

20 Lecture 20: March 17

Last time

- Unusual and influential data (JF chapter 11)

Today

- Anonymous internal midterm evaluations on canvas
- Added-variable plots
- Should unusual data be discarded
- Non-normality

Added-variable plots

Unlike the case of SLR, the scatterplot with the response variable and one predictor gives only the marginal effect in MLR. Instead, the added-variable plot (also called a partial-regression plot or a partial-regression leverage plot) gives a graphical inspection over each dimension.

Let $\tilde{Y}_i^{(1)}$ represent the residuals from the least-squares regression of Y on all the X s except X_1 , in other words, the residuals from the following fitted regression equation:

$$Y_i = \tilde{\beta}_0^{(1)} + \tilde{\beta}_2^{(1)} X_{i2} + \cdots + \tilde{\beta}_p^{(1)} X_{ip} + \tilde{Y}_i^{(1)}$$

where the parenthetical superscript (1) indicates the omission of X_1 from the right-hand side of the regression equation. Likewise, $\tilde{X}_i^{(1)}$ is the residual from the least-squares regression of X_1 on all the other X s:

$$X_{i1} = \tilde{\beta}_0^{(1)} + \tilde{\beta}_2^{(1)} X_{i2} + \cdots + \tilde{\beta}_p^{(1)} X_{ip} + \tilde{X}_i^{(1)}$$

Then, the residuals $\tilde{Y}_i^{(1)}$ and $\tilde{X}_i^{(1)}$ have the following interesting properties:

1. The slope from the least-squares regression of $\tilde{Y}_i^{(1)}$ on $\tilde{X}_i^{(1)}$ is simply the least-squares slope $\hat{\beta}_1$ from the *full* multiple regression.
2. The residuals from the simple regression of $\tilde{Y}_i^{(1)}$ on $\tilde{X}_i^{(1)}$ are the same as those from the full regression, that is

$$\tilde{Y}_i^{(1)} = \hat{\beta}_1 \tilde{X}_i^{(1)} + \hat{\epsilon}_i$$

3. The variation of $\tilde{X}_i^{(1)}$ is the *conditional variation* of X_1 holding the other X s constant.

Figure 20.1 shows that the conditional variation is smaller than its marginal variation – much smaller when X_1 is strongly collinear with other X s,

