unwanted experimental details. A classification rule developed from these 12 peaks using the $k$-nearest neighbor procedure correctly classified 90% of the PyGCs in the data set. Variance feature selection [23], combined

with the linear learning machine and the adaptive least-squares methods [24], was used to remove 6 of the 12 peaks found to be least relevant to the classification problem. A discriminant that misclassified only eight of the pyrochromatograms (136 correct of 144, 94%) was developed using the final set of only six peaks.

The contribution of the experimental parameters to the overall dichotomization power of the decision function based on the six peaks was assessed by reordering experiments. The set of PyGCs was first reordered in terms of donor gender, and classification results indistinguishable from random were obtained. Similar studies were done for passage number and column identity, and comparable results were obtained. The results of the reordering tests suggest that the decision function based on the six PyGC peaks incorporates mainly chemical information to separate the pyrochromatograms of the CF heterozygotes from those of the normals.

The ability of the decision function to classify a simulated unknown sample was tested using a procedure known as internal validation. Twelve sets of pyrochromatograms were developed by random selection where the training set contained 44 triplicates and the validation set contained the remaining 4 triplicates. Any particular triplicate was only present in one validation set of the 12 generated. Discriminants developed for the training sets were tested on the PyGCs that were held out. The average correct classification for the held-out pyrochromatograms was 87%. This same internal validation test was repeated except that members of the held-out sets included triplicate samples analyzed on the same column or grown in the same batch of growth medium. The average correct classification for the held-

546

out pyrochromatograms in this set of runs was 82%. Although the classification success rate of the decision function was diminished when we took into account these confounding effects, favorable results were still obtained.

## Recognition of Ants by Caste and Colony

Chemical communication among social insects can be studied with chromatographic methods. For example, evidence regarding the role of cuticular hydrocarbons in nestmate recognition came from a study of the Myrmecophilous beetle [25]. The data generated in such studies can be complex and may require multivariate statistical or pattern recognition methods for interpretation. Presently, we are analyzing gas chromatographic profiles of high molecular weight hydrocarbon extracts obtained from the cuticles of 179 red fire ant (*Solenopsin invicta*) samples. We are using pattern recognition methods to seek relations with social caste and colony. Each sample contains the hydrocarbons extracted with hexane from the cuticles of 100 individual ants. The hydrocarbon fraction analyzed by gas chromatography was isolated from the concentrated hexane washings by means of a silicic acid column. Evidence regarding the role of cuticular hydrocarbons in nestmate recognition came from a study of the Myrmecophilous beetle [25]. A gas chromatographic trace of the cuticular hydrocarbons from a *S. invicta* sample is shown in figure 5. The hydrocarbon extract was analyzed on a glass column packed with 3% OV-17 using temperature programming.

Five major hydrocarbon compounds were identified and quantified by GC/MS analysis: heptacosane $(n-C_{27}H_{56})$, 13-methylheptacosane, 13,15-dimethylheptacosane, 3-methylheptacosane, and 3,9-dimethyl-
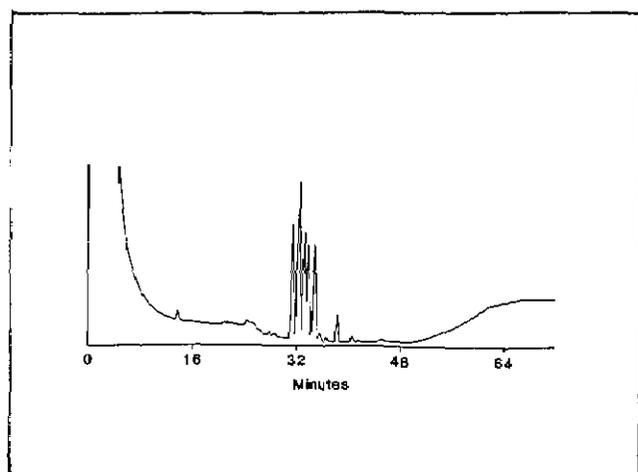
heptacosane in the order of elution from the OV-17 column used. An internal standard was used for quantification. Each chromatogram was normalized using the weight of the collected ants.

Several questions have been addressed in this study: 1) Are the hydrocarbon patterns characteristic of individual colonies? 2) Does the overall colony hydrocarbon pattern change with time? 3) Are the hydrocarbon patterns significantly different for the social castes? In this study, ant samples were obtained from five different colonies (E, J, P, Q, R), three different castes (foragers, broods, and reserves), and for four different time periods (the first three in spring and summer and time period four in the winter).

