

Statistical Data Analysis in the Computer Age

BRADLEY EFRON AND ROBERT TIBSHIRANI

Most of our familiar statistical methods, such as hypothesis testing, linear regression, analysis of variance, and maximum likelihood estimation, were designed to be implemented on mechanical calculators. Modern electronic computation has encouraged a host of new statistical methods that require fewer distributional assumptions than their predecessors and can be applied to more complicated statistical estimators. These methods allow the scientist to explore and describe data and draw valid statistical inferences without the usual concerns for mathematical tractability. This is possible because traditional methods of mathematical analysis are replaced by specially constructed computer algorithms. Mathematics has not disappeared from statistical theory. It is the main method for deciding which algorithms are correct and efficient tools for automating statistical inference.

MOST SCIENTISTS FACE PROBLEMS OF DATA ANALYSIS: What data should I collect? What can I conclude from my data? How far can I trust the conclusions? Statistics is the mathematical science that deals with these questions. Some statistical methods, such as linear regression, hypothesis testing, standard errors, and confidence intervals, have become familiar in the scientific literature over time. Most of the "classical" methods were developed between 1920 and 1950, by scientists such as R. A. Fisher, J. Neyman, and H. Hotelling, who were senior colleagues to statisticians still active today.

The 1980s produced a rising curve of new statistical theory and methods based on the power of electronic computation. Today's data analyst can afford to expend more computation on a single problem than the world's yearly total of statistical computation in the 1920s. How can such computational wealth be spent wisely, in a way that genuinely adds to the classical methodology without merely elaborating it? Answering this question has become a dominant theme of modern statistical theory.

Some promising developments in computer-intensive statistical methodology are described in this article. The examples involve bootstrap methods, nonparametric regression, generalized additive models, and classification and regression trees. The presentation here is mainly descriptive, without much mathematical development. However, we will try to indicate the crucial role that mathematics plays in tying the new statistical methods to their classical antecedents.

The Bootstrap

In almost every statistical data analysis, on the basis of a data set \mathbf{x} we calculate a statistic $t(\mathbf{x})$ for the purpose of estimating some quantity of interest. Box 1 shows the cholesterol reduction scores of

-21.0	3.25	10.75	13.75	32.50	39.50	41.75	56.75	80.0
-------	------	-------	-------	-------	-------	-------	-------	------

Box 1

nine men after taking cholestyramine; the scores are an ordered random sample from the scores of 164 men (1). The data set \mathbf{x} could be these nine scores, and $t(\mathbf{x})$ could be their mean value $\bar{x} = 28.58$, intended as an estimate of the true mean value of the cholesterol reduction scores. (The true mean value is the mean we would obtain if we observed a much larger set of scores.) The following fundamental question arises: how accurate is $t(\mathbf{x})$?

This question has a simple answer if $t(\mathbf{x})$ is the mean \bar{x} of numbers x_1, x_2, \dots, x_n . Then the standard error of \bar{x} , its root-mean-square error, is estimated by a formula made famous in elementary statistics courses

$$se(\bar{x}) = \left\{ \sum_{i=1}^n (x_i - \bar{x})^2 / [n(n-1)] \right\}^{1/2} \quad (1)$$

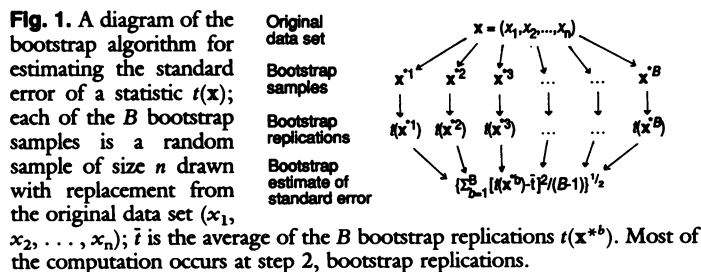
For the nine numbers in Box 1, Eq. 1 gives 10.13. The estimate of the true cholesterol reduction mean would usually be expressed as 28.58 ± 10.13 , or perhaps $28.58 \pm 10.13z$, where z is some constant, such as 1.645 or 1.960, relating to areas under a bell-shaped curve. With $z = 1.645$, the interval has approximately 90% chance of containing the true mean value. In other words, it is an approximate 90% confidence interval.

