

Unit 3: Linear models: Misspecification

Eugene Katsevich

October 11, 2021

In our discussion of linear model inference in Unit 2, we assumed the normal linear model throughout:

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}, \quad \text{where } \boldsymbol{\epsilon} \sim N(\mathbf{0}, \sigma^2 \mathbf{I}_n). \quad (1)$$

In this unit, we will discuss what happens when this model is misspecified:

- Non-normality (Section 1): $\boldsymbol{\epsilon} \sim (0, \sigma^2 \mathbf{I}_n)$ but not $N(0, \sigma^2 \mathbf{I}_n)$.
- Heteroskedastic errors (Section 2): $\epsilon_i \stackrel{\text{ind}}{\sim} N(0, \sigma_i^2)$, where it is not the case that $\sigma_1^2 = \dots = \sigma_n^2$.
- Correlated errors (Section 3): It is not the case that $(\epsilon_1, \dots, \epsilon_n)$ are independent.
- Model bias (Section 4): It is not the case that $\mathbb{E}[\mathbf{y}] = \mathbf{X}\boldsymbol{\beta}$ for some $\boldsymbol{\beta} \in \mathbb{R}^p$.
- Outliers (Section 5): For one or more i , it is not the case that $y_i \sim N(\mathbf{x}_{i*}^T \boldsymbol{\beta}, \sigma^2)$.

For each type of misspecification, we will discuss its origins, consequences, detection, and fixes (Sections 1-5). We conclude with an R demo (Section 7).

1 Non-normality

1.1 Origin

Non-normality occurs when the distribution of $y|\mathbf{x}$ is either skewed or has heavier tails than the normal distribution. This may happen, for example, if there is some discreteness in y .

1.2 Consequences

Non-normality is the most benign of linear model misspecifications. While we derived linear model inferences under the normality assumption, all the corresponding statements hold asymptotically without this assumption. Recall Homework 2 Question 1, or take for example the simpler problem of estimating the mean μ of a distribution based on n samples from it: We can test $H_0 : \mu = 0$ and build a confidence interval for μ even if the underlying distribution is not normal. So if n is relatively large and p is relatively small, you need not worry too much. If n is small and the errors are highly skewed or heavy-tailed, there might be an issue.

1.3 Detection

Non-normality is a property of the error-terms ϵ_i . We do not observe these directly, but we can approximate these using the residuals

$$\hat{\epsilon}_i = y_i - \mathbf{x}_{i*}^T \hat{\boldsymbol{\beta}}. \quad (2)$$

Recall from Unit 2 that $\text{Var}[\hat{\epsilon}] = \sigma^2(\mathbf{I} - \mathbf{H})$. Letting h_i be the i th diagonal entry of \mathbf{H} , it follows that $\hat{\epsilon}_i \sim (0, \sigma^2(1 - h_i))$. The *standardized residuals* are defined as

$$r_i = \frac{\hat{\epsilon}_i}{\hat{\sigma}\sqrt{1 - h_i}}. \quad (3)$$

Under normality, we would expect $r_i \sim N(0, 1)$. We can therefore assess normality by producing a histogram or normal QQ-plot of these residuals (see Figure 1).

