# The bimean: A measure of central tendency that accommodates outliers

PETER R. KILLEEN
*Department of Psychology, Arizona State University, Tempe, Arizona*

Observations that depart considerably from the center of a distribution demand special consideration. They may be retained, trimmed, or weighted less than other data. This article provides a BASIC program that implements Mosteller and Tukey's (1977) technique for weighting observations less as they depart from the middle of a distribution. Influence curves for this "bisquare-weighted mean," or "bimean," are displayed and compared with more traditional measures of central tendency.

When data are averaged, they are subjected to a sequence of mathematical operations to derive an index that is representative of their central tendency. One of the simplest operations is the calculation of the arithmetic mean, which yields a number with several desirable statistical properties. But the arithmetic mean is not always the best measure of central tendency. If the data are exponentially distributed, the geometric mean might be more appropriate, as it would be if probable error was proportional to the magnitude of the datum. In general, estimators are chosen based on the nature of the distribution of the data, and on the ease of computation and description of the estimator.

The use of most traditional averaging procedures presumes that all of the data belong in the sample. But there is always some possibility that an error has crept into the sample, due, for example, to clerical errors, aberrant subject responses, or equipment malfunctions. Ideally, these data should be excluded, but how are they discriminated from the good data? If the intrusive data are very different from the rest, they can be discarded upon inspection. Few experimentalists would retain a data point that is 5 or 6 SD from the mean, but what about one that is 3 SD from the mean? Confidence increases as the data points move toward the center of the distribution, but there is no discrete point at which confidence goes from 0% to 100%.

Some techniques that are employed to minimize the impact of outliers, such as trimming (e.g., deleting the upper and lower 5% of the data) or Winsorizing (e.g., reassigning the values of the most extreme data as those of their nearest neighbors), are based on step functions in confidence about the data. An extreme step function is embodied in the median, which gives 100% confidence to the middlemost datum (or 50% to each of a pair of middlemost data) and 0% confidence to the rest. Whereas the median is maximally robust against outliers, it is an inefficient measure of the central tendency of the data. Barnett and Lewis (1978) and Huber (1972) reviewed treatment of outliers by experimentalists and systematically evaluated the techniques designed to cope with them.

## THE BIMEAN

Hoaglin, Mosteller, and Tukey (1983) and Mosteller and Tukey (1977) have provided an alternate measure of central tendency, the "bisquare-weighted mean." This measure has also been referred to as the "biweight," but as it is a mean and not a weight, the contraction "bimean" seems more appropriate and will be used here. This procedure provides a continuous weighting of data as a function of their proximity to the adjusted center of the distribution. The weight is $w_i = (1-u_i^2)^2$ where $u_i = (x_i - BWM)/(cS)$ and $0 \leq u_i \leq 1$. BWM is the bimean, a weighted arithmetic mean of the data, each datum being weighted by $w_i$. Because the weights are defined in terms of the $u_i$ and depend on the value of the bimean, calculation of the bimean must be iterative, beginning with an estimated value for BWM, such as the median. S is a robust measure of spread, such as the semi-interquartile range, and c is a parameter that specifies how quickly the weights will roll off. For large values of c (e.g., $\geq 20$), all weights are close to 1.0, and the bimean approximates the simple arithmetic mean. Smaller values of c give less weight to the tails of the distribution. The iterative process converges quickly: New estimates are within .01% of the previous estimates after four or five iterations.

Figure 1 shows the "influence curves" generated by introducing an alien datum into the set 36, 43, 48, 52, 57, 64. These numbers are approximately normally distributed, with mean = 50 and SD = 10. A seventh number (x) is added, first with the value of 50, and then various measures of central tendency are calculated. Then, the value of x is incremented, and the process is repeated. The influence is symmetric around the mean, so that only
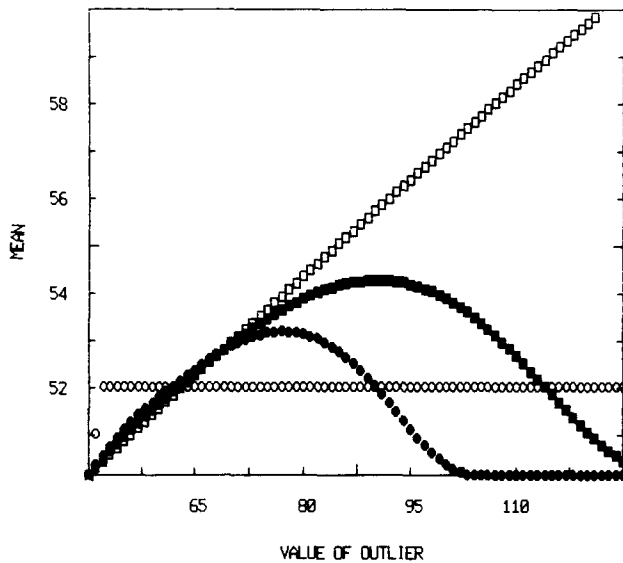
Figure 1. Influence curves for the median (open circles), mean (open squares), bimean with c=6 (filled circles) and bimean with c=9 (filled squares). These curves show the mean of the distribution of data (36, 43, 48, 52, 57, 64, x) as the value of the outlier, x, is incremented from 50 to 125.