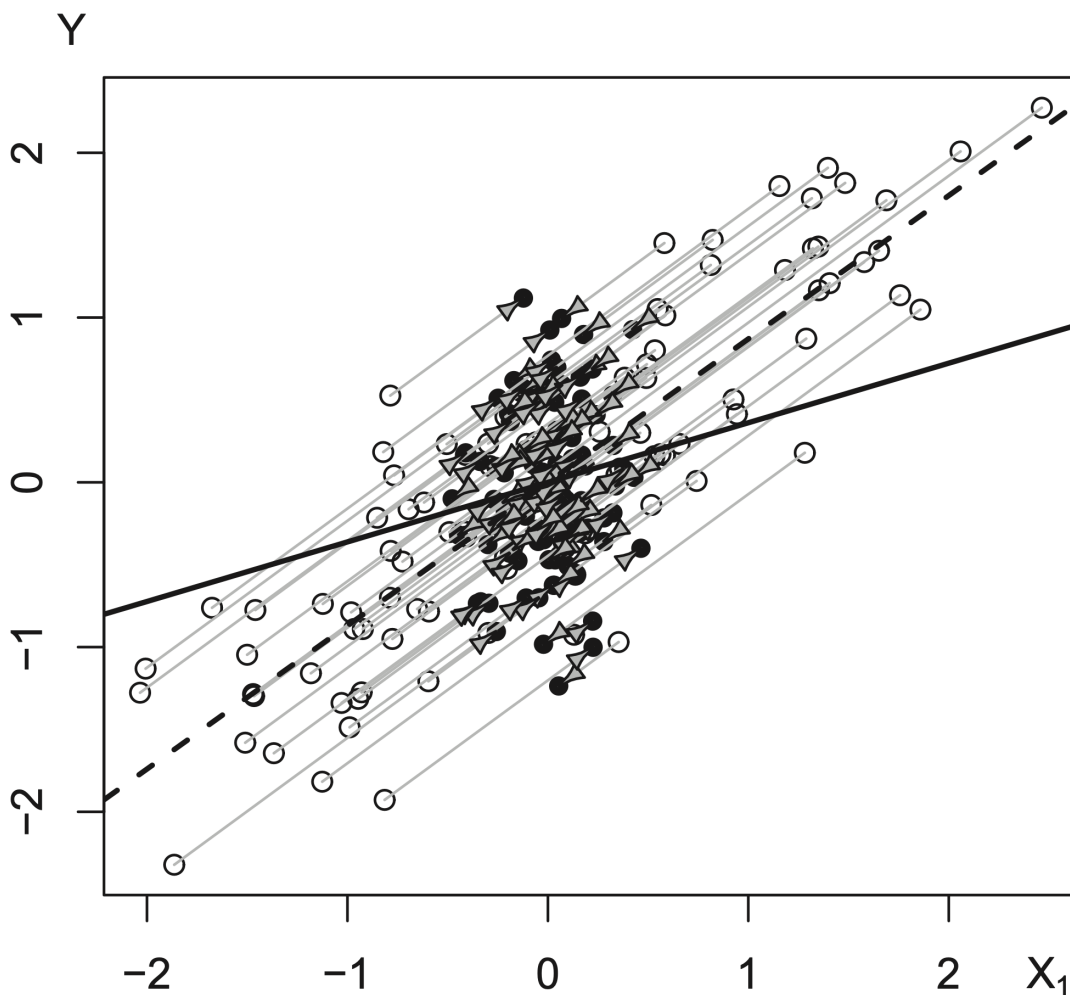


Figure 20.1: The marginal scatterplot (open circles) for Y and X_1 superimposed on the added-variable plot (filled circles) for X_1 in the regression of Y on X_1 and X_2 . The variables Y and X_1 are centered at their means to facilitate the comparison of the two sets of points. The arrows show how the points in the marginal scatterplot map into those in the AV plot. In this contrived data set, X_1 and X_2 are highly correlated ($r_{12} = 0.98$), and so the conditional variation in X_1 (represented by the horizontal spread of the filled points) is much less than its marginal variation (represented by the horizontal spread of the open points). The broken line gives the slope of the marginal regression of Y on X_1 alone, while the solid line gives the slope $\hat{\beta}_1$ of X_1 in the MLR of Y on both X s. JF Figure 11.9.

Figure 20.2 illustrates the added-variable plots using the Duncan's data.



Figure 20.2: Added-variable plots for Duncan’s regression of occupational prestige on the (a) income and (b) education levels of 45 US occupations in 1950. Three unusual observations, *ministers*, *conductors*, and *railroadengineers*, are identified on the plots. The added-variable plot for the intercept $\hat{\beta}_0$ is not shown. JF Figure 11.10.

The added-variable plot for income in Figure 20.2(a) reveals three observations that exert substantial leverage on the income coefficient:

- *minister*, whose income is unusually low given the educational level of the occupation
- *conductor*, whose income is unusually high given education
- *railroad engineer*, whose income is relatively high given education.

Remember that the horizontal variable in this added-variable plot is the residual from the regression of income on education, and thus values far from 0 in this direction are for occupations with incomes that are unusually high or low given their levels of education.

Should unusual data be discarded?

In practice, although problematic data should not be ignored, they also should not be deleted automatically and without reflection:

- It is important to investigate *why* an observation is unusual. Truly “bad” data (e.g., an error in data entry) can often be corrected or, if correction is not possible, thrown away. When a discrepant data point is correct, we may be able to understand why the observation is unusual. For Duncan’s data, for example, it makes sense that ministers enjoy prestige not accounted for by the income and educational levels of the occupation and for a reason not shared by other occupations. In a case like this, where an outlying observation has characteristics that render it unique, we may choose to set it aside from the rest of the data.

