where $j = 0, 1, 2, \ldots, k = 1, 2, 3, \ldots$, and when either integral exists. The result has been in the literature in various forms for some time, but it is certainly not well known.

## REFERENCES

ABRAMOWITZ, M., and STEGUN, I. (eds.) (1965), *Handbook of Mathematical Tables*, New York: Dover.

CHAO, M.T., and STRAWDERMAN, W.E. (1972), "Negative Moments of Positive Random Variables," *Journal of the American Statistical Association*, 67, 429–431.

FELLER, W. (1968), *Introduction to Probability Theory, Vol. 1* (3rd ed.), New York: John Wiley.

HALPERIN, M., and GURIAN, J. (1971), "A Note on Estimation in Straight-Line Regression When Both Variables Are Subject to Error," *Journal of the American Statistical Association*, 66, 587–589.

HOGG, R.V., and CRAIG, A.T. (1970), *Introduction to Mathematical Statistics* (3rd ed.), New York: Macmillan.

KABE, D.G. (1976), "Inverse Moments of Discrete Distributions," *Canadian Journal of Statistics*, 4, 133–141.

LAUE, G. (1980), "Remarks on the Relation Between Fractional Moments and Fractional Derivatives of Characteristic Functions," *Journal of Applied Probability*, 17, 456–466.

OBERHETTINGER, F. (1974), *Tables of Mellin Transforms*, Berlin: Springer.

OLDHAM, K.B., and SPANIER, J. (1974), *The Fractional Calculus*, New York: Academic Press.

SCHUH, H.-J. (1981), "Sums of iid Random Variables and an Application to the Explosion Criterion for Markov Branching Processes," *Journal of Applied Probability*, 18, in press.

WILLIAMS, J.D. (1941), "Moments of the Ratio of the Mean Square Successive Difference to the Mean Square Difference in Samples From a Normal Universe," *Annals of Mathematical Statistics*, 12, 239–241.

# Eye Fitting Straight Lines

FREDERICK MOSTELLER, ANDREW F. SIEGEL, EDWARD TRAPIDO, AND CLEO YOUTZ*

Because little is known about properties of lines fitted by eye, we designed and carried out an empirical investigation. Inexperienced graduate and postdoctoral students instructed to locate a line for estimating $y$ from $x$ for four sets of points tended to choose slopes near that of the first principal component (major axis) of the data, and their lines passed close to the centroids. Students had a slight tendency to choose consistently either steeper or shallower slopes for all sets of data.

KEY WORDS: Least squares; Regression; Subjective fitting; Principal components.

## 1. INTRODUCTION

The properties of least squares and other computed lines are well understood, but surprisingly little is known about the commonly used method of fitting by eye. This method involves maneuvering a string, black thread, or ruler until the fit seems satisfactory, and then drawing the line. We report one systematic investigation of eye fitting lines.

Students fitted lines by eye to four sets of points given in an experimental design to help us discover the properties of their fitted lines and whether order of fitting or practice made a difference. Other populations of subjects may produce different results. These sets of data were not unusual in curvature or in having outlying points or patterns. Thus additional populations of data sets could profitably be investigated.

The principal quantitative reference on fitting straight lines by eye is Finney (1951). He found that a mathematical iteration starting with slopes provided by scientists, inexperienced with probit analysis, gave satisfactory approximations to the relative potency in a bioassay.
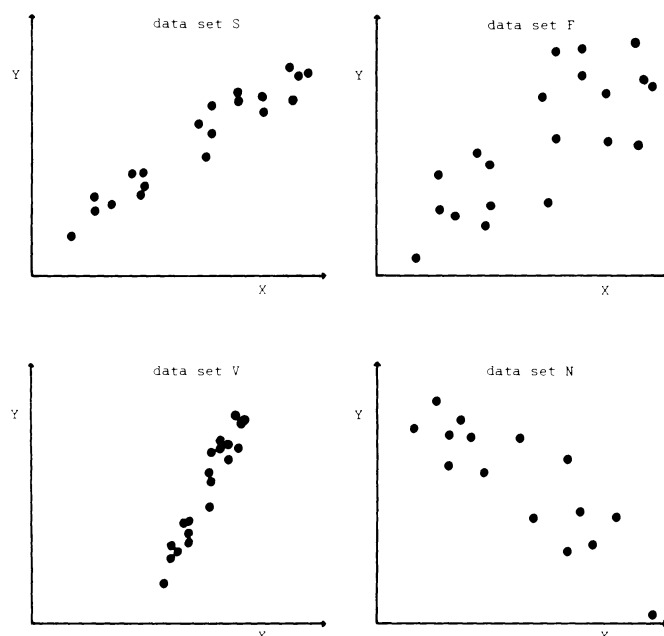
Figure 1. The Data Sets of S, F, V, and N

