

COM 531  
 Neuendorf  
 Transforming Data  
 Errata 3/1/04—replaces old page 1

General guidelines:

- \*Always check the effect of the transform before proceeding with further analyses.
- \*Transformations should be applied to the IV except in the case of heteroscedasticity. This is particularly true for multivariate situations, such as multiple regression (Fox, 1997).
- \*As Fox (1997) says, an effective transformation can be selected analytically or by trial and error. (“It’s a creative process” some have said.)
- \*As noted by Hair et al. (1998), for a noticeable effect from transformations, the ratio of a variable’s mean to its standard deviation should be less than 4.0. When the transformation may be performed on either of two variables, select the variable with the smallest such ratio.
- \*Always remember that a transformed variable is now a *new* variable, and needs to be interpreted as such. (See Figure D below.)

Types of transformation:

- I. Transforming data to correct for deviations from normality in a univariate distribution. As noted by Fox, descending the ladder of powers (e.g., to the square root of X or log X) tends to correct a positive skew; ascending the ladder of powers (e.g., to  $X^2$  or  $X^3$ ) tends to correct a negative skew. (NOTE: Hair et al. are *not* correct in their advice to use a square root with negatively skewed data, p. 77) (See Figure D below.)

<u>Problem</u>	<u>Remedy</u>
Positive skew	take a root (e.g., square root) or logarithm (e.g., $\log_{10}$ , ln)
Negative skew	use an exponent (e.g., square, cube)
Positive kurtosis (lepto, pointy)	??—in search of a good transform
Negative kurtosis (platy, flat)	take the inverse (1/X or 1/Y)

- II. Transforming data to correct for heteroscedasticity (DV exhibits unequal variance across the range of an IV). As Hair et al. note, heteroscedasticity is also due to the distribution(s) of the variable(s); hence, the first step is to check for non-normality of each of the variables, and transform accordingly.

- III. Transforming data to correct for heteroscedasticity of residuals in multiple regression (residuals exhibit unequal variance of residuals (errors of prediction) across the range of the predicted DV). Hair et al. offer the following:

<u>Problem</u>	<u>Remedy</u>
Cone opens to the right	take the inverse ( $1/X$ )
Cone opens to the left	take a root (e.g., square root of $X$ )

- IV. Transforming data to linearize a nonlinear relationship. (See Figures A through D below.)

\*The introduction of nonlinear components in a multiple regression model—*polynomials*. For example,  $X^2$  and  $X^3$  are both polynomial versions of  $X$  that may be included in a regression model.

\*Monotonic vs. non-monotonic relationships

#### References:

Fox, J. (1997). *Applied regression analysis, linear models, and related methods*. Thousand Oaks, CA: Sage Publications.

Hair, J. F. Jr., Anderson, R. E., Tatham, R. L., & Black. W. C. (1998). *Multivariate data analysis* (5<sup>th</sup> ed.). Upper Saddle River, NJ: Prentice Hall.