- Alternatively, outliers, high-leverage points, or influential data may motivate model respecification, and the pattern of unusual data may suggest the introduction of additional explanatory variables. We noticed, for example, that both conductors and railroad engineers had high leverage in Duncan's regression because these occupations combined relatively high income with relatively low education. Perhaps this combination of characteristics is due to a high level of unionization of these occupations in 1950, when the data were collected. If so, and if we can ascertain the levels of unionization of all of the occupations, we could enter this as an explanatory variable, perhaps shedding further light on the process determining occupational prestige.
- Except in clear-cut cases, we are justifiably reluctant to delete observations or to re-specify the model to accommodate unusual data. Some researchers reasonably adopt alternative estimation strategies, such as robust regression, which continuously down-weights outlying data rather than simply discarding them. Because these methods assign zero or very small weight to highly discrepant data, however, the result is generally not very different from careful application of least squares, and, indeed, robust-regression weights can be used to identify outliers.
- Finally, in large samples, unusual data substantially alter the results only in extreme instances. Identifying unusual observations in a large sample, therefore, should be regarded more as an opportunity to learn something about the data not captured by the model that we have fit, rather than as an occasion to reestimate the model with the unusual observations removed.

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 α_3 is defined as $\alpha_3 \equiv \frac{\mu_3}{(\mu_2)^{3/2}}$ where μ_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 μ and the interquartile range/1.39 as a robust estimator of σ .
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 20.3 plots a sample of $n = 100$ observations from a normal distribution with mean $\mu = 50$ and standard deviation $\sigma = 10$.

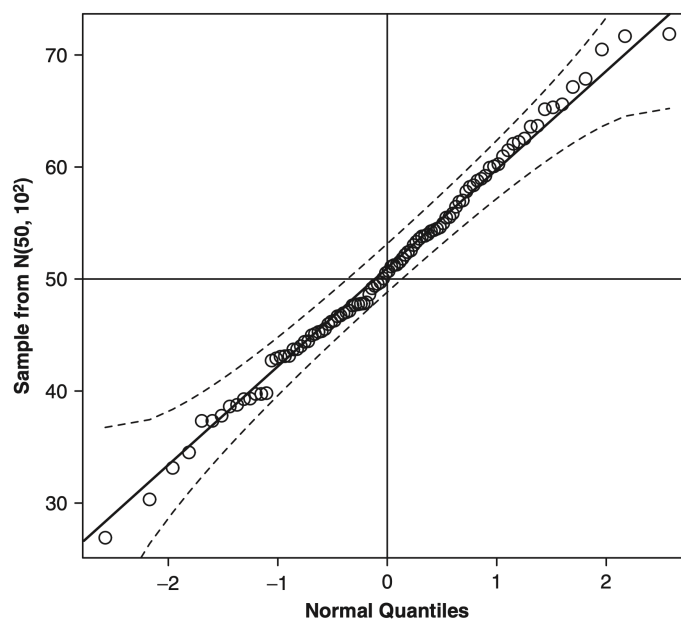


Figure 20.3: Normal quantile-comparison plot for a sample of 100 observations drawn from a normal distribution with mean 50 and standard deviation 10. The fitted line is through the quantiles of the distribution, the broken lines give a pointwise 95% confidence interval around the fit. JF Figure 3.8.

The plotted points are reasonably linear and stay within the rough 95% confidence envelope.

- Figure 20.4 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.)

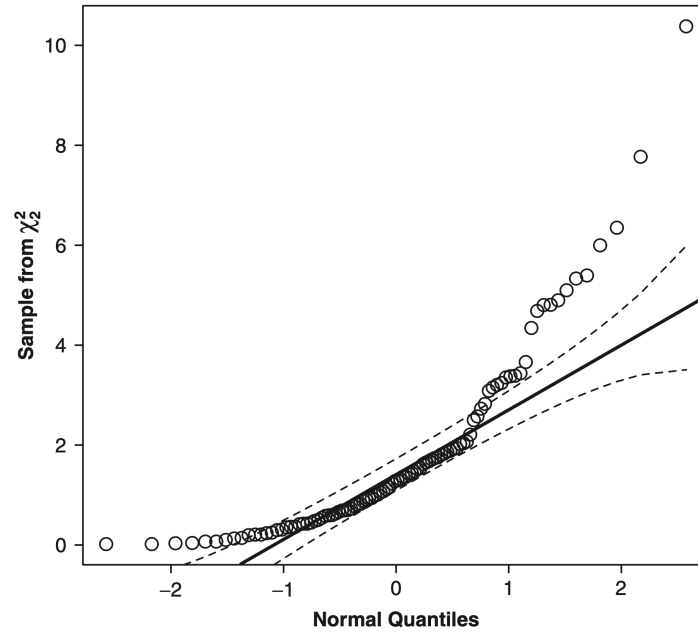


Figure 20.4: Normal quantile-comparison plot for a sample of 100 observations drawn from the positively skewed chi-square distribution with 2 degrees of freedom. JF Figure 3.9.