## 2. METHOD

We conducted this experiment in a class of graduate and postdoctoral students in introductory biostatistics. Most students had not studied statistics before and had not yet been shown formal methods for fitting lines. In a previous class session we illustrated the idea of using a regression line fitted to a set of points to estimate the vertical value, $y$, from the horizontal value, $x$.

Each student was given the same set of four scatter diagrams and an $8\frac{1}{2} \times 11$ inch transparency with a straight line etched completely across the middle. Students moved the transparency over the scatter diagram until satisfied with the fit of the etched line, and then marked an $\times$ on the scatter diagram at each end of the line. This transparency method is preferable to the black-thread method, which requires three hands.

The four scatter diagrams were labeled S for standard, F for fat, V for vertical, and N for negative; these are shown in Figure 1. Data sets S, F, and V are linear transformations of each other, so that F has more vertical error than S and V has a steeper slope

than S. Data sets S, F, and V come from a table of random numbers in Beyer (1971), whereas data set N is a linear transformation of the fiber strength data in Dunn and Clark (1974, p. 224).

To assess the effect of the order of presentation, we used a Latin square design with packets stapled in four different orders: SNFV, NSVF, FVSN, and VFNS. We distributed them systematically in that sequence so that students sitting side by side had different kinds of packets. We laid out on desks before class 175 packets and collected 153 at the end of the hour.

## 3. RESULTS

Table 1 summarizes the averages, variabilities, and actual (least squares) values for the slope and intercept of each data set. We have reported medians and interquartile ranges to reduce the effect of the few outlying values. The $y$ intercept at $\bar{x}$ measures the height of a line as well as does the $y$ intercept at zero, and is less correlated with the slope. To get Table 1, we pooled results from the four orders of presentation because we found no trend in the differences due to order.

Comparing the students average slope to the actual slope, we see that the slope of the least squares regression of $y$ on $x$ is close to the average in each data set except F. One possible explanation might be that students tended to fit the slope of the first principal component or major axis (the line that minimizes the sum of squares of perpendicular rather than vertical distances). The principal component slope is closer to the median slope in every case except N, and is notably closer for F.

Because the $y$ intercept at $\bar{x}$ is the same for the regression and major axis lines, the conclusion here is simply that the students placed their lines near the centroid of the cloud of points in each case.

By computing the correlation matrix for the students' slopes for the four data sets, we see in Table 2 that students who gave steep slopes for one data set also tended to give steep slopes on the others. This effect seems slight but is definite. The negative values arise because data set N has negative slope.

The individual-to-individual variability in slope and

Table 1. Averages, Variabilities, and Actual Values for Slopes and Intercepts

|  | S | F | V | N |
|---|---|---|---|---|
| Slope | | | | |
| median (interquartile range) | .70 (.04) | .84 (.14) | 2 07 ( 14) | −.73 (.20) |
| actual least squares | | | | |
| regression | .66 | .66 | 1.98 | −.70 |
| principal component | .68 | .82 | 2.11 | −.79 |
| y intercept at x̄ | | | | |
| median (interquartile range) | 3.88 ( 10) | 3.86 (.17) | 3.95 (.18) | 4 04 ( 24) |
| actual least squares | 3.88 | 3.90 | 3.89 | 4.11 |

*Table 2. Correlations Between Slope Estimates*

|   | F | V | N |
|---|---|---|---|
| S | .18 | .14 | −.14 |
| F |   | .28 | −.08 |
| V |   |   | −.05 |

in intercept was near the standard error provided by least squares for the four data sets. When comparable measures of variability were used, that for slopes was .9 times and that for intercepts was .7 times the least squares standard error. Admittedly, no theory

encourages us to believe in such relations, and further empirical work might be instructive.