The bootstrap (2) was introduced primarily as a device for extending Eq. 1 to estimators other than the mean. For example suppose $t(\mathbf{x})$ is the 25% trimmed mean, $\bar{x}\{0.25\}$, defined as the average of the middle 50% of the data. We order the observations x_1, x_2, \dots, x_n , discard the lower and upper 25% of them, and take the mean of the remaining 50%. Interpolation is required for cases where $0.25n$ is not an integer. For the cholesterol data

$$\bar{x}\{0.25\} = \frac{3/4(10.75) + (13.75) + (32.5) + (39.5) + 3/4(41.25)}{3/4 + 1 + 1 + 1 + 3/4} = 27.81 \quad (2)$$

There is no neat algebraic formula such as Eq. 1 for the standard error of a trimmed mean or for almost any estimate other than the mean. That is why the mean is so popular in statistics courses. In lieu of a formula, the bootstrap uses computational power to get a numerical estimate of the standard error. The bootstrap algorithm depends on the notion of a bootstrap sample, which is a sample of

B. Efron is in the Department of Statistics, Stanford University, Stanford, CA 94305. R. Tibshirani is in the Department of Preventive Medicine and Biostatistics and Department of Statistics, University of Toronto, Toronto, Ontario, Canada M5S 1A8.



size n drawn with replacement from the original data set $\mathbf{x} = (x_1, x_2, \dots, x_n)$. The bootstrap sample is denoted $\mathbf{x}^* = (x_1^*, x_2^*, \dots, x_n^*)$. Each x_i^* is one of the original x values, randomly selected (perhaps $x_1^* = x_7, x_2^* = x_5, x_3^* = x_5, x_4^* = x_9, x_5^* = x_7$, and so forth). The name “bootstrap” refers to the use of the original data set to generate new data sets \mathbf{x}^* .

The bootstrap estimate of standard error for $\bar{x}\{0.25\}$ is computed as follows: (i) a large number B of independent bootstrap samples, each of size n , is generated using a random number device, (ii) the 25% trimmed mean is calculated for each bootstrap sample, and (iii) the empirical standard deviation of the B bootstrap trimmed means is the bootstrap estimate of standard error for $\bar{x}\{0.25\}$. A schematic diagram of the bootstrap algorithm, applied to a general statistic $t(\mathbf{x})$, is shown in Fig. 1.

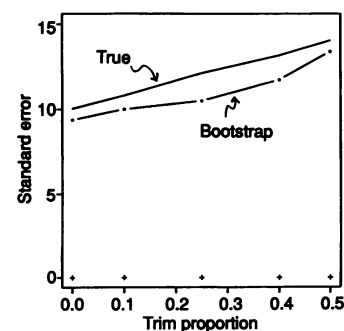
These bootstrap estimates of standard error for the 25% trimmed mean, applied to the cholesterol data, were obtained for different values of B : $B = 25$, bootstrap estimate = 12.44; $B = 50$, bootstrap estimate = 9.71; $B = 100$, bootstrap estimate = 11.50; $B = 200$, bootstrap estimate = 10.70; $B = 400$, bootstrap estimate = 10.48. Ideally, B would go to infinity. However, randomness in the bootstrap standard error that comes from using a finite value of B is usually negligible for B greater than 200; that is, this randomness would be small relative to the randomness caused by variations in the original data set \mathbf{x} . Even values of B as small as 25 often give satisfactory results. This can be important if the statistic $t(\mathbf{x})$ is difficult to compute because the bootstrap algorithm requires about B times as much computation as $t(\mathbf{x})$.

The bootstrap algorithm can be applied to almost any statistical estimation problem: (i) The individual data points x_i need not be single numbers; they can be vectors, matrices, or more general quantities, such as maps or graphs. (ii) The statistic $t(\mathbf{x})$ can be anything at all, as long as we can compute $t(\mathbf{x}^*)$ for every bootstrap data set \mathbf{x}^* . (iii) The data set \mathbf{x} does not have to be a simple random sample from a single distribution. Other data structures, for example, regression models, time series, or stratified samples, can be accommodated by appropriate changes in the definition of a bootstrap sample. (iv) Measures of statistical accuracy other than the standard error, for instance, biases, mean absolute value errors, and confidence intervals, can be calculated at the final stage of the algorithm (3). The example below illustrates some of these points.

There is one statistic $t(\mathbf{x})$ for which one does not need the computer to calculate the bootstrap standard error, namely the mean \bar{x} . In this case, it can be proved that, as B goes to infinity, the bootstrap standard error estimate goes to $\sqrt{(n-1)/n}$ times Eq. 1. The factor $\sqrt{(n-1)/n}$, which equals 0.943 for $n = 9$, could be removed by redefinition of the last step of the bootstrap algorithm, but there is no general advantage to doing so. For the statistic \bar{x} , using the bootstrap algorithm gives about the same result as Eq. 1.

