

## 18 Lecture 18: March 1

### Last time

- Unusual and influential data (JF chapter 11)

### Today

- HW2 deadline extends to end of this week.
- Diagnosing non-normality, non-constant error variance, and nonlinearity (JF chapter 12)
- Data transformation (JF chapter 4)

### Central Limit Theorem

Let  $X_1, X_2, \dots$  be a sequence of iid random variables whose mgfs exist in a neighborhood of 0 (that is,  $M_{X_i}(t)$  exists for  $|t| < h$ , for some positive  $h$ ). Let  $EX_i = \mu$  and  $\text{Var}X_i = \sigma^2 > 0$ . (Both  $\mu$  and  $\sigma^2$  are finite since the mgf exists.). Define  $\bar{X}_n = (1/n) \sum_{i=1}^n X_i$ . Let  $G_n(x)$  denote the cdf of  $\sqrt{n}(\bar{X}_n - \mu)/\sigma$ . Then, for any  $x$ ,  $-\infty < x < \infty$ ,

$$\lim_{n \rightarrow \infty} G_n(x) = \int_{-\infty}^x \frac{1}{\sqrt{2\pi}} e^{-y^2/2} dy$$

that is,  $\sqrt{n}(\bar{X}_n - \mu)/\sigma$  has a limiting standard normal distribution. (Refer to Casella & Berger p.237 - p.238 for a proof.)

### Delta Method

Let  $Y_n$  be a sequence of random variables that satisfies  $\sqrt{n}(Y_n - \theta) \rightarrow N(0, \sigma^2)$  in distribution. For a given function  $g$  and a specific value of  $\theta$ , suppose  $g'(\theta)$  exists and is not 0. Then

$$\sqrt{n}[g(Y_n) - g(\theta)] \rightarrow N(0, \sigma^2[g'(\theta)]^2) \text{ in distribution.}$$

(Refer to Casella & Berger p.243 for a proof using Taylor expansion.)

### Second-order Delta Method

Let  $Y_n$  be a sequence of random variables that satisfies  $\sqrt{n}(Y_n - \theta) \rightarrow N(0, \sigma^2)$  in distribution. For a given function  $g$  and a specific value of  $\theta$ , suppose that  $g'(\theta) = 0$  and  $g''(\theta)$  exists and is not 0. Then

$$\sqrt{n}[g(Y_n) - g(\theta)] \rightarrow \sigma^2 \frac{g''(\theta)}{2} \chi_1^2 \text{ in distribution.}$$

## Non-normally distributed errors

The assumption of normally distributed errors is almost always arbitrary. Nevertheless, the central limit theorem ensures that, under very broad conditions, inference based on the least-squares estimator is approximately valid in all but small samples. Why concern about non-normal errors?

- For some types of error distributions, particularly those with heavy tails, the efficiency of least-squares estimation decreases markedly.
- Highly skewed error distributions, aside from their propensity to generate outliers in the direction of the skew, compromise the interpretation of the least-squares fit. This fit is a conditional mean (of  $Y$  given the  $X$ s), and the mean is not a good measure of the center of a highly skewed distribution.
- A multimodal error distribution suggests that omission of one or more discrete explanatory variables that divide the data naturally into groups. An examination of the distribution of the residuals may motivate respecification of the model.

Note: The skewness  $\alpha_3$  is defined as  $\alpha_3 \equiv \frac{\mu_3}{(\mu_2)^{3/2}}$  where  $\mu_n$  denotes the  $n$ th central moment of a random variable  $X$ . The skewness measures the lack of symmetry in the pdf.

### Quantile-comparison plot, JF 3.1.3

*Quantile-comparison plots* are useful for comparing an empirical sample distribution with a theoretical distribution, such as the normal distribution.