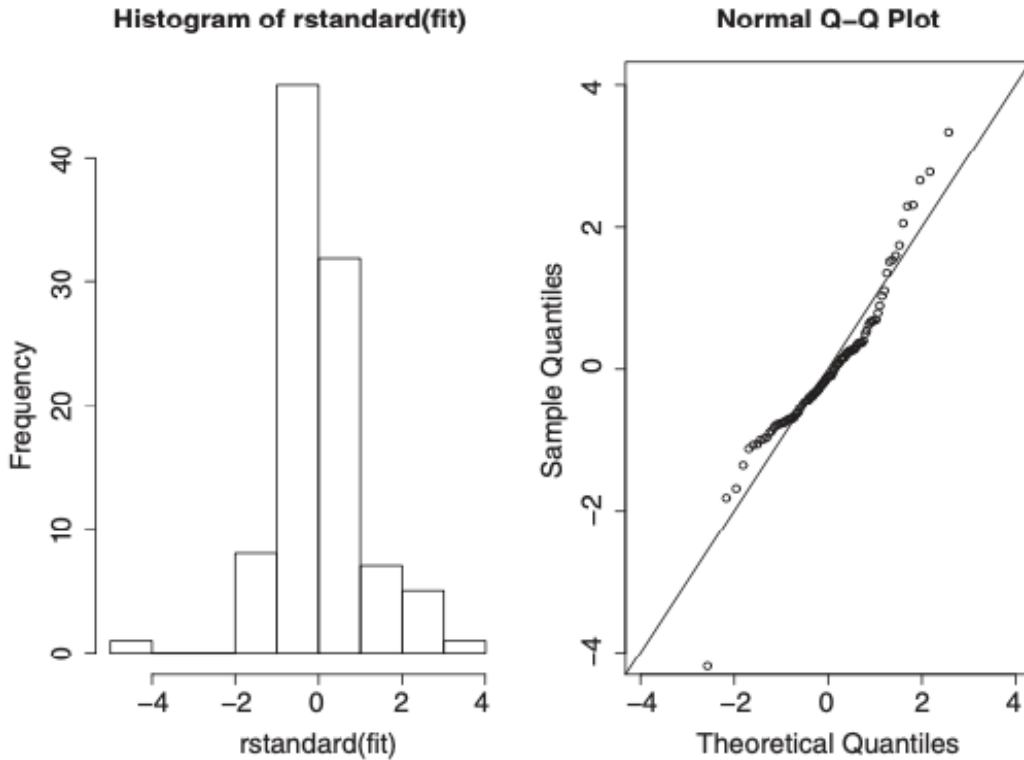


Figure 1: Histogram and normal QQ plot of standardized residuals.

1.4 Fixes

As mentioned in Section 1.2, non-normality is not necessarily a problem that needs to be fixed, except in small samples. In small samples, we can apply the bootstrap (Section 6.2.2) for robust standard error computation and a few different strategies (Section 6.3) for robust hypothesis testing.

2 Heteroskedastic errors

2.1 Origin

Suppose each observation y_i is actually the average of n_i underlying observations, each with variance σ^2 . Then, the variance of y_i is σ^2/n_i , which will differ across i if n_i differ. It is also common to see the variance of a distribution increase as the mean increases (as in Figure 2), whereas for a linear model the variance of y stays constant as the mean of y varies.

2.2 Consequences

All normal linear model inference from Unit 2 hinges on the assumption that $\epsilon_i \stackrel{\text{i.i.d.}}{\sim} N(0, \sigma^2)$. The coverage of confidence intervals and the levels of hypothesis tests may depart from their nominal levels. This is easiest to see if we consider the width of confidence intervals for $\mathbf{x}_0^T \boldsymbol{\beta}$; see Figure 2 for intuition.

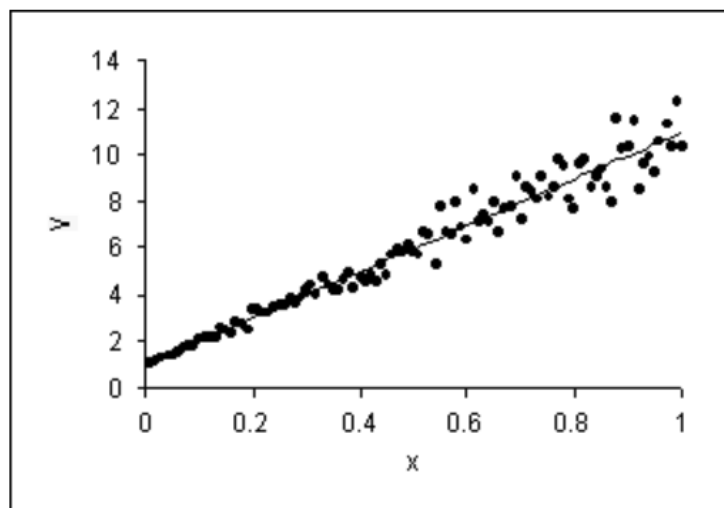


Figure 2: Heteroskedasticity in a simple bivariate linear model ([image source](#)).

2.3 Detection

Heteroskedasticity is usually assessed via the *residual plot* (Figure 3). In this plot, the standardized residuals r_i (3) are plotted against the fitted values $\hat{\mu}_i$. In the absence of heteroskedasticity, the spread of the points around the origin should be roughly constant as a function of $\hat{\mu}$ (Figure 3(a)). A common sign of heteroskedasticity is the fan shape where variance increases as a function of $\hat{\mu}$ (Figure 3(c)).