At a deeper level, the logic that makes Eq. 1 a reasonable assessment of standard error for \bar{x} applies equally well to the bootstrap as an assessment of standard error for a general statistic $t(\mathbf{x})$. In both cases, the standard error of the statistic of interest is assessed by the true standard error that would apply if the unknown probability distribu-

Fig. 2. Bootstrap estimates of standard error for five different trimmed means $\bar{x}\{p\}$, $p = 0, 0.10, 0.25, 0.40, 0.5$, applied to the cholesterol data of Box 1, based on $B = 400$ bootstrap samples. Also shown is the true standard error of $\bar{x}\{p\}$, obtained by taking random samples of size 9 from the population of 164 cholesterol reduction scores. In this case, the bootstrap correctly indicates that $\bar{x}\{0\}$, the ordinary mean, gives the smallest standard error.



tion yielding the data exactly equaled the empirical distribution of the data. The efficacy of this simple estimation principle has been verified by a large amount of theoretical work in the statistics literature of the past decade; see (3–5) and references within.

Why use a trimmed mean rather than \bar{x} ? The theory of robust statistics, developed since 1960, shows that if the data \mathbf{x} comes from a long-tailed probability distribution, then the trimmed mean can be substantially more accurate than \bar{x} . That is, it can have substantially smaller standard error (6, 7). In practice, however, one does not know a priori if the true probability distribution is long-tailed. The bootstrap can help answer this question.

The bootstrap estimates of standard error for five different trimmed means $\bar{x}\{p\}$, where p is the proportion of the data trimmed off each end of the sample before the mean is taken are shown in Fig. 2. (So $\bar{x}\{0\}$ is \bar{x} , the usual mean, whereas $\bar{x}\{0.5\}$ is the median.) These were computed with the use of the bootstrap algorithm in Fig. 1 ($B = 400$), except that at step 2, bootstrap replication, five different statistics were evaluated for each bootstrap sample \mathbf{x}^* , namely $\bar{x}\{0\}$, $\bar{x}\{0.10\}$, $\bar{x}\{0.25\}$, $\bar{x}\{0.40\}$, and $\bar{x}\{0.50\}$.

According to the bootstrap standard errors in Fig. 2, the ordinary mean has the smallest standard error among the five trimmed means. This seems to indicate that there is no advantage to trimming for this particular data set.

The nine cholesterol reduction scores in Box 1 were a random sample from a larger data set: 164 scores, corresponding to the 164 men in the Stanford arm of a large clinical trial designed to test the efficiency of the cholesterol-reducing drug cholestyramine (8). With all of this extra data available, the bootstrap standard errors can be checked. The solid line in Fig. 2 indicates the true standard errors for each of the five trimmed means, that is, the standard errors of random samples of size 9 taken from the population of 164 scores.

We see that the true standard errors confirm the bootstrap conclusion that the ordinary mean is the estimator of choice in this case. The main point here is that the bootstrap estimates use only the data in Box 1, whereas the true standard errors require extra data that usually is not available in a real data analysis problem.

Theoretical work on properties of the bootstrap is proceeding at a vigorous pace (4, 5). We have emphasized standard errors here, but the main theoretical thrust has been toward confidence intervals. Getting dependable confidence intervals from bootstrap calculations is challenging, in theory and in practice, but progress on both fronts has been considerable.

Nonparametric Regression

The data for all 164 men in the Stanford arm of the cholestyramine experiment are shown in Fig. 3. The vertical axis plots the cholesterol reduction scores, nine of which appear in Box 1. The horizontal axis plots compliance, the proportion of the intended dose each man actually took (measured by counting of the packets of

unused cholestyramine returned to the clinic). Better compliance tends to be associated with a greater reduction in cholesterol, as might be hoped.

The smooth curve in Fig. 3 is a quadratic regression curve fit to the 164 data points. In other words, it is the quadratic function of compliance that minimizes the sum of the 164 squared distances from the curve to the data points, where distance is measured in the vertical direction. Least-squares regression is a classical estimation method dating back to Gauss and Legendre in the early 1800s (9). The height of the quadratic curve at 60% compliance is 27.72 ± 3.08 . The standard error 3.08 is provided by a formula much like Eq. 1, which is not surprising because the average \bar{x} is the simplest example of a least-squares estimate.

The value 27.72 estimates the true amount of cholesterol reduction at the average compliance (60%), a quantity of particular importance in assessing the true cholesterol-reducing powers of cholestyramine (8). One might worry that a quadratic function of compliance does not accurately model cholesterol reduction as a function of compliance. If not, the estimate 27.72 will be biased, a form of statistical error not included in the formula that gave 3.08.

The irregular curve in Fig. 3 was obtained using loess (pronounced "low ess") (10), a computer-based fitting method that does not attempt to fit a simple model, like a quadratic curve, over the entire compliance range. Instead, loess fits a series of local regression curves for different values of compliance, in each case using only data points near the compliance value of interest.