Let  $P(x)$  represent the theoretical cumulative distribution function (cdf) with which we want to compare the data, that is  $P(x) = \Pr(X \leq x)$ . The quantile-comparison plot is constructed by:

1. Order the data values from smallest to largest,  $X_{(1)}, X_{(2)}, \dots, X_{(n)}$ . The  $X_{(i)}$  are called the order statistics of the sample.
2. By convention, the cumulative proportion of the data “below”  $X_{(i)}$  is given by

$$P_i = \frac{i - \frac{1}{2}}{n}$$

3. Use the inverse of the cdf to find the value  $z_i$  corresponding to the cumulative probability  $P_i$ , that is

$$z_i = P^{-1}\left(\frac{i - \frac{1}{2}}{n}\right)$$

4. Plot the  $z_i$  as horizontal coordinates against the  $X_{(i)}$  as vertical coordinates. If  $X$  is sampled from the distribution  $P$ , then  $X_{(i)} \approx z_i$ .
  - if the distributions are identical except for location, then the plot is approximately linear with nonzero intercept,  $X_{(i)} \approx \mu + z_i$

- if the distributions are identical except for scale, then the plot is approximately linear with a slope different from 1,  $X_{(i)} \approx \sigma z_i$
  - if the distributions differ both in location and scale but have the same shape, then  $X_{(i)} \approx \mu + \sigma z_i$
5. It is often helpful to place a comparison line on the plot to facilitate the perception of departures from linearity. For a normal quantile-comparison plot (comparing the distribution of the data with the standard normal distribution), we can alternatively use the median as a robust estimator of  $\mu$  and the interquartile range/1.39 as a robust estimator of  $\sigma$ .
  6. We expect some departure from linearity because of sampling variation. It therefore assists interpretation to display the expected degree of sampling error in the plot. The standard error of the order statistic  $X_{(i)}$  is

$$\text{SE}(X_{(i)}) = \frac{\hat{\sigma}}{p(z_i)} \sqrt{\frac{P_i(1 - P_i)}{n}}$$

where  $p(z_i)$  is the probability density function, pdf, corresponding to the CDF  $P(z)$ . The values along the fitted line are given by  $\hat{X}_{(i)} = \hat{\mu} + \hat{\sigma} z_i$ . An approximate 95% confidence “envelope” around the fitted line is, therefore,

$$\hat{X}_{(i)} \pm 2 \times \text{SE}(X_{(i)})$$

- Figure 18.1 plots a sample of  $n = 100$  observations from a normal distribution with mean  $\mu = 50$  and standard deviation  $\sigma = 10$ . The plotted points are reasonably linear and stay within the rough 95% confidence envelope.
- Figure 18.2 plots a sample of  $n = 100$  observations from the positively skewed chi-square distribution with 2 degrees of freedom. The positive skew of the data is reflected in points that lie *above* the comparison line in both tails of the distribution. (In contrast, the tails of negatively skewed data would lie *below* the comparison line.)
- Figure 18.3 plots a sample of  $n = 100$  observations from the heavy-tailed  $t$  distribution with 2 degrees of freedom. In this case, values in the upper tail lie above the corresponding normal quantiles, the values in the lower tail below the corresponding normal quantiles.
- Figure 18.4 shows the normal quantile-comparison plot for the distribution of infant mortality. The positive skew of the distribution is readily apparent.

## Nonconstant error variance

One of the assumptions of the regression model is that the variation of the response variable around the regression surface (the error variance) is everywhere the same:

$$\text{Var}(\epsilon) = \text{Var}(Y|x_1, \dots, x_p) = \sigma_\epsilon^2$$

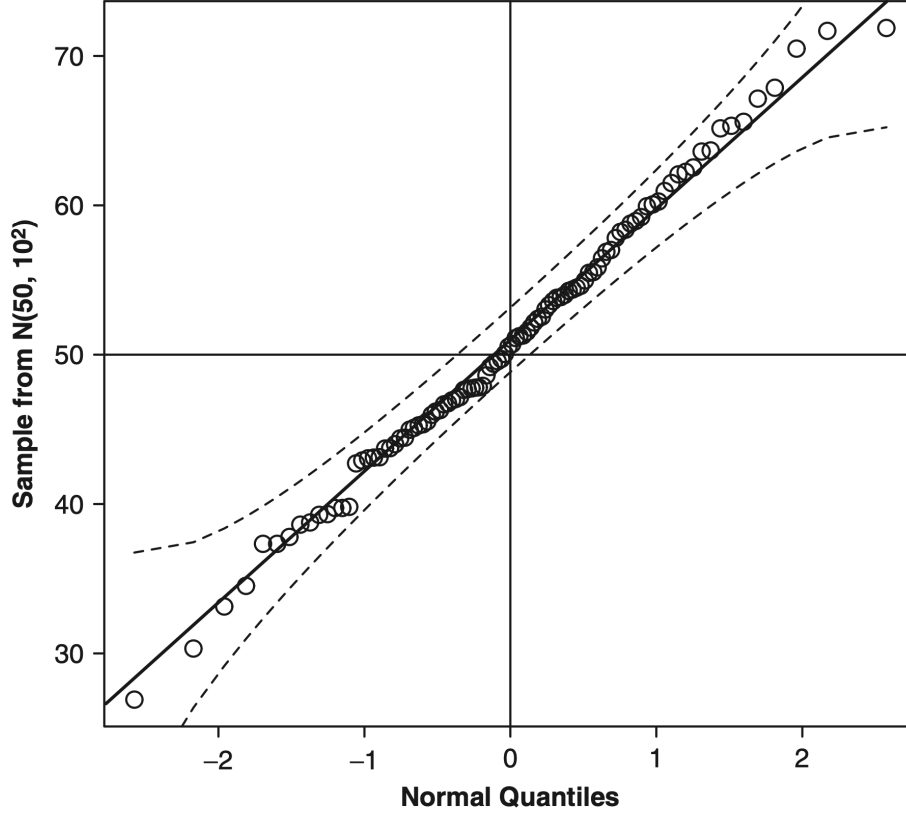


Figure 18.1: Normal quantile-comparison plot for a sample of 100 observations drawn from a normal distribution with mean 50 and standard deviation 10. The fitted line is through the quartiles of the distribution, the broken lines give a pointwise 95% confidence interval around the fit. JF Figure 3.8.

Constant error variance is often termed homoscedasticity, and similarly, nonconstant error variance is termed heteroscedasticity. We detect nonconstant error variances through graphical methods.

### Residual plots

Because the least square residuals have unequal variance even when the constant variance assumption is correct:

$$\text{Var}(\hat{\epsilon}_i) = \sigma^2(1 - h_i).$$

It is preferable to plot studentized residuals against fitted values. A pattern of changing spread is often more easily discerned in a plot of absolute studentized residuals,  $|\hat{\epsilon}_i^*|$ , or squared studentized residuals,  $\hat{\epsilon}_i^{*2}$ , against  $\hat{Y}$ . If the values of  $\hat{Y}$  are all positive, then we can plot  $\log|\hat{\epsilon}_i^*|$  against  $\log \hat{Y}$ . Figure 18.5 shows a plot of studentized residuals against fitted values and spread-level plot of studentized residuals, several points with negative fitted values were omitted. It is apparent from both graphs that the residual spread tends to increase with the level of the response, suggesting a violation of constant error variance assumption.