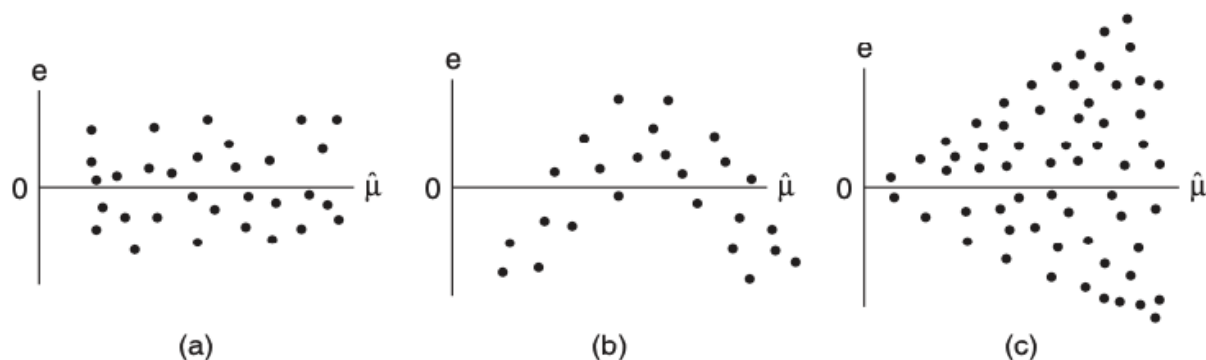


Figure 3: Residuals plotted against linear-model fitted values that reflect (a) model adequacy, (b) quadratic rather than linear relationship, and (c) nonconstant variance (image source: Agresti Figure 2.8).

2.4 Fixes

Heteroskedasticity-robust standard errors for hypothesis testing and confidence intervals can be obtained using a number of strategies, including the Huber-White sandwich estimator 6.2.1, the bootstrap 6.2.2, and permutation tests 6.3.1.

3 Correlated errors

3.1 Origin

Correlated errors can arise when observations have group, spatial, or temporal structure. Below are examples:

- Group/clustered structure: We have 10 samples (\mathbf{x}_{i*}, y_i) each from 100 schools.
- Spatial structure: We have 100 soil samples from a 10×10 grid on a $1\text{km} \times 1\text{km}$ field.
- Temporal structure: We have 366 COVID positivity rate measurements, one from each day of the year 2020.

The issue arises because there are common sources of variation among sample that are in the same group or spatially/temporally close to one another.

3.2 Consequences

Like with heteroskedastic errors, correlated errors can cause invalid standard errors. In particular, positively correlated errors typically cause standard errors to be smaller than they should be, leading to inflated Type-I error rates. For intuition, consider estimating the mean of a distribution based on n samples. Consider the cases when these samples are independent, compared to when they are perfectly correlated. The effective sample size in the former case is n and in the latter case is 1.

3.3 Detection

Residual plots once again come in handy to detect correlated errors. Instead of plotting the standardized residuals against the fitted values, we should plot the residuals against whatever variables we think might explain variation in the response that the regression does not account for. In the presence of group structures, we can plot residuals versus group (via a boxplot); in the presence of spatial or temporal structure, we can plot residuals as a function of space or time. If the residuals show a dependency on these variables, this suggests they are correlated.

3.4 Fixes

There are a few approaches to addressing correlated errors:

1. Estimate the covariance matrix Σ of the observations, so that $\mathbf{y} \sim N(\mathbf{X}\beta, \Sigma)$. This is a *generalized least squares* problem for which inference can be carried out. The generalized least squares estimate is $\hat{\beta} = (\mathbf{X}^T \Sigma^{-1} \mathbf{X})^{-1} \mathbf{X}^T \Sigma^{-1} \mathbf{y}$, which is distributed as $\hat{\beta} \sim N(\beta, (\mathbf{X}^T \Sigma^{-1} \mathbf{X})^{-1})$. We can carry out inference based on the latter distributional result analogously to how we did so in Unit 2. A special case of this is the *linear mixed effects model*, which hopefully we will have time to discuss in Unit 6.
2. Use the Liang-Zeger variance estimator; see Section 6.2.1.
3. Apply a clustered or block bootstrap; see Section 6.2.2.