values ≥ 50 are examined. The results of these calculations are shown in Figure 1, which displays the influence curves for the median, arithmetic mean, and bimeans with c=6 and c=9.

The arithmetic mean (open squares) is of course a linear function of the value of the outlier, whereas the median (open circles) increases with the value of the outlier until the outlier is no longer the middlemost number and, thereafter, remains fixed and independent of the value of the outlier. Both bimeans are close to the arithmetic mean until the value of the outlier exceeds two standard deviations (70, in terms of the original six data), at which point they begin to diverge. For c=6 (filled circles) the weight on the outlier decreases to 50% at 3 SD, and to 1% at 5 SD. For c=9 (filled squares) the weight on the outlier decreases to 50% between 4 and 5 SD, and to 1% at 7 SD.

Table 1 shows the bimean calculated for a range of values of c when the outlier is 3 and 4 SD from the mean. Note that for small and moderate values of c, the outlier

**Table 1**
**The Effect of c on the Bimean When the Outlier is 3 (x = 80) or 4 (x = 90) SD Removed From the Mean**

| | BWM | |
|---|---|---|
| c | x = 80 | x = 90 |
| 4 | 50.9 | 50.0 |
| 5 | 52.3 | 50.4 |
| 6 | 53.0 | 51.9 |
| 7 | 53.4 | 53.0 |
| 8 | 53.7 | 53.7 |
| 9 | 53.8 | 54.2 |
| 10 | 53.9 | 54.5 |
| 15 | 54.1 | 55.2 |
| ∞ | 54.3 | 55.7 |

has a greater impact on the bimean when it is smaller. This increased impact is because the larger weight for the outlier that is closer to the center of the distribution more than offsets its smaller value. For values of c greater than 15, the bimean is very close to the arithmetic mean. The bimean assumes the value of the arithmetic mean as c approaches infinity.

Because the bimean is more efficient than the median, its use effectively increases the number of data in terms of that statistic (by 40%, when the data are drawn from a Gaussian distribution). Because it is less efficient than the mean, use of the bimean effectively diminishes (by 4%) the number of data in terms of that statistic. It is clear from Figure 1 that the bimean is very robust against outliers, and 4% is small insurance to pay for that protection. Routine use of the bimean saves the investigator from the painful and ad hoc decision of when to reject data as outliers. It is acceptable to editors (see, e.g., Roberts & Holder, 1985).