- Figure 20.5 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.

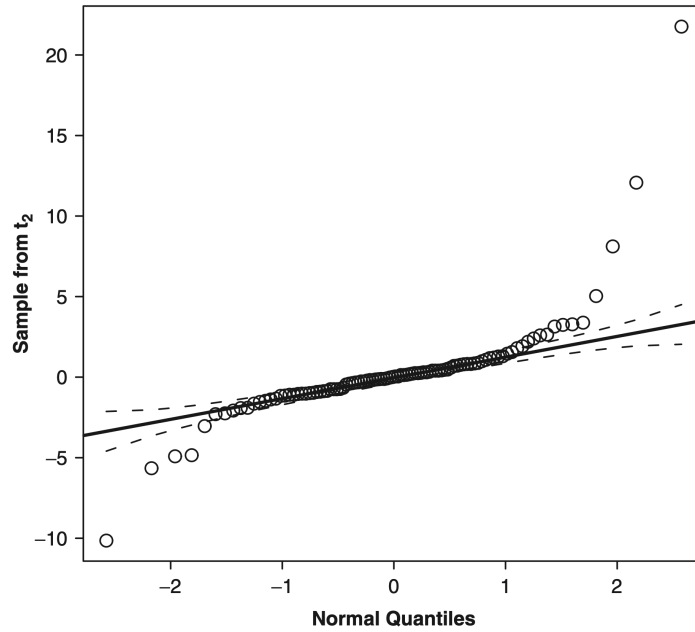


Figure 20.5: Normal quantile-comparison plot for a sample of 100 observations drawn from heavy-tailed t -distribution with 2 degrees of freedom. JF Figure 3.10.

- Figure 20.6 shows the normal quantile-comparison plot for the distribution of infant mortality. The positive skew of the distribution is readily apparent.

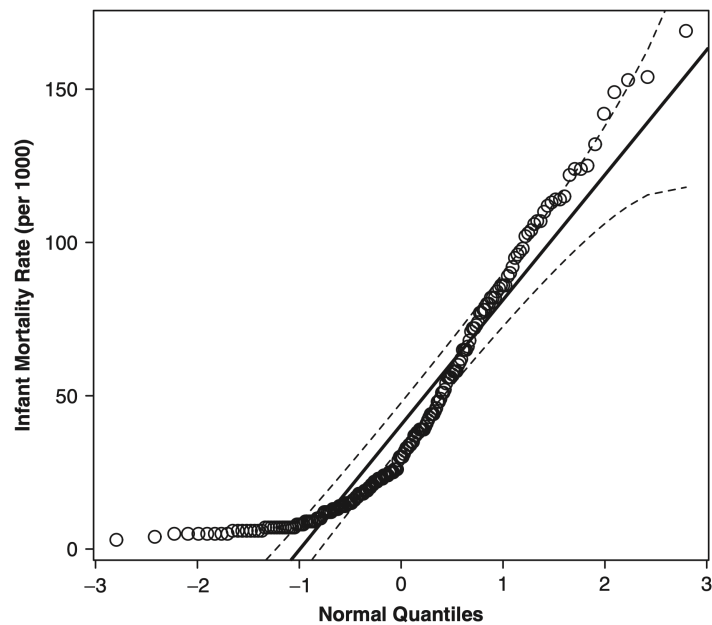


Figure 20.6: Normal quantile-comparison plot for the distribution of infant mortality. Note the positive skew. JF Figure 3.11.

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:

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

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 20.7 shows a plot of studentized residuals against fitted values and spread-level plot of studentized residuals, several points with negative fitted values were omitted.