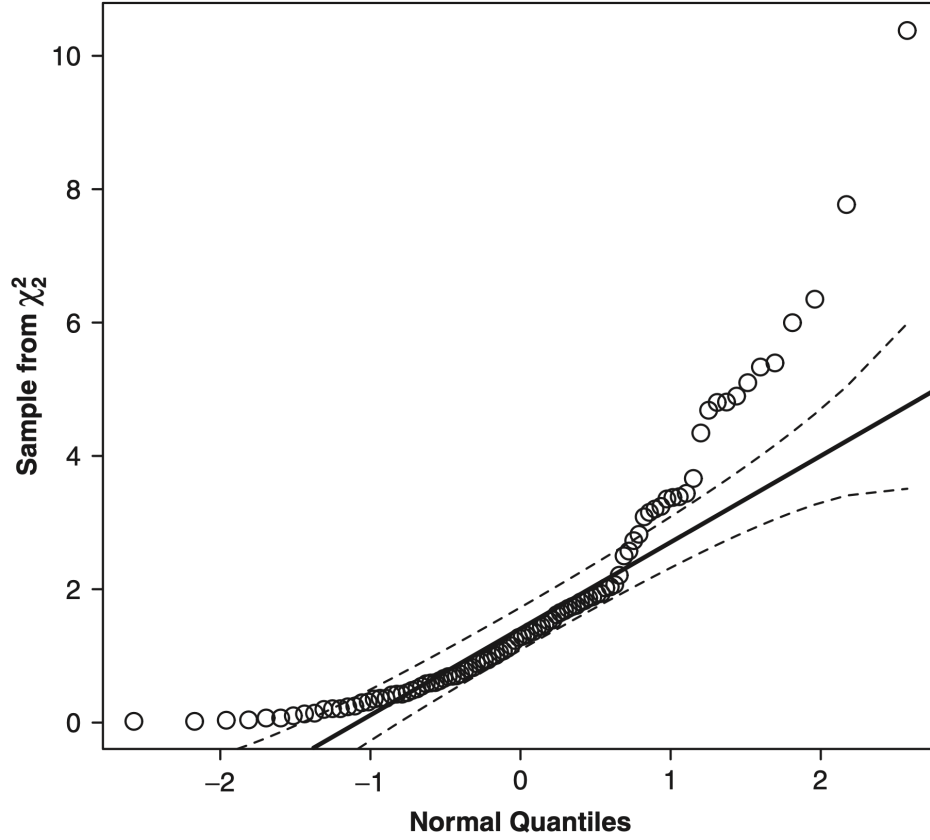


Figure 18.2: Normal quantile-comparison plot for a sample of 100 observations drawn from the positively skewed chi-square distribution with 2 degrees of freedom. JF Figure 3.9.

### Weighted-least-squares estimation

Weighted-least-squares (WLS) regression provides an alternative approach to estimation in the presence of nonconstant error variance. Suppose that the errors from the linear regression model  $\mathbf{Y} = \mathbf{X}\beta + \epsilon$  are independent and normally distributed, with zero means but *different* variances:  $\epsilon_i \sim N(0, \sigma_i^2)$ . Suppose further that the variances of the errors are known up to a constant of proportionality  $\sigma_\epsilon^2$ , so that  $\sigma_i^2 = \sigma_\epsilon^2/w_i^2$ . Then the likelihood for the model is

$$L(\beta, \sigma_\epsilon^2) = \frac{1}{(2\pi)^{n/2} |\Sigma|^{1/2}} \exp \left[ -\frac{1}{2} (\mathbf{Y} - \mathbf{X}\beta)^T \Sigma^{-1} (\mathbf{Y} - \mathbf{X}\beta) \right]$$

where  $\Sigma$  is the covariance matrix of the errors,

$$\Sigma = \sigma_\epsilon^2 \times \text{diag}\{1/w_1^2, \dots, 1/w_n^2\} \equiv \sigma_\epsilon^2 \mathbf{W}^{-1}$$

The maximum-likelihood estimators of  $\beta$  and  $\sigma_\epsilon^2$  are then

$$\begin{aligned} \hat{\beta} &= (\mathbf{X}^T \mathbf{W} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{W} \mathbf{Y} \\ \hat{\sigma}_\epsilon^2 &= \frac{\sum (w_i \hat{\epsilon}_i)^2}{n} \end{aligned}$$