The first step was to use mapping and display methods [12,17] to examine the structure of the data set. Methods used included principal components mapping and nonlinear mapping [14]. In figures 6 and 7 the results of principal component mapping experiments for colonies J and Q are shown. Colony J includes samples from time periods one through three, whereas colony Q is represented by ants from all four time periods. Colony J has 9 and colony Q has 12 members from each social caste. Pattern groupings according to time period and caste can be seen in figures 6 and 7. The first two principal components account for 96.2% and 97% respectively of the total cumulative variance in the two plots shown. Mapping experiments of this nature were also carried out for samples from a particular caste or time period, and pattern groupings with respect to colony identity, social caste, and temporal period were observed.
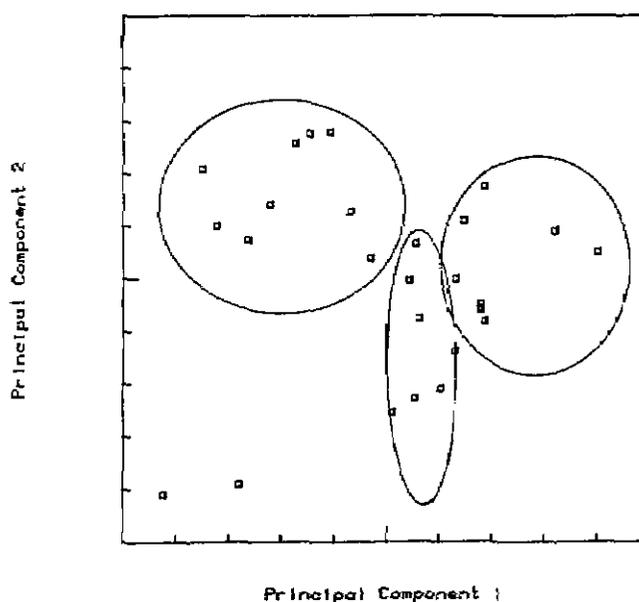


Figure 5–Gas chromatographic trace of cuticular hydrocarbons from *S. invicta* (Reprinted with permission from ref. [25]).



Figure 6–Plot of the two principal components of the five GC peaks for colony J. The elipses show groupings of samples by time period.
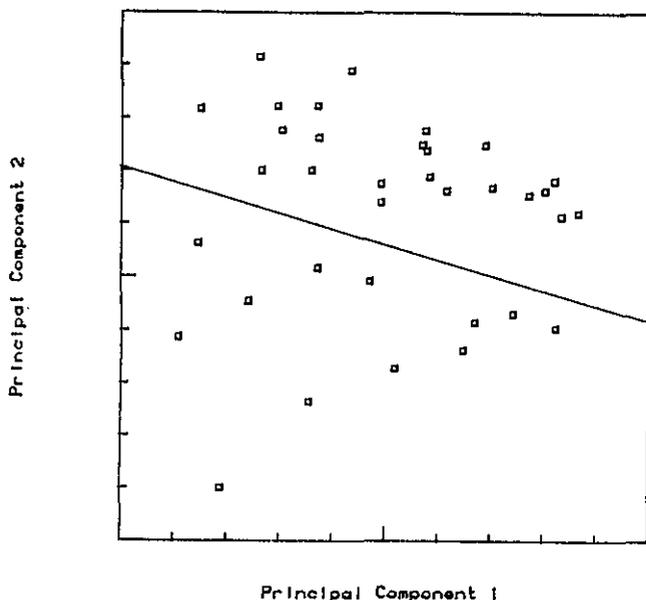
Figure 7–Plot of the two principal components of the five GC peaks for colony Q. The foragers are separated from the reserves and broods by the linear discriminant.

**Table 1.** Percentage of chromatograms correctly classified by colony for several two-way classifications.

| Colony | No. in Colony | Colony in Second Group | | | | | |
|--------|---------------|---|---|---|---|---|---|
|        |               | E | J | P | Q | R | All |
| E | 36 | — | 98 | 100 | 100 | 100 | 95 |
| J | 27 |   | — | 100 | 100 | 100 | 98 |
| P | 36 |   |   | — | 100 | 100 | 99 |
| Q | 35 |   |   |   | — | 73 | 85 |
| R | 36 |   |   |   |   | — | 82 |

gression have been employed in this study. The results obtained using these techniques support the conclusions drawn from the pattern recognition experiments. In summary, the GC traces representing ant cuticle extracts could be related to colony identity, social caste, and time period using pattern recognition methods.

Discriminant analysis studies were also performed. In one study the data set was divided into three categories according to the social caste of the pooled ant sample. Linear discriminants were developed using the areas of the five GC peaks. The hydrocarbon patterns of the foragers were found to be very different from the broods and reserves. In fact, information necessary to discriminate foragers from broods and reserves was primarily encoded in the concentration pattern of the first GC peak. A similar study was undertaken for time period, and the fourth time period was found to be very different from time periods one, two, and three. During time period four the ants are in a state of hibernation, whereas time periods one, two, and three correspond to the spring and summer months.

The hydrocarbon profiles were also found to be characteristic of the individual colonies. Linear decision surfaces were developed from the five GC peaks, using an iterative least-squares method. The purpose was to separate one colony from another or one colony from all other colonies. The results of these discriminant analysis experiments are summarized in table 1. The first row of the table shows that colony E could be separated from colony J by a discriminant that achieved 98% correct classifications (63 correct out of 64 samples) and that colony E could be separated from all the remaining colonies by a discriminant that achieved 95% correct classifications (162 correct out of 170). Colonies Q and R could not be separated well by this method. In addition, multivariate statistical methods such as multivariate analysis of variance and stepwise logistic re-

## References

[1] Zlatkis A.; R. S. Brazell and C. F. Poole, The role of organic volatile profiles in clinical diagnosis, Clin. Chem. 27, 789–797 (1981).