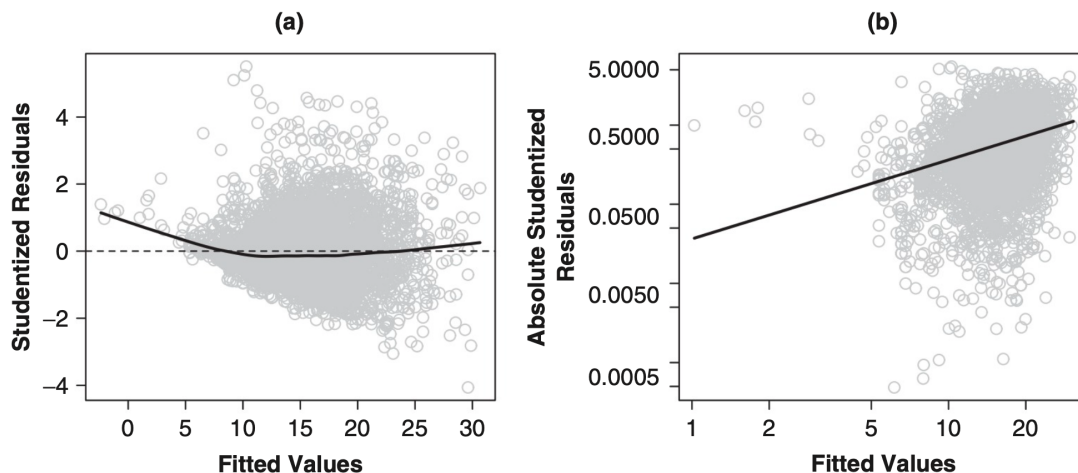


Figure 20.7: (a) Plot of studentized residuals versus fitted values and (b) spread-level plot for studentized residuals. JF Figure 12.3.

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.

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 σ_ϵ^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 Σ 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 β and σ_ϵ^2 are then

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

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 .

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

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

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 $\tilde{\mathbf{Var}}(\hat{\beta})$ and $\tilde{\mathbf{Var}}^*(\hat{\beta})$ 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).

Nonlinearity

If $\mathbf{E}(\mathbf{Y}|\mathbf{X})$ is not linear in \mathbf{X} (in other words, $\mathbf{E}(\epsilon|\mathbf{X}) \neq 0$ for some x), $\hat{\beta}$ may be biased and inconsistent. Usually we employ “linearity by default” but we should try to make sure this is appropriate: **detect** non-linearities and **model** them accurately.

Lowess smoother, JF 2.3

We can employ local averaging plots to help with diagnostics. Lowess method is in many respects similar to local-averaging smoothers, except that instead of computing an average Y -value within the neighborhood of a focal x , the lowess smoother computes a *fitted* value based on a locally weighted least-squares line, giving more weight to observations in the neighborhood that are close to the focal x than to those relatively far away. The name “lowess” is an acronym for *locally weighted scatterplot smoother* and is sometimes rendered as *loess*, for *local regression*.

Component-plus-residual plots

Added-variable plots, introduced for detecting influential data, can reveal nonlinearity. However, the added-variable plots are not always useful for locating a transformation:

- The added-variable plot adjusts X_j for the other X s.
- The *unadjusted* X_j is transformed in respecifying the model.

Moreover, Cook (1998, Section 14.5) shows that added-variable plots are biased toward linearity when the correlations among the explanatory variables are large.

Component-plus-residual plots (also called *partial-residual plots*) are often an effective alternative. The component-plus-residual plots are not as suitable as added-variable plots for revealing leverage and influence, though. The component-plus-residual plots are constructed by

1. Compute residuals from full regression:

$$\hat{\epsilon}_i = Y_i - \hat{Y}_i$$

2. Compute “linear component” of the partial relationship:

$$C_i = \hat{\beta}_j X_{ij}$$

3. Add linear component to residual to get partial residual for the j th explanatory variable

$$\hat{\epsilon}_i^{(j)} = \hat{\epsilon}_i + C_i = \hat{\epsilon}_i + \hat{\beta}_j X_{ij}$$

4. Plot $\hat{\epsilon}_i^{(j)}$ against $X_{.j}$

Figure 20.8 shows the component-plus-residual plots for the regression of log wages on variables (age, education and sex) of the 1994 wave of Statistics Canada’s Survey of Labour and Income Dynamics (SLID) data. The SLID data set includes 3997 employed individuals who were between 16 and 65 years of age and who resided in Ontario.