Loess works in the following way (Fig. 4). First, a window of points (the shaded region) closest to the target point (arrow) is formed; in this case, the window contains the nearest 20% of the data points. Then a smooth weight function (dotted curve) known as the tricube function is constructed so that it is highest at the target point and falls to zero at the edges of the shaded region. Finally, a weighted linear regression (dashed line) is estimated for the points in the shaded region, with the weights determined by the tricube function. This process defines the estimate at the target point. Repeating the process for all possible target points gives the solid curve in Fig. 4. This curve is called a nonparametric regression

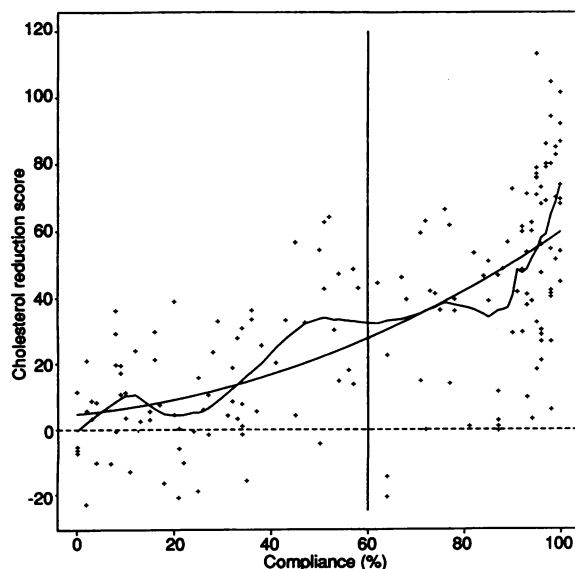
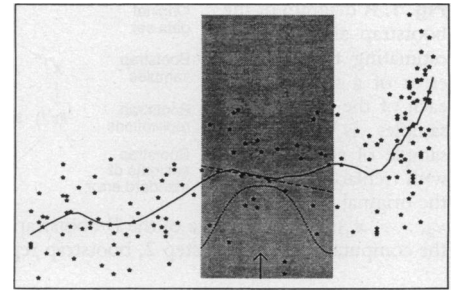


Fig. 3. Cholesterol reduction scores of 164 men in the Stanford arm of experiment LRC-CPPT plotted against compliance, measured as the percentage of intended cholestyramine dose that was actually taken. The average compliance was 60%. The smooth curve is a quadratic regression fit to the 164 points by least squares; the irregular curve is loess, a scatterplot smoother that uses local regressions fit to a moving window of 20% of the points.

Fig. 4. How the loess smoother works. The shaded region indicates the window of points around the target point (arrow). A weighted linear regression (dashed line) is computed, with weights given by the tricube function (dotted curve). Repetition of this process for all target points gives the solid curve.



estimate because it does not assume a particular parametric form (such as quadratic) for the regression.

The height of the loess curve at 60% compliance is 32.38, which indicates substantially greater cholesterol-reducing power than the quadratic estimate 27.72. But how dependable is the loess answer? It is bound to be less biased than the quadratic estimate because it makes fewer assumptions about the form of the dependence between compliance and cholesterol reduction. However, one cannot assess its value as an estimate without some idea of its standard error, and there is nothing like Eq. 1 for loess.

The bootstrap algorithm for standard error can be applied exactly as described in Fig. 1. Now $n = 164$, and each x_i is the pair of numbers (compliance, cholesterol reduction score) for patient i . The function $t(x^*)$ takes any data set x^* consisting of 164 pairs, applies the loess algorithm to it, and reads off the height of the loess function evaluated at 60% compliance. Knowledge of the complicated details of the loess algorithm, is not necessary. All one need do is call the same loess subroutine that gave the estimate 32.38 for the original data set.

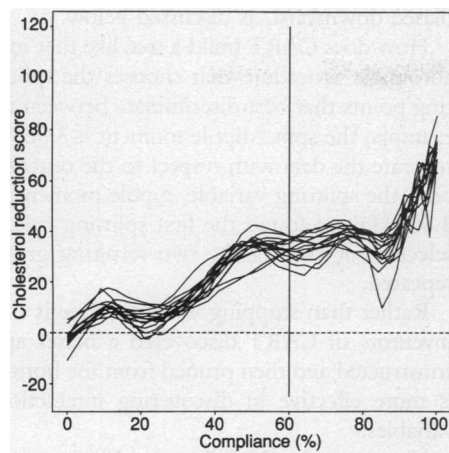
The bootstrap algorithm was run with $B = 50$, and the first 15 of the bootstrap loess curves are shown (Fig. 5). There is considerable variability in the intercepts of these curves at 60% compliance. The bootstrap estimate of standard error for the intercept, based on all 50 bootstrap replications, was 5.71, nearly twice the standard error of the quadratic fit. On balance, the quadratic estimate should probably be preferred in this case. It would have to have an unusually large bias to undo its superiority in standard error.