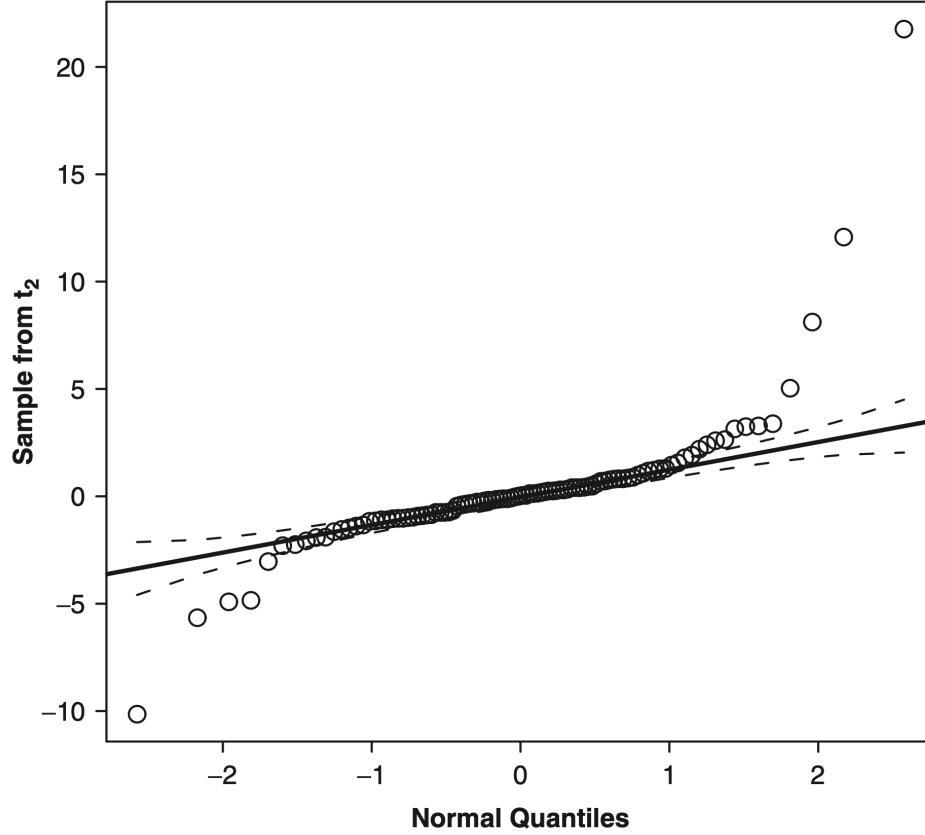


Figure 18.3: Normal quantile-comparison plot for a sample of 100 observations drawn from heavy-tailed  $t$ -distribution with 2 degrees of freedom. JF Figure 3.10.

#### Correcting OLS standard errors for nonconstant variance

The covariance matrix of the ordinary-least-squares (OLS) estimator is

$$\begin{aligned}\mathbf{Var}(\hat{\beta}) &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Var}(\mathbf{Y}) \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \\ &= \sigma_\epsilon^2 (\mathbf{X}^T \mathbf{X})^{-1}\end{aligned}$$

under the standard assumptions, including the assumption of constant error variance,  $\mathbf{Var}(\mathbf{Y}) = \sigma_\epsilon^2 \mathbf{I}_n$ . If, however, the errors are heteroscedastic but independent then  $\Sigma \equiv \mathbf{Var}(\mathbf{Y}) = \text{diag}\{\sigma_1^2, \dots, \sigma_n^2\}$ , and

$$\mathbf{Var}(\hat{\beta}) = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \Sigma \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1}$$

White (1980) shows that the following is a consistent estimator of  $\mathbf{Var}(\hat{\beta})$

$$\tilde{\mathbf{Var}}(\hat{\beta}) = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \hat{\Sigma} \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1}$$

with  $\hat{\Sigma} = \text{diag}\{\hat{\sigma}_1^2, \dots, \hat{\sigma}_n^2\}$ , where  $\hat{\sigma}_i^2$  is the OLS residual for observation  $i$ .