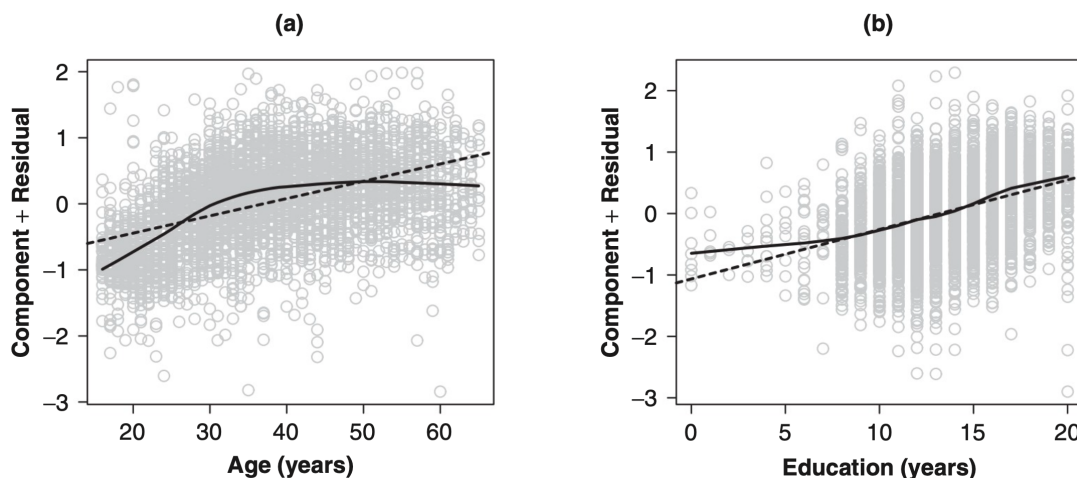


Figure 20.8: Component-plus-residual plots for age and education in SLID regression of log wages on these variables and sex. The solid lines are for lowess smooths with spans of 0.4, and the broken lines are for linear least-squares fits. JF Figure 12.6.

Data transformation

The family of powers and Roots, JF 4.1

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 20.9

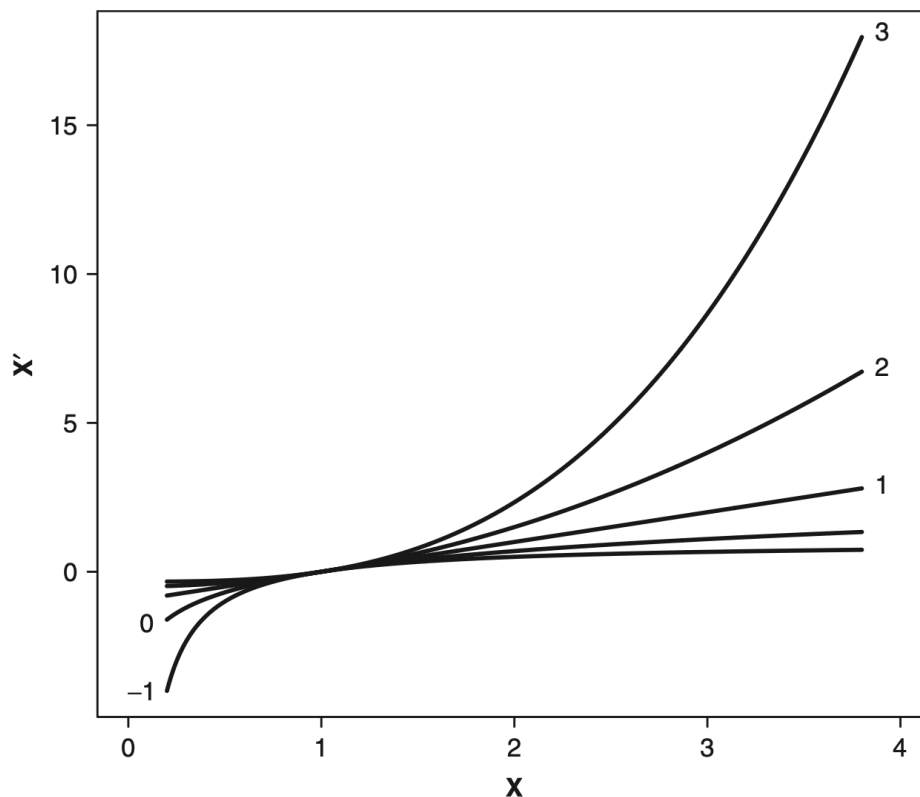


Figure 20.9: The Box-Cox family of power transformations X' of X . The curve labeled p is the transformation $X^{(p)}$, that is $(X^p - 1)/p$; $X^{(0)}$ is $\log_e(X)$. JF Figure 4.1.