4 Model bias

4.1 Origin

Model bias arises when predictors are left out of the regression model:

$$\text{assumed model: } \mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}; \quad \text{actual model: } \mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\boldsymbol{\gamma} + \boldsymbol{\epsilon}. \quad (4)$$

We may not always know about or measure all the variables that impact a response \mathbf{y} .

Model bias can also arise when the predictors do not impact the response on the linear scale. For example:

$$\text{assumed model: } \mathbb{E}[\mathbf{y}] = \mathbf{X}\boldsymbol{\beta}; \quad \text{actual model: } g(\mathbb{E}[\mathbf{y}]) = \mathbf{X}\boldsymbol{\beta}. \quad (5)$$

4.2 Consequences

In cases of model bias, the parameters $\boldsymbol{\beta}$ in the assumed linear model lose their meanings. The least squares estimate $\hat{\boldsymbol{\beta}}$ will be a biased estimate for the parameter we probably actually want to estimate. In the case (4) when predictors are left out of the regression model, these additional predictors \mathbf{Z} will act as confounders and create bias in $\hat{\boldsymbol{\beta}}$ as an estimate of the $\boldsymbol{\beta}$ parameters in the true model, unless $\mathbf{X}^T \mathbf{Z} = 0$. As discussed in Unit 2, this can lead to misleading conclusions.

4.3 Detection

Similarly to the detection of correlated errors, we can try to identify model bias by plotting the standardized residuals against predictors that may have been left out of the model. A good place to start is to plot standardized residuals against the predictors \mathbf{X} (one at a time) that are in the model, since nonlinear transformations of these might have been left out. In this case, you would see something like Figure 3(b).

It is possible to formally test for model bias in cases when we have repeated observations of the response for each value of the predictor vector. In particular, suppose that $\mathbf{x}_{i*} = \mathbf{x}_c$ for $c = c(i)$ and predictor vectors $\mathbf{x}_1, \dots, \mathbf{x}_C \in \mathbb{R}^p$. Then, consider testing the following hypothesis:

$$H_0 : y_i = \mathbf{x}_{i*}^T \boldsymbol{\beta} + \epsilon_i \quad \text{versus} \quad H_1 : y_i = \beta_{c(i)} + \epsilon_i. \quad (6)$$

The model under H_0 (the linear model) is nested in the model for H_1 (the saturated model), and we can test this hypothesis using an F -test called the *lack of fit F-test*.

4.4 Fixes

To fix model bias in the case (4), ideally we would identify the missing predictors \mathbf{Z} and add them to the regression model. This may not always be feasible or possible. To fix model bias in the case (5), it is sometimes advocated to find a transformation g (e.g. a square root or a logarithm) of \mathbf{y} such that $\mathbb{E}[g(\mathbf{y})] = \mathbf{X}\boldsymbol{\beta}$. However, a better solution is to use a *generalized linear model*, which we will discuss starting in Unit 4.

5 Outliers

5.1 Origin

Outliers often arise due to measurement or data entry errors. An observation can be an outlier in \mathbf{x} , in y , or both.

5.2 Consequences

An outlier can have the effect of biasing the estimate $\hat{\beta}$. This occurs when an observation has outlying \mathbf{x} as well as outlying y .

5.3 Detection

There are a few measures associated to an observation that can be used to detect outliers, though none are perfect. The first quantity is called the *leverage*, defined as

$$\text{leverage of observation } i \equiv \text{corr}(y_i, \hat{\mu}_i)^2. \quad (7)$$

This quantity measures the extent to which the fitted value $\hat{\mu}_i$ is sensitive to the (noise in the) observation y_i . It can be derived that

$$\text{leverage of observation } i = h_{ii}, \quad (8)$$

which is the i th diagonal element of the hat matrix \mathbf{H} . This is related to the fact that $\text{Var}[\hat{\epsilon}_i] = \sigma^2(1 - h_{ii})$. The larger the leverage, the smaller the variance of the residual, so the closer the line passes to the i th observation. The leverage of an observation is larger to the extent that \mathbf{x}_{i*} is far from $\bar{\mathbf{x}}$. For example, in the bivariate linear model $y_i = \beta_0 + \beta_1 x_i + \epsilon_i$,

$$h_{ii} = \frac{1}{n} + \frac{(x_i - \bar{x})^2}{\sum_{i'=1}^n (x_{i'} - \bar{x})^2}.$$