3/1/04

Figure A: An example of linearizing a monotonic, nonlinear relationship

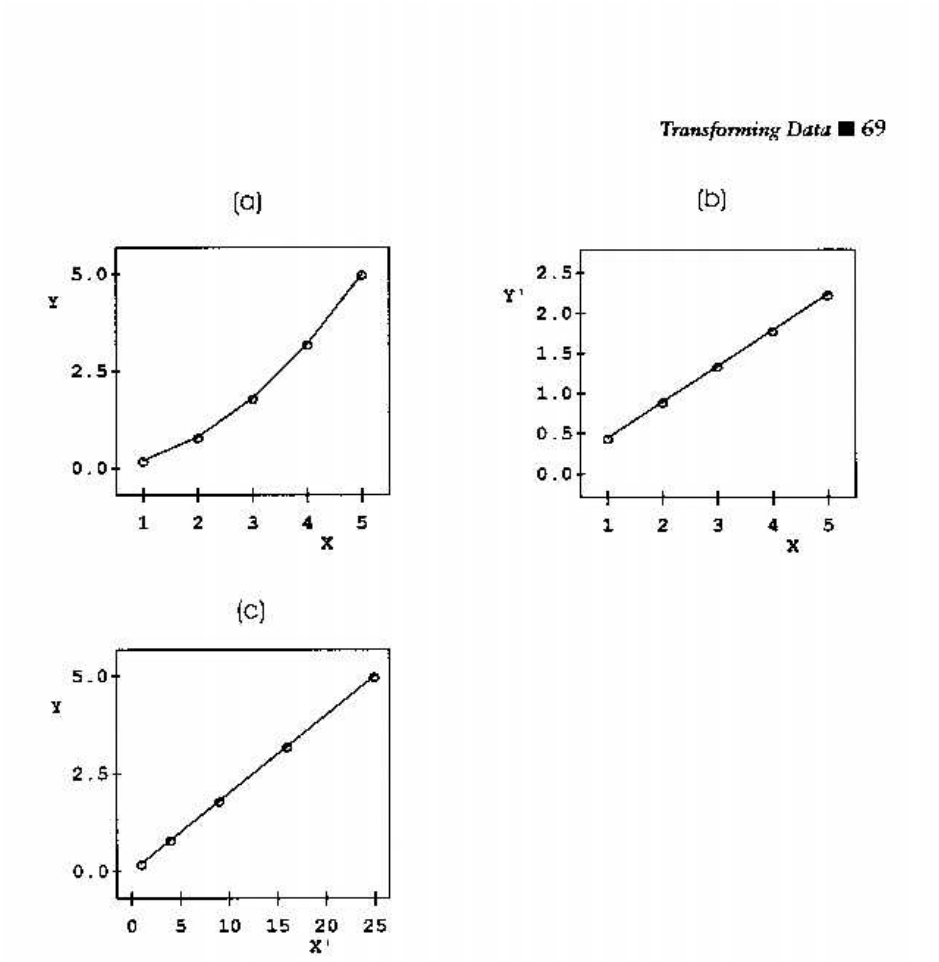


Figure 4.4. How a power transformation of Y or X can make a simple monotone nonlinear relationship linear. Panel (a) shows the relationship  $Y = \frac{1}{3}X^2$ . In panel (b), Y is replaced by the transformed value  $Y' = Y^{1/2}$ . In panel (c), X is replaced by the transformed value  $X' = X^2$ .

Figure B: Examples of monotonic and non-monotonic nonlinear relationships

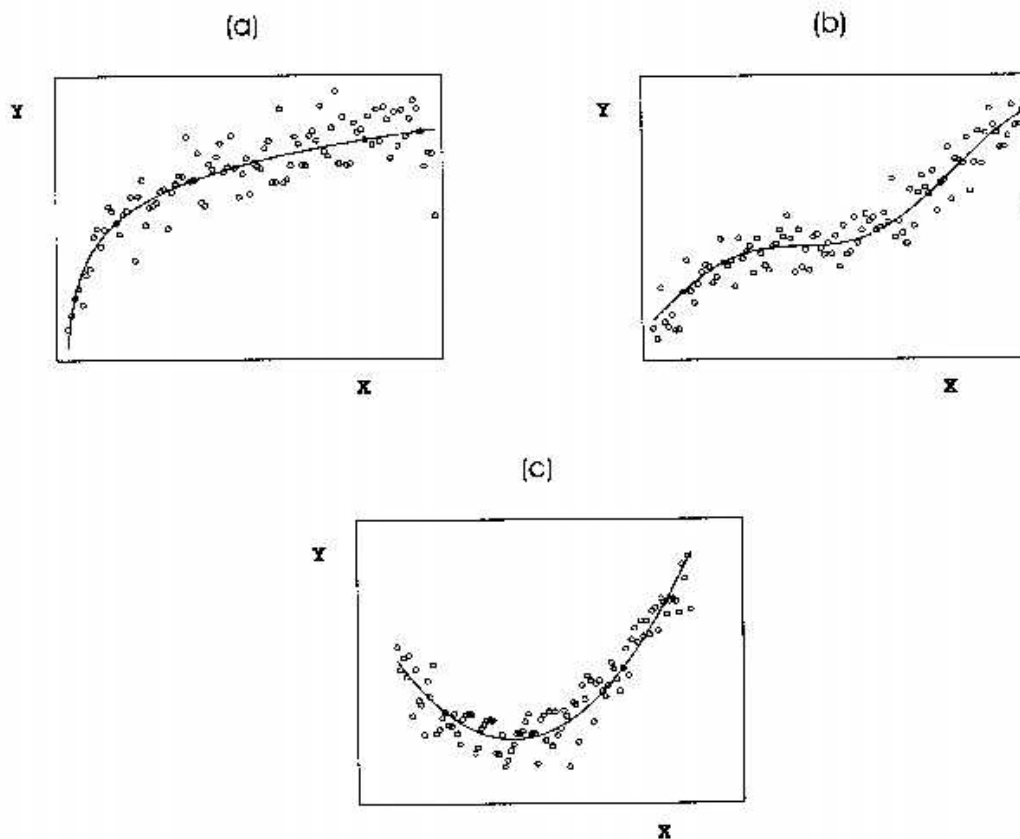


Figure 4.5. (a) A simple monotone relationship between  $Y$  and  $X$ ; (b) a monotone relationship that is not simple; (c) a relationship that is simple but not monotone. A power transformation of  $Y$  or  $X$  can straighten (a), but not (b) or (c).