- 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. $1^{(p)} = 0$, 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.
- 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, $X^{(0)} \equiv \log_e X$.

Box-Cox transformation of Y

Box and Cox (1964) suggested a power transformation of Y with the object of normalizing the error distribution, stabilizing the error variance, and straightening the relationship of Y to the X s. The general Box-Cox model is

$$Y_i^{(\lambda)} = \beta_0 + \beta_1 X_{i1} + \cdots + \beta_p X_{ip} + \epsilon_i$$

where $\epsilon_i \stackrel{iid}{\sim} N(0, \sigma_\epsilon^2)$, and

$$Y_i^{(\lambda)} = \begin{cases} \frac{Y_i^\lambda - 1}{\lambda} & \text{for } \lambda \neq 0 \\ \log_e Y_i & \text{for } \lambda = 0 \end{cases}$$

Note: in statistics, \log_e is often written as \log .

For a particular choice of λ , the conditional maximized log-likelihood (see JF 12.5.1 p.324 footnote 55) is

$$\begin{aligned} \log_e L(\beta_0, \beta_1, \dots, \beta_p, \sigma_\epsilon^2 | \lambda) &= -\frac{n}{2}(1 + \log_e 2\pi) \\ &\quad - \frac{n}{2} \log_e \hat{\sigma}_\epsilon^2(\lambda) + (\lambda - 1) \sum_{i=1}^n \log_e Y_i \end{aligned}$$

where $\hat{\sigma}_\epsilon^2(\lambda) = \sum \hat{\epsilon}_i^2(\lambda)/n$ and where $\hat{\epsilon}_i(\lambda)$ are the residuals from the least-squares regression of $Y^{(\lambda)}$ on X s. The least-squares coefficients from this regression are the maximum-likelihood estimates of β s conditional on the values of λ .

A simple procedure for finding the maximum-likelihood estimator $\hat{\lambda}$ is to evaluate the maximized $\log_e L$ (called the profile log-likelihood) for a range of values of λ . To test: $H_0 : \lambda = 1$, calculated the likelihood-ratio statistic

$$G_0^2 = -2[\log_e L(\lambda = 1) - \log_e L(\lambda = \hat{\lambda})]$$

which is asymptotically distributed as χ_1^2 with one degree of freedom under H_0 . A 95% confidence interval for λ includes those values for which

$$\log_e L(\lambda) > \log_e L(\lambda = \hat{\lambda}) - 1.92$$

The number 1.92 comes from $\frac{1}{2}\chi_{1,0.05}^2 = 0.5 \times 1.96^2$.

Figure 20.10 shows a plot of the profile log-likelihood against λ for the original SLID regression of composite hourly wages on sex, age, and education. The maximum-likelihood estimate of λ is $\hat{\lambda} = 0.09$, and a 95% confidence interval runs from 0.04 to 0.13. Although 0 is outside of the CI (confidence interval), it is essentially the same as log transformation of wages (the correlation between log wages and wages^{0.09} is 0.9996).