Generalized Additive Models

Nonparametric regression procedures like loess can be used to model complex data in a flexible manner. This allows the data analyst to make new discoveries about the data. As an example, Williams and colleagues from Toronto's Hospital for Sick Children collected data on the survival of 497 infants after cardiac surgery for heart defects, for the years 1983 to 1988 (11). This was an observational study rather than randomized clinical trials. A warm-blood cardioplegia (WBC) arrest of the heart, thought to improve chances for survival, was introduced in February 1988. The procedure was used on those infants for whom it was thought appropriate and only by those surgeons who liked the procedure. The main question was whether the introduction of WBC improved survival relative to the standard treatment; the importance of risk factors age (in days) and weight (in kilograms) was also of interest. Of the 57 infants who received WBC, 7 died; of the 440 infants who received the standard procedure, 133 died. WBC seemed to be improving the survival rate considerably.

A linear logistic model is the standard way to approach problems of this kind. This model assumes that the log of the odds ratio, probability (death)/probability (survival), is a linear function of the

Fig. 5. The first 15 of the 50 bootstrap loess curves, based on the data for 164 men (Fig. 3). The intercept at 60% compliance has empirical standard deviation 5.71, based on all $B = 50$ bootstrap replications.



age and weight of the infant, plus a term indicating if WBC was used. The results of fitting a linear logistic model to these data suggested that WBC had a strong beneficial effect on survival, with an odds ratio of 3.8 ± 1.8 . Thus the odds of dying were 3.8 times as high with the standard treatment as with WBC. Furthermore, the risk of death decreased with weight, but the age of the infant did not have a significant effect on survival.

Using nonparametric regression procedures, one can learn more from the data. Rather than assuming that the log-odds of survival is a linear function of age and weight, one assumes only that it is a sum of a smooth function of age and a smooth function of weight. This is an example of a generalized additive model (12). The data analyst is not required to specify the form of these smooth functions (such as linear, quadratic, or logarithmic); instead, the form of each of these functions is estimated by a computer-intensive algorithm that makes repeated use of nonparametric regression procedure such as loess.

The curves that resulted from a fit of the generalized additive model are shown in Fig. 6. The shaded regions are approximate confidence bands for the curves. The left curve, for example, represents the log-odds of death as a function of the weight of the infant. The log-odds is highest for the lighter infants (~ -1) and lowest for the heavier infants (~ -3). Hence the odds ratio for light versus heavy infants is the exponential of $[1 - (-3)] \approx 55$. The log-odds does not start to decrease until the infant is at least 3 or 3.5 kg.

The log-odds curve for age is, perhaps, surprising. The operation is least dangerous for infants who are about 200 days old and is more risky for younger or older infants. In a traditional logistic

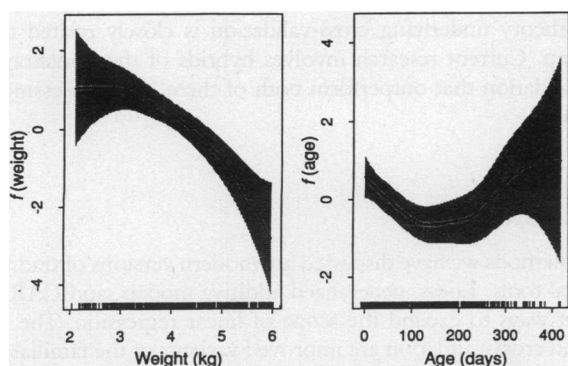


Fig. 6. Function estimates from the heart data (11). The curve on the left represents the log-odds of death as a function of weight; the curve on the right is the log-odds of death as a function of age. The shaded regions are approximate confidence bands.

regression, these curves might be forced to be straight lines, and one would not discover the effects seen in these pictures. The danger of oversimplified regressions becomes more acute in more complicated situations where there are large numbers of explanatory variables.

The generalized additive model also provides an assessment of WBC. The estimated odds ratio for the standard treatment versus WBC was 4.2 ± 1.9 , almost the same as the linear logistic estimate.

Modern statistical tools that are powerful and flexible also tend to be more difficult to analyze mathematically. For example, because of the complexity of the generalized additive model, many approximations were used to obtain the value 1.9 for the standard error reported above. With so much at stake medically, some additional effort to check the accuracy of this value is worthwhile. The bootstrap can be used to accomplish this. A bootstrap sample is created by random drawing of 497 patients with replacement from the original set of 497 patients. A generalized additive model is fit to the bootstrap sample and the estimated odds ratio for WBC is recorded. This entire process is repeated a large number of times, in this case 100. The standard deviation of the 100 odds ratios equaled 2.0, just slightly larger than the approximate value 1.9. The agreement of the bootstrap with the approximate standard error strengthens our belief in both of them.

Generalized additive models can be applied in a wide variety of settings, providing a flexible tool for discovering the underlying structure of scientific processes. Although the algorithm to fit these models required a mainframe computer 10 years ago, now the computations can be carried out on a personal computer. Generalized additive models are just one example of flexible modeling tools that exploit the power of the computer. The development of such tools is an active area of statistical research.