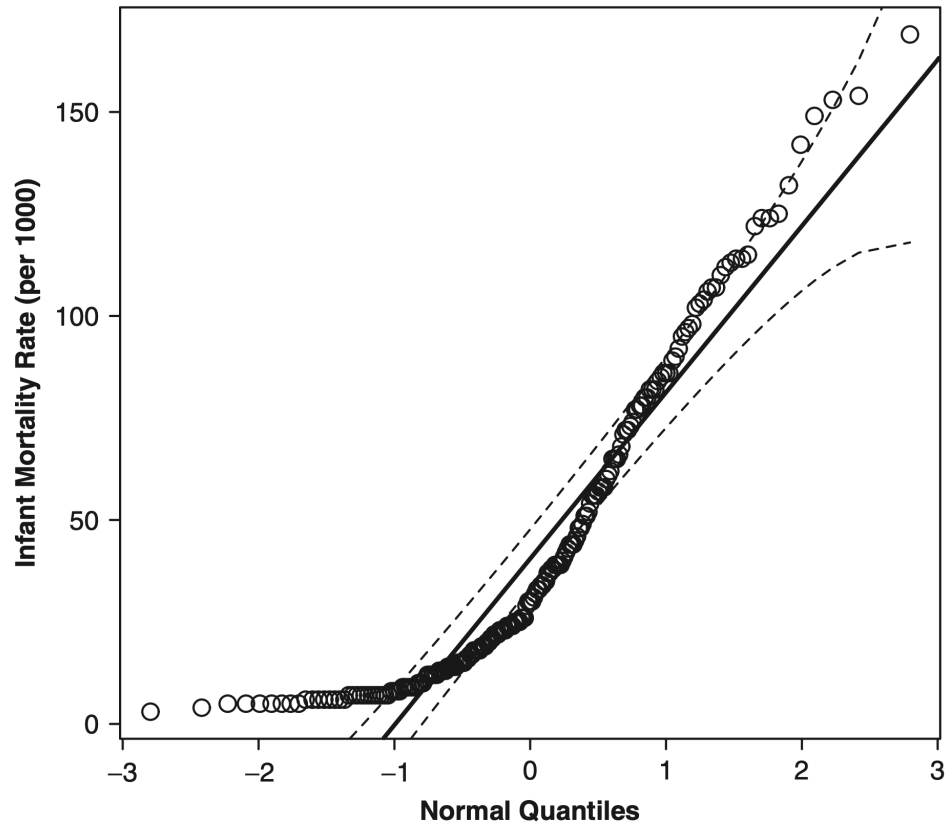


Figure 18.4: Normal quantile-comparison plot for the distribution of infant mortality. Note the positive skew. JF Figure 3.11.

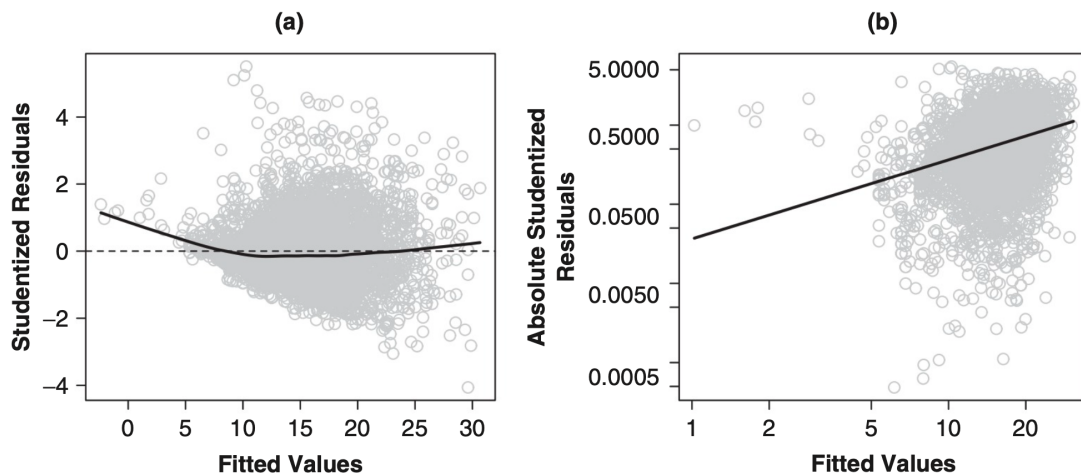


Figure 18.5: (a) Plot of studentized residuals versus fitted values and (b) spread-level plot for studentized residuals. JF Figure 12.3.

Subsequent work suggested small modifications to White's coefficient-variance estimator,

and in particular simulation studies by Long and Ervin (2000) support the use of

$$\tilde{\text{Var}}^*(\hat{\beta}) = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \hat{\Sigma}^* \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1}$$

where  $\hat{\Sigma}^* = \text{diag}\{\hat{\sigma}_i^2/(1 - h_i)^2\}$  and  $h_i$  is the hat-value associated with observation  $i$ . In large samples, where  $h_i$  is small, the distinction between  $\tilde{\text{Var}}(\hat{\beta})$  and  $\tilde{\text{Var}}^*(\hat{\beta})$  essentially disappears.

A rough *rule* is that nonconstant error variance seriously degrades the least-squares estimator only when the ratio of the largest to smallest variance is about 10 or more (or, more conservatively, about 4 or more).

## Data transformation

### The family of powers and Roots

A particularly useful group of transformations is the “family” of powers and roots:

$$X \rightarrow X^p$$

where the arrow indicates that we intend to replace  $X$  with the transformed variable  $X^p$ . If  $p$  is negative, then the transformation is an inverse power. For example,  $X^{-1} = 1/X$ . If  $p$  is a fraction, then the transformation represents a root. For example,  $X^{1/3} = \sqrt[3]{X}$ .