**Table 2**
**Program Listing**

```
10 '              "BIMEAN":  Calculates bisquare-weighted means
20 DIM X(64),U(64),W(64)
30 DEFINT I-L,N
40 C = 7'       Set C between 5 and 10; smaller values increase weighting
45 CRIT=.0001'  Criterion for ending iteration (rel change in BWM)
50 S=0:S2=0:BWMOLD=0:CHNG=1:F53$ = "#####.###"
100 '           Input
105 PRINT
110 INPUT "# DATA";N
120 FOR I=1 TO N
130   PRINT I;:INPUT X(I)
140   IF X(I)=-9999 THEN I=I-2: GOTO 170'     Redo previous input
150   S=S+X(I)
160   S2=S2+X(I)*X(I)
170 NEXT I
180 MEAN = S/N
190 SD=SQR((S2-S*S/N)/(N-1))'   Sample std dev
200 CLS
210 '           Rank data
220 FOR I=1 TO N-1
230   FOR J=N TO I+1 STEP -1
240     IF X(J)>X(J-1) THEN SWAP X(J),X(J-1)
250   NEXT J:NEXT I
260 NM=INT(N/2)'               Calculate median
270 IF 2*NM<>N THEN MED = X(NM+1) ELSE MED = (X(NM)+X(NM+1))/2
280 DPTH = INT((N+1)/4)
290 W1=((DPTH+.5)/N-.25)*N:W2=1-W1:'   Interpolate for quartiles
300 Q1=X(DPTH)*W1+X(DPTH+1)*W2
310 Q3=X(N-DPTH+1)*W1+X(N-DPTH)*W2
320 SIQR = (Q1-Q3)/2'          Semi-interquartile range
330 SPRD = C*SIQR
340 BWM = (MED+MEAN)/2'        Initial estimate of BWM
350 PRINT "C =";C;"  BWM =";
370 WHILE ABS(CHNG) > CRIT'    Loop until rel change falls to criterion
380   BWS=0:WS=0
390   PRINT USING F53$;BWM;
400   FOR I=1 TO N
410     U(I)=(X(I)-BWM)/SPRD
420     IF ABS(U(I))>1 THEN W(I)=0:GOTO 440'   Keep weights positive
430     W(I)=(1-U(I)*U(I))^2'   Bi-square weights
440     BWS=BWS+W(I)*X(I)'      Weighted sum
450     WS=WS+W(I)'             Sum of weights
460   NEXT I
470   BWM=BWS/WS
475   CHNG=(BWMOLD-BWM)/BWM'    Relative change in BWM
478   BWMOLD=BWM
480 WEND
485 '          Calculate Median Absolute Deviation
490 FOR I=1 TO N-1'     Rank (normalized) deviations
500   FOR J=N TO I+1 STEP -1
510     IF ABS(U(J))>ABS(U(J-1)) THEN SWAP U(J),U(J-1)
520   NEXT J:NEXT I
530 IF 2*NM<>N THEN MAD=SPRD*ABS(U(NM+1))
                 ELSE MAD=SPRD*(ABS(U(NM))+ABS(U(NM+1)))/2
600 '          Output
620 PRINT:PRINT:PRINT "DATA","  WEIGHTS"
630 FOR I=N TO 1 STEP -1:PRINT X(I),W(I):NEXT I
640 PRINT
650 PRINT "BIMEAN =";BWM;TAB(30);"SEMI-INTRQRTL RANGE =";SIQR
660 PRINT "MEDIAN =";MED;TAB(30);"MEDIAN ABSOLUTE DEV =";MAD
670 PRINT "MEAN   =";MEAN;TAB(30);"STANDARD DEVIATION  =";SD
680 RUN
700 END
```

A program that calculates bimeans, ranks the data, and reports other measures of location and spread is listed in Table 2. Sample output is printed in Table 3.

I have used bimeans for several years to average behavioral data both within and across subjects, and to estimate representative parameters of models whose effect on predictions is often strongly nonlinear. This estimator cannot, of course, reduce the importance of good experimental control and sampling techniques. As users grow more comfortable with this estimator, they tend to use smaller values for c; values around 7 are most commonly employed.

## PROGRAM LANGUAGE AND REQUIREMENTS

The program is written in MICROSOFT BASIC and is run on a Tandy 2000. With the exceptions of the "swap" and "while/wend" commands, which are easily modified, the program should be readily transportable to other computers and other operating systems.

### Table 3
### Sample Output

| c = 7 | BWM = | 53.143 | 53.376 | 53.418 | 53.425 |
|---|---|---|---|---|---|
| Data | Weights | | | | |
| 36 | .8528539 | | | | |
| 43 | .9459867 | | | | |
| 48 | .9852256 | | | | |
| 52 | .9989772 | | | | |
| 57 | .9935695 | | | | |
| 64 | .9444401 | | | | |
| 80 | .6757846 | | | | |

| Bimean | = 53.42612 | Semi-Intrqrtl Range | = 9 |
|---|---|---|---|
| Median | = 52 | Median Absolute Dev | = 10.42484 |
| Mean | = 54.28572 | Standard Deviation | = 14.54549 |

# Data?

### REFERENCES

BARNETT, V., & LEWIS, T. (1978). Outliers in statistical data. New York: Wiley.

HOAGLIN, D. C., MOSTELLER, F., & TUKEY, J. W. (1983). Understanding robust and exploratory data analysis. New York: Wiley.

HUBER, P. J. (1972). Robust statistics: A review. Annals of Mathematical Statistics, 43, 1041-1067.

MOSTELLER, F., & TUKEY, J. W. (1977). Data analysis and regression: A second course in statistics. Reading, MA: Addison-Wesley.

ROBERTS, S., & HOLDER, M. D. (1985). Effect of classical conditioning on an internal clock. Journal of Experimental Psychology: Animal Behavior Processes, 11, 194-214.