Note that the leverage is not a function of y_i , so a high-leverage point might or might not be an outlier in y_i and therefore might or might not have a strong impact on the regression. To assess more directly whether an observation is *influential*, we can compare the least squares fits with and without that observation. To this end, we define the *Cook's distance*

$$D_i = \frac{\sum_{i'=1}^n (\hat{\mu}_{i'} - \hat{\mu}_{i'}^i)^2}{p\hat{\sigma}^2}, \quad (9)$$

where $\hat{\mu}_{i'}^i = \mathbf{x}_{i'*}^T \hat{\beta}^i$ and $\hat{\beta}^i$ is the least squares estimate based on $(\mathbf{X}_{-i,*}, \mathbf{y}_{-i})$. An observation is considered influential if it has Cook's distance greater than one.

There is a connection between Cook's distance and leverage:

$$D_i = \left(\frac{y_i - \hat{\mu}_i}{\hat{\sigma}\sqrt{1 - h_{ii}}} \right)^2 \cdot \frac{h_{ii}}{p(1 - h_{ii})}. \quad (10)$$

We recognize the first term as the standardized residual; therefore a point is influential if its residual and leverage are large.

Note that Cook's distance may not successfully identify outliers. For example, if there are groups of outliers, then they will *mask* each other in the calculation of Cook's distance.

5.4 Fixes

If outliers can be detected, then the fix is to remove them from the regression. But, we need to be careful. Definitively determining whether observations are outliers can be tricky. Outlier detection can even be used as a way to commit fraud with data, as now-defunct blood testing start-up [Theranos is alleged to have done](#).

As an alternative to removing outliers, we can fit estimators $\hat{\beta}$ that are less sensitive to outliers; see Section 6.1.

6 Robust inference

There are a number of strategies designed to address one or more of the misspecification issues listed above. These fall into the categories of robust estimation (to get better estimates of $\hat{\beta}$ in the presence of outliers; see Section 6.1), robust standard error computation (to get more reliable standard errors in the presence of heteroskedasticity or correlated errors; see Section 6.2), and robust hypothesis testing (to get more reliable hypothesis tests in the presence of heteroskedasticity, correlated errors, and sometimes even model bias; see Section 6.3).

6.1 Robust estimation

The squared error loss $\sum_{i=1}^n (y_i - \mathbf{x}_{i*}^T \beta)^2$ is sensitive to outliers in the sense that a large value of $y_i - \mathbf{x}_{i*}^T \beta$ can have a significant impact on the loss function. The least squares estimate, as the minimizer of this loss function, is therefore sensitive to outliers. One way of addressing this challenge is to replace the squared error loss by a different loss that does not grow so quickly in $y_i - \mathbf{x}_{i*}^T \beta$. A popular choice for such a loss function is the Huber loss:

$$L_\delta(y_i - \mathbf{x}_{i*}^T \beta) = \begin{cases} \frac{1}{2}(y_i - \mathbf{x}_{i*}^T \beta)^2, & \text{if } |y_i - \mathbf{x}_{i*}^T \beta| \leq \delta; \\ \delta(|y_i - \mathbf{x}_{i*}^T \beta| - \delta), & \text{if } |y_i - \mathbf{x}_{i*}^T \beta| > \delta. \end{cases} \quad (11)$$

This function is differentiable, like the squared error loss, but grows linearly as opposed to quadratically. We can then define

$$\hat{\beta}^{\text{Huber}} \equiv \arg \min_{\beta} \sum_{i=1}^n L_\delta(y_i - \mathbf{x}_{i*}^T \beta).$$

This is an *M-estimator*; it is consistent and has an asymptotic normal distribution that can be used for inference.

6.2 Robust standard error computation

When the error terms in a regression are not homoskedastic and independent, the usual standard errors are invalid. There are several strategies to computing valid standard errors in such situations.