Classification and Regression Trees

In an experiment designed to provide information about the causes of duodenal ulcers (13), one of 56 model alkyl nucleophiles was administered to each of a sample of 745 rats. Each rat was later autopsied to check for the development of duodenal ulcer and the outcome was classified as 1, 2, or 3 in increasing order of severity. There were 535 class 1, 90 class 2, and 120 class 3 outcomes. The objective in the analysis of these data was to ascertain which of 67 characteristics of these compounds was associated with the development of duodenal ulcers.

The CART (Classification and Regression Trees) method (14) is a computer-intensive approach to this problem. When applied to this data, CART produced the classification tree shown in Fig. 7.

At each node of the tree a question is asked; data points for which the answer is "yes" are assigned to the left branch and other data points are assigned to the right branch. The leaves of the tree in Fig. 7 are called terminal nodes. Each observation is assigned to one of the terminal nodes on the basis of the answers to the questions. For example, a rat that received a compound with dipole moment ≤ 3.56 D and melting point $> 98.1^\circ\text{C}$ would go left, then right, and would end up in the terminal node [13, 7, 41]. Triplets of numbers such as [13, 7, 41] below each terminal node number indicate the membership at that node, that is, 13 class 1, 7 class 2, and 41 class 3 observations.

In the CART procedure, each terminal node is assigned a class (1, 2, or 3). The most obvious way to assign classes to the terminal nodes is to use a majority rule and assign the class that is most numerous in the node. With a majority rule, node [13, 7, 41] would be assigned to class 3 and all of the other terminal nodes would be assigned to class 1. In this study, however, the investigators decided that it would be less desirable to misclassify an animal with a severe

ulcer than one with a milder ulcer, and hence they prescribed a higher penalty to errors of the former type. Through the use of the prescribed penalties, a best rule for each terminal node can then be worked out. The assigned class is underlined at each terminal node in Fig. 7; for example, the node at the bottom left ([10, 0, 5]) has the number 5 underlined and hence is a class 3 node.

We can summarize the tree as follows. The top (root) node was split on dipole moment. A high dipole moment indicates the presence of electronegative groups. This split separates the class 1 and 2 compounds: the ratio of class 2 to class 1 in the right split (66/180) is more than five times as large as the ratio in the left split (24/355). However, the class 3 compounds are divided equally, 60 on each side of the split. If, in addition, the sum of squared atomic charges is low, then CART finds that all compounds are class 1. Hence ionization is a major determinant of biologic action in compounds with high dipole moments. Moving further down the right side of the tree, the solubility in octanol then partially separates class 3 from class 2 compounds. High octanol solubility probably reflects the ability to cross membranes and to enter the central nervous system.

On the left side of the root node, compounds with low dipole moment and high melting point were found to be class 3 (severe). Compounds at this terminal node are related to cysteamine (2-aminoethanethiol). Compounds with low melting points and high polarizability, all thiols in this study, were classified as class 2 or 3, with the partition coefficient separating these two classes. Of those chemicals with low polarizability, those of high density were classified as class 1. These chemicals have high molecular weight and volume, and this terminal node contains the highest number of observations. The low-density side of the split is composed of all short chain amines.

In statistical terminology, the data set of 745 observations is called a learning sample. It is easy to work out the misclassification rate for each class when the tree of Fig. 7 is applied to the learning sample. Looking at the terminal nodes that predict classes 2 or 3, the number of errors for class 1 is $13 + 89 + 50 + 10 + 25 + 25 = 212$, so the apparent misclassification rate for class 1 is $212/535 = 39.6\%$. Similarly, the apparent misclassification rates for classes 2 and 3 are 56.7% and 18.3%. "Apparent" is an important qualifier here because misclassification rates in the learning sample can be badly

biased downward, as discussed below.

How does CART build a tree like that in Fig. 7? CART is a fully automatic procedure that chooses the splitting variables and splitting points that best discriminate between the outcome classes. For example, the split "dipole moment ≤ 3.56 " was determined to best separate the data with respect to the outcome classes. CART chose both the splitting variable, dipole moment, and the splitting value, 3.56. Having found the first splitting rule, new splitting rules are selected for each of the two resulting groups, and this process is repeated.

Rather than stopping when the tree is some reasonable size, the inventors of CART discovered a better approach: a large tree is constructed and then pruned from the bottom. This latter approach is more effective in discovering interactions that involve several variables.

How large should the tree be? If we were to build a very large tree with only one observation in each terminal node, then the apparent misclassification rate would be 0%. However, this tree would probably poorly predict the outcomes for a new sample of rats because it is too much geared to the learning sample; in statistical terminology, it is overfit.

The tree of best size would have the lowest misclassification rate for some new data. Thus, if one had a second data set available (a test sample), one could apply the trees of various sizes to it and then choose the one with lowest misclassification rate.