It is more convenient to define the family of power transformations in a slightly more complex manner, called the Box-Cox family of transformations (introduced in a seminal paper on transformations by Box & Cox, 1964):

$$X \rightarrow X^{(p)} = \frac{X^p - 1}{p}$$

Because  $X^{(p)}$  is a linear function of  $X^p$ , the two transformations have the same essential effect on the data, but, as is apparent in Figure 18.6

- Dividing by  $p$  preserves the direction of  $X$ , which otherwise would be reversed when  $p$  is negative.
- The transformations  $X^{(p)}$  are “matched” above  $X = 1$  both in level and in slope:
  1.  $X^{(p)} = 0$ , for all values of  $p$
  2. each transformation has a slope of 1 at  $X = 1$ .
- Descending the “ladder” of powers and roots towards  $X^{(-1)}$  compresses the large values of  $X$  and spreads out the small ones. Ascending the ladder of powers and roots towards  $X^{(2)}$  has the opposite effect. As  $p$  moves further from  $p = 1$  (i.e. no transformation) in either direction, the transformation grows more powerful, increasingly “bending” the data.



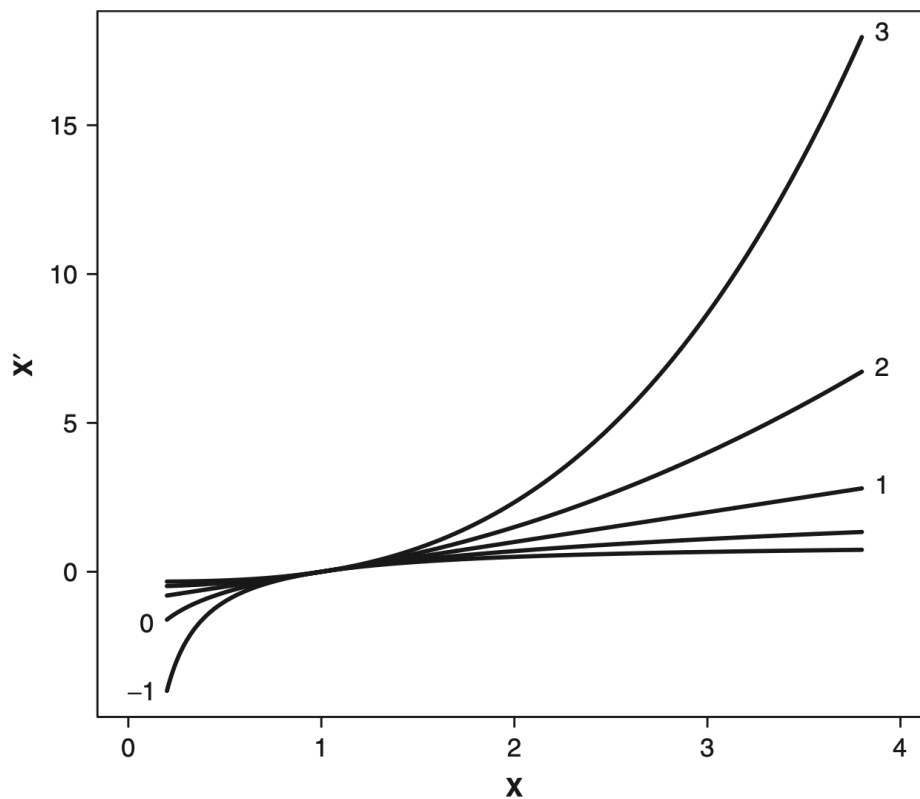


Figure 18.6: The Box-Cox family of power transformations  $X'$  of  $X$ . The curve labeled  $p$  is the transformation  $X^{(p)}$ , that is  $(X^p - 1)/p$ ;  $X^{(0)}$  is  $\log_e(X)$ . JF Figure 4.1.

- The power transformation  $X^0$  is useless because it changes all values to 1, but we can think of the log transformation as a kind of “zeroth” power:

$$\lim_{p \rightarrow 0} \frac{X^p - 1}{p} = \log_e X$$

and by convention,  $X^{(0)} \equiv \log_e X$ .