6.2.1 Huber-White and Liang-Zeger sandwich estimators

Let's say that $\mathbf{y} = \mathbf{X}\beta + \epsilon$, where $\epsilon \sim N(\mathbf{0}, \Sigma)$. Then, we can compute that the covariance matrix of the least squares estimate $\hat{\beta}$ is

$$\text{Var}[\hat{\beta}] = (\mathbf{X}^T \mathbf{X})^{-1} (\mathbf{X}^T \Sigma \mathbf{X}) (\mathbf{X}^T \mathbf{X})^{-1}. \quad (12)$$

Note that this expression reduces to the usual $\sigma^2 (\mathbf{X}^T \mathbf{X})^{-1}$ when $\Sigma = \sigma^2 \mathbf{I}$. It is called the sandwich variance because we have the $(\mathbf{X}^T \Sigma \mathbf{X})$ term sandwiched between two $(\mathbf{X}^T \mathbf{X})^{-1}$ terms. If we have some estimate $\hat{\Sigma}$ of the covariance matrix, we can construct

$$\widehat{\text{Var}}[\hat{\beta}] \equiv (\mathbf{X}^T \mathbf{X})^{-1} (\mathbf{X}^T \hat{\Sigma} \mathbf{X}) (\mathbf{X}^T \mathbf{X})^{-1}. \quad (13)$$

Different estimates $\hat{\Sigma}$ are appropriate in different situation. Below we consider two of the most common choices: one for heteroskedasticity (due to Huber-White) and one for correlated errors (due to Liang-Zeger).

Huber-White standard errors. Now, suppose $\Sigma = \text{diag}(\sigma_1^2, \dots, \sigma_n^2)$ for some variances $\sigma_1^2, \dots, \sigma_n^2 > 0$. The Huber-White sandwich estimator is defined by (12), with

$$\hat{\Sigma} \equiv \text{diag}(\hat{\sigma}_1^2, \dots, \hat{\sigma}_n^2), \quad \text{where} \quad \hat{\sigma}_i^2 = (y_i - \mathbf{x}_{i*}^T \hat{\beta})^2. \quad (14)$$

While each estimator $\hat{\sigma}_i^2$ is very poor, Huber and White's insight was that the resulting estimate of the (averaged) quantity $\mathbf{X}^T \hat{\Sigma} \mathbf{X}$ is not bad.

Liang-Zeger standard errors. Next, let's consider the case of correlated errors. Specifically, suppose that the observations are *clustered*, with correlated errors among clusters but not between clusters (recall Section 3.1). Suppose there are C clusters of observations, with the i th observation belonging to cluster $c(i) \in \{1, \dots, C\}$. Suppose for the sake of simplicity that the observations are ordered so that clusters are contiguous. Let $\hat{\epsilon}_c$ be the vector of residuals in cluster c , so that $\hat{\epsilon} = (\hat{\epsilon}_1, \dots, \hat{\epsilon}_C)$. Then, the true covariance matrix is $\Sigma = \text{block-diag}(\Sigma_1, \dots, \Sigma_C)$ for some positive definite $\Sigma_1, \dots, \Sigma_C$. The Liang-Zeger estimator is then defined by (12), with

$$\hat{\Sigma} \equiv \text{block-diag}(\hat{\Sigma}_1, \dots, \hat{\Sigma}_C), \quad \text{where} \quad \hat{\Sigma}_c \equiv \hat{\epsilon}_c \hat{\epsilon}_c^T. \quad (15)$$

Note that the Liang-Zeger estimator is a generalization of the Huber-White estimator. Its justification is similar as well: while each $\hat{\Sigma}_c$ is a poor estimator, the resulting estimate of the (averaged) quantity $\mathbf{X}^T \hat{\Sigma} \mathbf{X}$ is not bad as long as the number of clusters is large. Liang-Zeger standard errors are sometimes referred to as “clustered standard errors.”

6.2.2 Bootstrap

A completely different approach to constructing robust standard errors is the *bootstrap*. The core idea of the bootstrap is to use the data to construct an approximation to the data-generating distribution, and then to approximate the sampling distribution of any test statistic by simulating from this approximate data-generating distribution. This approach, pioneered by Brad Efron in 1979, replaces mathematical derivations with computation. The bootstrap is extremely flexible, and can be adapted to apply in a variety of settings.