Figure C: The “bulging rule” for transforming nonlinear, monotonic relationships

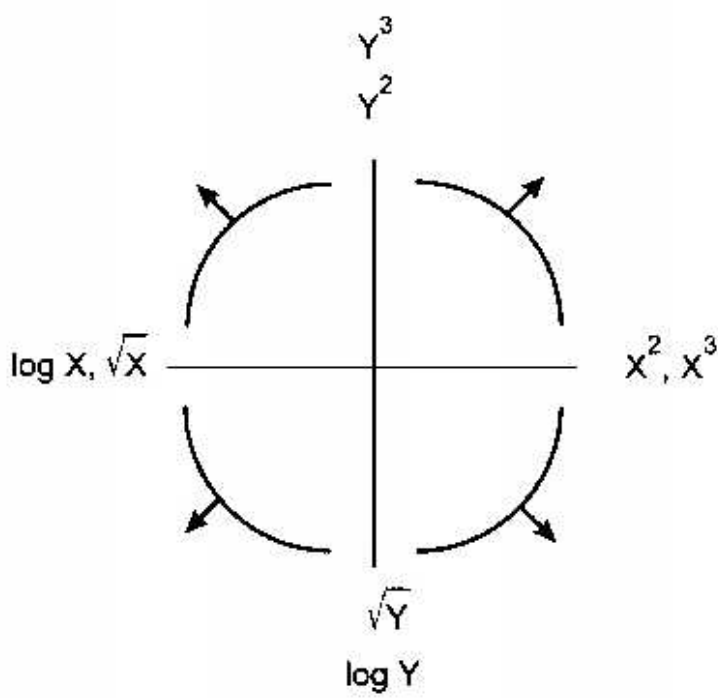
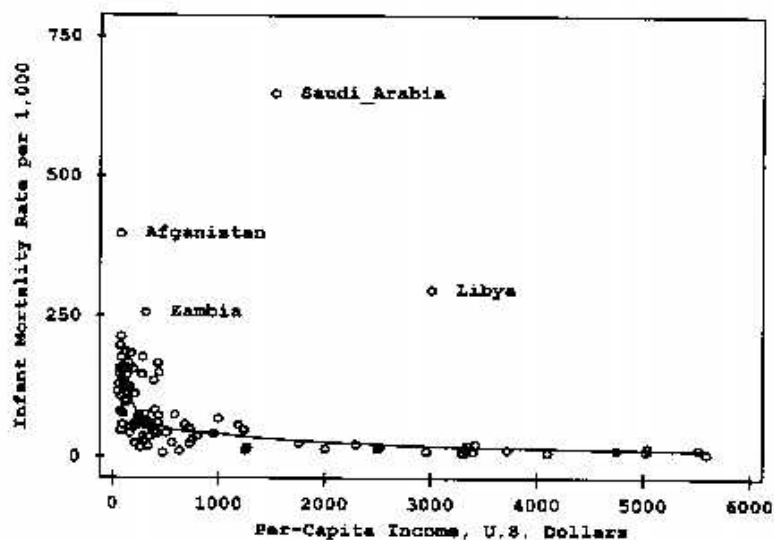


Figure 4.6. Tukey and Mosteller’s “bulging rule”: The direction of the “bulge” indicates the direction of the power transformation of  $Y$  and/or  $X$  to straighten the relationship between them.

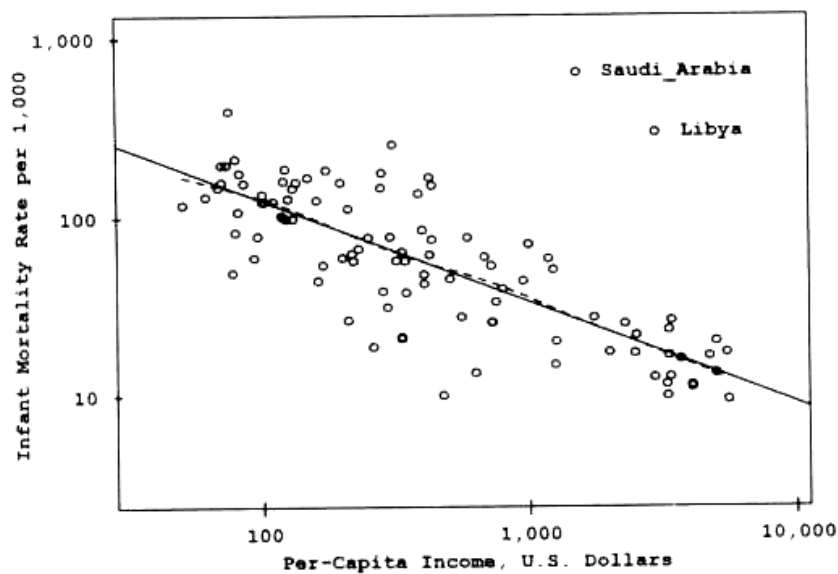
Simple monotone nonlinearity can often be corrected by a power transformation of  $X$ , of  $Y$ , or of both variables. Mosteller and Tukey’s “bulging rule” assists in the selection of a transformation.

Figure D: An application of the bulging rule. . . that also addressed positive skewness in both variables



**Figure 4.9.** Scatterplot of infant mortality rate versus income in U.S. dollars, for 101 nations *circa* 1970. The nonparametric regression shown on the plot was calculated by robust local regression. Several outlying observations are flagged.

*Source of Data:* Leinhardt and Wasserman (1978).



**Figure 4.10.** Scatterplot of  $\log_{10}$  infant mortality rate versus  $\log_{10}$  per-capita income for 101 nations. The solid line was calculated by least-squares linear regression, omitting Saudi Arabia and Libya; the broken line was calculated by robust local regression.