## REFERENCES

BEYER, WILLIAM H. (ed.) (1971), *Basic Statistical Tables*, Cleveland, Ohio: Chemical Rubber Co.
DUNN, OLIVE JEAN, and CLARK, VIRGINIA A. (1974), *Applied Statistics: Analysis of Variance and Regression*, New York: John Wiley.
FINNEY, D.J. (1951), "Subjective Judgment in Statistical Analysis: An Experimental Study," *Journal of the Royal Statistical Society*, Ser. B, 13, 284–297.

# The Sample Coefficient of Determination in Simple Linear Regression

G.B. RANNEY AND C.C. THIGPEN*

The behavior of the sample coefficient of determination is examined for some arrangements of independent variable values in a simple linear regression with normally distributed error terms. Numerical values of means and standard deviations are presented that provide some insight into the influence of range and arrangement of independent variable values and sample size on the size of the sample coefficient of determination. Some asymptotic results are given.

KEY WORDS: Sample coefficient of determination; Simple linear regression.

When we observe pairs of values, $(X, Y)$, randomly sampled from a bivariate population, $r^2$, the sample coefficient of determination is an estimate of $\rho^2$, the population coefficient of determination. Thus $r^2$ is an estimator of strength of linear relationship between the two variables. The same type of inference is not available when the independent variable values are fixed.

In a descriptive sense, $r^2$ measures the proportion of variation in $Y$ attributable to variation in $X$ in a sample. Some authors (e.g., Neter and Wasserman 1974) warn that when the $X$ values are fixed, $r^2$ does not estimate any single model parameter and the value of $r^2$ is influenced by the spacing of the $X$ values.

* G.B. Ranney is Assistant Professor and C.C. Thigpen is Professor and Head, both with the Department of Statistics, The University of Tennessee, Knoxville, TN 37916.

In spite of these warnings, the student may assume that $r^2$ should be used here, as in the bivariate case, to estimate strength of linear relationship in a population sense and usefulness of $X$ in predicting $Y$. In this context we wish to examine the behavior of $r^2$ with regard to the choice of $X$ values and sample size. We seek to examine three sources of possible effects on $r^2$: the number of repeated observations of $Y$ for a given set of $X$ values, the range of the $X$ values, and the arrangement of $X$ values within a given range.

Consider the simple linear model $Y_i = \alpha + \beta X_i + \epsilon_i$, where the $X_i$ are fixed and the components $\epsilon_i$ are independent normal random variables, each with mean zero and variance $\sigma^2$. Then $r^2$ from a sample of size $n$ has a noncentral beta distribution with one and $n - 2$ degrees of freedom (df) and noncentrality parameter $(\beta/\sigma)^2 \sum (X_i - \bar{X})^2/2$. We will examine the behavior of the mean of $r^2$ under the preceding assumptions. In general, $E(r^2)$ cannot be expressed in a simple closed form in terms of the model parameters. Therefore, we will numerically examine some special cases, using series representations of the mean and variance of a noncentral beta random variable presented by Wishart (1932).

First, let the $X$ values be the first 10 positive integers. To obtain a sample of size $n = 10k$, $k$ observations are made on $Y$ at each of the 10 $X$ values ($k = 1, 2, \ldots$). We wish to see the effect on $E(r^2)$ of increasing the number of replicate measurements of $Y$ on a given set of $X$ values. In Table 1, values of the mean and standard deviation of $r^2$ are shown for $\beta/\sigma$, ranging from zero to five and $k = 1, 2,$ and 10. For each value of $\beta/\sigma$ in Table 1, $E(r^2)$ is largest when $k$ is one and is nonincreasing as $k$ in-