[2] Jellum, E. J., Profiling of human body fluids in healthy and diseased states using gas chromatography and mass spectrometry, with special reference to organic acids, Jour. Chromatog. 143, 427–462 (1977).

[3] Reiner, E., and F. L. Bayer, Botulism: a pyrolysis-gas-liquid chromatographic study, Jour. Chrom. Sci. 16, 623–629 (1978).

[4] Reiner, E., and J. J. Hicks, Differentiation of normal and pathological cells by pyrolysis-GLC, Chromatographia 5, 525–528 (1972).

[5] Jellum, E.; I. Bjoernson, R. Nesbakken, E. Johansson, and S. Wold, Classification of human cancer cells by means of capillary gas chromatography and pattern recognition analysis, Jour. Chromatog. 217, 231–237 (1981).

[6] Soderstrom, B.; W. Wold and G. Blomquist, Pyrolysis-gas chromatography combined with SIMCA pattern recognition for classification of fruit-bodies of some ectomycorrhizal suillus species, Jour. Gen. Micro. 128, 1773–1784 (1982).

[7] McConnell, M. L.; G. Rhodes, U. Watson, and M. Novotny, Application of pattern recognition and feature extraction techniques to volatile constitutent metabolic profiles obtained by capillary gas chromatography, Jour. Chromatog. 162, 495–506 (1979).

[8] Wold, S.; E. Johansson, E. Jellum, I. Bjoernson, and R. Nesbakken, Application of SIMCA multivariate data analysis to the classification of gas chromatographic profiles of human brain tissues, Anal. Chim. Acta 133, 251–259 (1981).

[9] Scoble, H. A.; J. L. Fashing and P. R. Brown, Chemometrics and liquid chromatography in the study of acute lymphocytic leukemia, Anal. Chim. Acta 150, 171–181 (1983).

[10] Rhodes, G.; M. Miller, M. L. McConnell, and M. Novotny, Metabolic abnormalities associated with diabetes mellitus, as investigated by gas chromatography and pattern recognition analysis of profiles of volatile metabolites, Clin. Chem. **27**, 580–585 (1981).

[11] Jurs, P. C., and T. L. Isenhour, Chemical applications of pattern recognition, Wiley-Interscience: New York (1975).

[12] Varmuza, K., Pattern recognition in chemistry, Springer-Verlag: Berlin (1980).

[13] Kryger, L., Interpretation of analytical chemical information by pattern recognition methods—a survey, Talanta **28**, 871–887 (1981).

[14] Fukunaga, K., Introduction to statistical pattern recognition, Academic Press: New York (1972).

[15] Stuper, A. J.; W. E. Brugger and P. C. Jurs, Computer assisted studies of chemical structure and biological function, Wiley-Interscience: New York (1979).

[16] Nilsson, N. J., Learning machines, McGraw-Hill: New York (1965).

[17] Tou, J. T., and R. C. Gonzalez, Pattern recognition principles, Addison-Wesley Pub. Co.: Reading, MA (1974).

[18] Stuper, A. J., and P. C. Jurs, Reliability of nonparametric linear classifiers, Jour. Chem. Inf. Comp. Sci. **16**, 238–241 (1976).

[19] Whalen-Pedersen, E. K., and P. C. Jurs, The probability of dichotomization by a binary linear classifier as a function of training set population distribution, Jour. Chem. Inf. Comp. Sci. **19**, 264–266 (1979).

[20] Stouch, T. R., and P. C. Jurs, Monte Carlo studies of the classification made by *nonparametric linear discriminant functions*, Jour. Chem. Inf. Comp. Sci. **25**, 45–50 (1985).

[21] Pino, J. A.; J. E. McMurry, P. C. Jurs, B. K. Lavine, and A. M. Harper, Application of pyrolysis/gas chromatography/pattern recognition to the detection of cystic fibrosis heterozygotes, Anal. Chem. **57**, 295–302 (1985).

[22] Pino, J. A., Pyrochromatography of human skin fibroblasts: normal subjects vs cystic fibrosis heterozygotes, Ph.D. Thesis, Cornell University (1984).

[23] Zander, G. S.; A. J. Stuper and P. C. Jurs, Nonparametric feature selection in pattern recognition applied to chemical problems, Anal. Chem. **47**, 1085–1093 (1975).

[24] Moriguchi, I.; K. Komatsu and Y. Matsushita, Adaptive least-squares method applied to structure-activity correlation of hypotensive N-Alkyl-N″-cyano-N′-pyridylguanidines, Jour. Med. Chem. **23**, 20–26 (1980).

[25] Vander Meer, R. K., and D. P. Wojcik, Chemical mimicry in the myrmecophilous beetle *myrmecaphodius excavaticollis*, Science **218**, 806–808 (1982).