In most situations, one does not have extra data to work with. Data is so precious that all of it is used to estimate the best possible tree. CART uses the method of cross-validation to choose the tree size; this method attempts to mimic the use of a test sample. It works by dividing the data into ten groups of equal size, building a tree on 90% of the data, and then assessing the tree's misclassification rate on the remaining 10% of the data. This is done for each of the ten groups in turn, and the total misclassification rate is computed over the ten runs. The best tree size is determined to be that tree size giving the lowest misclassification rate. This size is used in constructing the final tree from all of the data. The crucial feature of cross-validation is the separation of data for building and assessing the trees: each one-tenth of the data acts as a test sample for the other nine-tenths.

The process of cross-validation not only provides an estimate of the best tree size, it also gives a realistic estimate of the misclassification rate of the final tree. The learning sample misclassification rates computed above are often unrealistically low because the training sample is used both for building and for assessing the tree. For the tree of Fig. 7, the cross-validated misclassification rates were about 10% higher than the learning sampling misclassification rates. It is the cross-validated rates that provide an honest assessment of how effective the tree will be in classifying a new sample of animals.

The theory underlying cross-validation is closely related to the bootstrap. Current research involves hybrids of the bootstrap and cross-validation that outperform both of them in the assessment of error rates.

Conclusion

The methods we have discussed are modern versions of traditional statistical tools. Loess, generalized additive models, and CART are different ways to expand the scope of linear regression. The bootstrap and cross-validation are improved variants of the familiar error estimate in Eq. 1. All of these developments, and a host of others we have not mentioned, differ in one important way from their classical predecessors: they substitute computer algorithms for the traditional mathematical ways of getting a numerical answer. One immediate

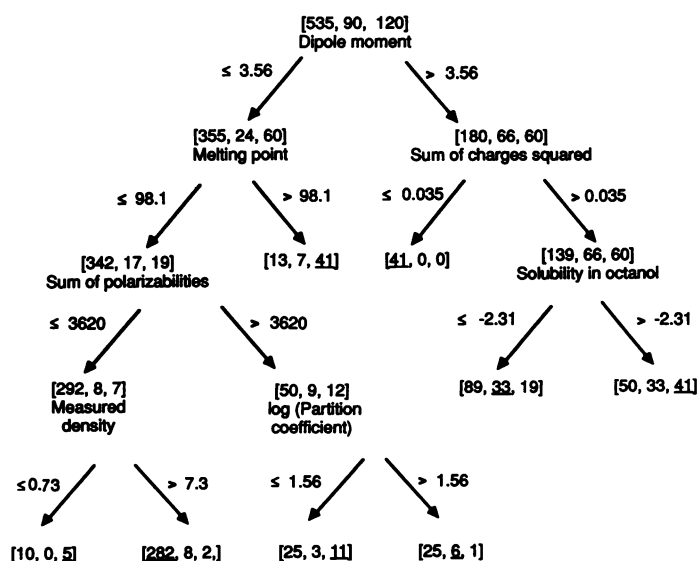


Fig. 7. CART tree. Classification tree from the CART analysis of data on duodenal ulcers (13). At each node of the tree, a question is asked; data points for which the answer is "yes" are assigned to the left branch, and other data points are assigned to the right branch.

reward is freedom from the bell-shaped curve assumptions of the traditional approach. More importantly, the new methods free the scientist to choose statistical methodology appropriate to the problem at hand, rather than choosing on the basis of mathematical tractability.

None of this means that mathematics has disappeared from statistical theory, only that it is disappearing from routine statistical applications. The question of which computer-based method to use, and when to use it, is becoming a central concern of mathematical statistics.

REFERENCES AND NOTES

1. The data in Box 1, from the Stanford arm of the LRC-CPPT experiment, is courtesy of D. Feldman and J. Farquhar, Stanford University; see (7).
2. B. Efron, *Am. Stat.* **40**, 1 (1986).
3. ——— and R. Tibshirani, *Stat. Sci.* **1**, 54 (1986).
4. T. DiCiccio and J. Romano, *J. R. Stat. Soc. B* **50**, 338 (1988).
5. D. V. Hinkley, *ibid.*, p. 321.
6. P. Huber, *Robust Statistics* (Wiley, New York, 1981), p. 5.
7. F. Hampel, E. Ronchetti, P. Rousseeuw, W. Stahel, *Robust Statistics: The Approach Based on Influence Functions* (Wiley, New York, 1986), p. 29.
8. B. Efron and D. Feldman, *J. Am. Stat. Assoc.*, in press.
9. B. Efron, *SLAM Rev.* **30**, 421 (1988).
10. W. S. Cleveland, *J. Am. Stat. Assoc.* **74**, 829 (1979).
11. W. G. Williams *et al.*, *J. Thorac. Cardiovasc. Surg.*, in press.
12. T. Hastie and R. Tibshirani, *Generalized Additive Models* (Chapman and Hall, London, 1990), p. 136.
13. C. Giampaolo, A. Gray, R. Olshen, S. Szabo, *Proc. Natl. Acad. Sci. U.S.A.*, in press.
14. L. Breiman, J. H. Friedman, R. Olshen, C. J. Stone, *Classification and Regression Trees* (Wadsworth, Belmont, CA, 1984).
15. We thank R. Olshen for allowing us to use his CART example.