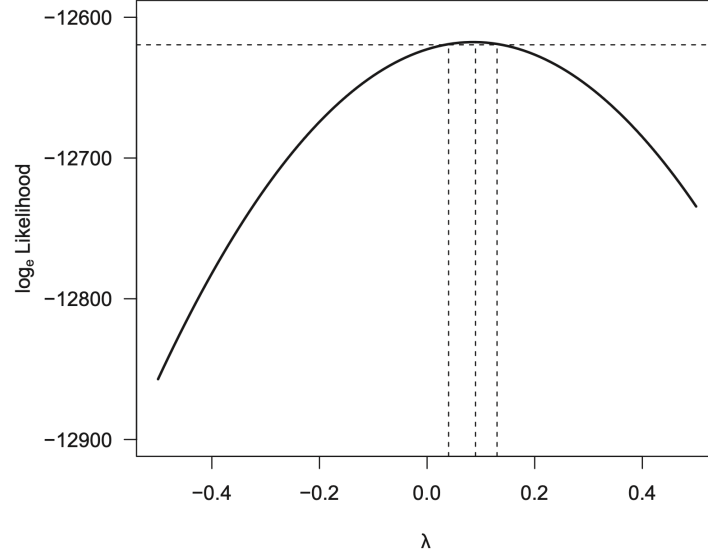


Figure 20.10: Box-Cox transformations for the SLID regression of wages on sex, age, and education. The maximized (profile) log-likelihood is plotted against the transformation parameter λ . The intersection of the line near the top of the graph with the profile log-likelihood curve marks off a 95% confidence interval for λ . The maximum of the log-likelihood corresponds to the MLE of λ . JF Figure 12.14.

Box-Tidwell transformation of X s

Now, consider the model

$$Y_i = \beta_0 + \beta_1 X_{i1}^{\gamma_1} + \cdots + \beta_p X_{ip}^{\gamma_p} + \epsilon_i$$

where the errors are independently distributed as $\epsilon_i \stackrel{iid}{\sim} N(0, \sigma_\epsilon^2)$ and all the X_{ij} are positive.

The parameters of this model ($\beta_0, \beta_1, \dots, \beta_p, \gamma_1, \dots, \gamma_p$, and σ_ϵ^2) could be estimated by general nonlinear least squares. Box and Tidwell (1962) suggested the following computationally more efficient procedure (also yields a constructed-variable diagnostic):

1. Regress Y on X_1, \dots, X_p , obtaining $\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_p$. (“Regress A on B s” is the same as “fitting the linear regression model with A as the response variable and B s as the explanatory variables”.)
2. Regress Y on X_1, \dots, X_p and the constructed variables $X_1 \log_e X_1, \dots, X_p \log_e X_p$ (again, by fitting the model of $Y = \beta_0 + \beta_1 X_1 + \cdots + \beta_p X_p + \delta_1 X_1 \log_e X_1 + \cdots + \delta_p X_p \log_e X_p + \epsilon_i$) to obtain $\tilde{\beta}_0, \tilde{\beta}_1, \dots, \tilde{\beta}_p, \tilde{\delta}_1, \dots, \tilde{\delta}_p$. In general $\hat{\beta}_i \neq \tilde{\beta}_i$. (The constructed variables result from the first-order Taylor-series approximation to $X_j^{\gamma_j}$ evaluated at $\gamma_j = 1$: $X_j^{\gamma_j} \approx X_j + (\gamma_j - 1)X_j \log_e X_j$.)
3. The constructed variable $X_j \log_e X_j$ can be used to assess the need for a transformation of X_j by testing the null hypothesis $H_0 : \delta_j = 0$. Added-variable plots for the constructed variables are useful for assessing leverage and influence on the decision to transform the X s.

4. A preliminary estimate of the transformation parameter γ_j (not the MLE) is

$$\tilde{\gamma}_j = 1 + \frac{\tilde{\delta}_j}{\hat{\beta}_j}$$

where $\tilde{\delta}_j$ is from step 2 and $\hat{\beta}_j$ is from step 1.

Polynomial regression

A machinery of multiple regression to fit non-linear relationships between predictor(s) and response.

- Linear: $y = \beta_0 + \beta_1 x + \epsilon$
- Quadratic: $y = \beta_0 + \beta_1 x + \beta_2 x^2 + \epsilon$
- Cubic: $y = \beta_0 + \beta_1 x + \beta_2 x^2 + \beta_3 x^3 + \epsilon$
- k^{th} order polynomial: $y = \beta_0 + \beta_1 x + \beta_2 x^2 + \cdots + \beta_k x^k + \epsilon$

Question:

Does quadratic model provide a significantly better fit than linear model?

Solution: Test $H_0 : \beta_2 = 0$ vs. $H_a : \beta_2 \neq 0$.

Alternatively, compare the corresponding adjusted- R^2 values.