Parametric bootstrap. The parametric bootstrap proceeds by fitting a parametric model, and then by resampling from this model. In the linear regression case, we use the original data to fit $(\hat{\beta}, \hat{\sigma}^2)$. Then, we sample new response vectors

$$y_i^b = \mathbf{x}_{i*}^T \hat{\beta} + \epsilon_i^b, \quad \epsilon_i^b \stackrel{\text{i.i.d.}}{\sim} N(0, \hat{\sigma}^2) \quad \text{for } b = 1, \dots, B. \quad (16)$$

We then fit a least squares coefficient vector $\hat{\beta}^b$ to $(\mathbf{X}, \mathbf{y}^b)$ for each b , and then get variance estimates by treating $\{\hat{\beta}^b\}_{b=1}^B$ as though it were the sampling distribution of $\hat{\beta}$. For example, we could use the sample standard deviation of $\hat{\beta}_j^b$ as the standard error for β_j .

This is the most model-based of the bootstrap variants. It assumes a completely well-specified model, and gives equivalent results to traditional parametric inference. It is typically not applied in regression settings, and presented here mainly for pedagogical purposes.

Residual bootstrap. We can weaken the assumptions of the parametric bootstrap by assuming only that $y_i = \mathbf{x}_{i*}^T \beta + \epsilon_i$, where $\epsilon_i \stackrel{\text{i.i.d.}}{\sim} F$ for some distribution F . Then, the data-generating distribution is specified by (β, F) , which we approximate by substituting $\hat{\beta}$ for β and the empirical

distribution of the residuals $\hat{\epsilon}_i$ (call it \hat{F}) for F . We can then sample new response vectors based on this approximate data-generating distribution:

$$y_i = \mathbf{x}_{i*}^T \hat{\boldsymbol{\beta}} + \epsilon_i^b, \quad \epsilon_i^b \stackrel{\text{i.i.d.}}{\sim} \hat{F} \quad \text{for } b = 1, \dots, B. \quad (17)$$

Note that i.i.d. sampling ϵ_i^b from \hat{F} amounts to sampling $(\epsilon_1^b, \dots, \epsilon_n^b)$ with replacement from $(\hat{\epsilon}_1, \dots, \hat{\epsilon}_n)$. Then, as with the parametric bootstrap, we fit a least squares coefficient vector $\hat{\boldsymbol{\beta}}^b$ to $(\mathbf{X}, \mathbf{y}^b)$ for each b and obtain standard errors by treating $\{\hat{\boldsymbol{\beta}}^b\}_{b=1}^B$ as though it were the sampling distribution of $\hat{\boldsymbol{\beta}}$.

The residual bootstrap corrects for non-normality, but not heteroskedasticity or correlated errors, since it assumes that the noise terms are i.i.d. from some distribution.

Pairs bootstrap. Weakening the assumptions further, let's assume only that $(\mathbf{x}_{i*}, y_i) \stackrel{\text{i.i.d.}}{\sim} F$ for some joint distribution F . We then resample our observations by sampling with replacement from the original observations.

Note that, unlike the parametric or residual bootstrap, the pairs bootstrap treats the predictors \mathbf{X} as random rather than fixed. The benefit of the pairs bootstrap is that it does not assume homoskedasticity, since the error variance is allowed to depend on \mathbf{x}_{i*} . Therefore, the pairs bootstrap addresses both non-normality and heteroskedasticity, though it does not address correlated errors (though variants of the pairs bootstrap do; see below). Note that the pairs bootstrap does not even assume that $\mathbb{E}[y_i] = \mathbf{x}_{i*}^T \boldsymbol{\beta}$ for some $\boldsymbol{\beta}$. However, in the presence of model bias, it is unclear for what parameters we are even doing inference. While the pairs bootstrap assumes less than the residual bootstrap, it may be somewhat less efficient in the case when the assumptions of the latter are met.

The pairs bootstrap has several variants that help it overcome correlated errors, in addition to heteroskedasticity. The *cluster bootstrap* is applicable in the case when errors have a clustered/grouped structure. In this case, we sample entire clusters of observations, with replacement, from the original set of clusters. The *moving blocks bootstrap* is applicable in the case of spatially or temporally structured errors. In this variant of the pairs bootstrap, we resample spatially or temporally adjacent blocks of observations together to preserve their joint correlation structure.

6.3 Robust hypothesis testing

6.3.1 Permutation tests

6.3.2 Rank-based tests

6.3.3 Bootstrap-based tests

7 R demo