Enols and Other Reactive Species

YVONNE CHIANG AND A. JERRY KRESGE

Rapid advances in the chemistry of enols and other reactive species have been made possible recently by the development of methods for generating these short-lived substances in solution under conditions where they can be observed directly and their reactions can be monitored accurately. New laboratory techniques are described and a sample of the new chemistry they have made available is provided; special attention is given to ynols and ynamines and the remarkable effects that the carbon-carbon triple bonds of these substances have on their acid-base properties.

THE CHEMISTRY OF ENOLS IS CURRENTLY EXPERIENCING A renaissance (1) primarily because of the development of methods for generating these usually very reactive substances in solution under conditions where their reactions can be studied in detail. Such studies are worthwhile because enols and enolate ions are essential intermediates in many important reactions of carbonyl compounds, and a number of biological reactions also involve enol formation; if we wish to understand these processes, and through understanding to control them, we must understand the chemistry of enols.

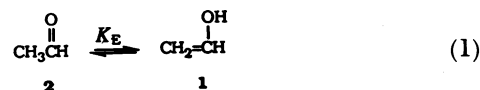
We began work in this area by examining enol isomers of simple aldehydes and ketones. That work, however, soon led to the investigation of other reactive species, such as enols of carboxylic acids and their derivatives, ketenes, carbenes, ynols, and ynamines. The latter are especially fascinating substances: they are believed to exist in interstellar space and are postulated as prebiotic molecules. We have discovered that the carbon-carbon triple bond in ynols and ynamines exerts a remarkable influence on the acid-base properties of their hydroxyl and amino groups; theoretical calculations at the ab initio level have helped us understand the origins of this effect.

This article begins with an account of our work on enols and continues with a description of what we have learned about ynols and

ynamines. Although the discussion is limited largely to research done in our own laboratory, we owe much to stimulation provided by the pioneering work of Guthrie *et al.* (2), Capon *et al.* (3), Dubois, Toulecc, and co-workers (4), and Rappoport and co-workers (5), and we are indebted as well to an early review by Hart (6).

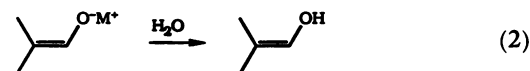
Generation of Enols

Simple enols such as vinyl alcohol, **1**, can be formed readily from their keto isomers, **2**, Eq. 1.



The reaction, however, is reversible, and the position of equilibrium generally lies strongly on the keto side; the amount of enol present at equilibrium is consequently seldom sufficient to permit direct observation, even by the most sensitive spectroscopic methods. Investigation of enol chemistry therefore requires generation of the enol in greater than the equilibrium amount in the medium of interest. We have developed a number of ways of accomplishing this in aqueous solution.

We first made enols by hydrolysis of their alkali metal salts, Eq. 2,



using solutions of these salts in aprotic solvents prepared by standard synthetic methodology (7). Addition of a small quantity of such a solution to a large amount of water resulted in a very fast oxygen-to-oxygen proton transfer and produced the enol in an essentially wholly aqueous medium. Conversion of the enol to its keto isomer then proceeded at a slower rate, which we could monitor accurately by following the marked change in the ultraviolet spectrum that accompanies the ketonization reaction.

This method of generating enols requires mixing two solutions and consequently cannot be applied to substances with lifetimes shorter than the mixing time. This limitation unfortunately excludes

The authors are in the Department of Chemistry, University of Toronto, Toronto, Ontario, Canada M5S 1A1.



Statistical Data Analysis in the Computer Age
BRADLEY EFRON and ROBERT TIBSHIRANI (July 12, 1991)
Science **253** (5018), 390-395. [doi: 10.1126/science.253.5018.390]

Editor's Summary

This copy is for your personal, non-commercial use only.

- Article Tools** Visit the online version of this article to access the personalization and article tools:
<http://science.sciencemag.org/content/253/5018/390>
- Permissions** Obtain information about reproducing this article:
<http://www.sciencemag.org/about/permissions.dtl>

Science (print ISSN 0036-8075; online ISSN 1095-9203) is published weekly, except the last week in December, by the American Association for the Advancement of Science, 1200 New York Avenue NW, Washington, DC 20005. Copyright 2016 by the American Association for the Advancement of Science; all rights reserved. The title *Science* is a registered trademark of AAAS.