

7.2: Estimating Linear Models

With stochastic models we don't know if the error assumptions are met, nor do we know the values of α and β ; therefore we must estimate them, as denoted by a hat (e.g., $\hat{\alpha}$ is the estimate for α). The stochastic model as shown in Equation (7.4) is estimated as:

$$Y_i = \alpha + \beta X_i + \epsilon_i \quad (7.4) \quad \hat{Y}_i = \hat{\alpha} + \hat{\beta} X_i + \epsilon_i$$

where ϵ_i is the **residual term** or the estimated error term. Since no line can perfectly pass through all the data points, we introduce a residual, ϵ , into the regression equation. Note that the predicted value of Y is denoted \hat{Y} (y-hat).

$$Y_i = \alpha + \beta X_i + \epsilon_i = \hat{Y}_i + \epsilon_i \quad \epsilon_i = Y_i - \hat{Y}_i = Y_i - \alpha - \beta X_i \quad \hat{Y}_i = \alpha + \beta X_i + \epsilon_i = Y_i + \epsilon_i = Y_i - Y_i = Y_i - \alpha - \beta X_i$$

7.2.1 Residuals

Residuals measure prediction errors of how far observation Y_i is from predicted \hat{Y}_i . This is shown in Figure 7.2.3.

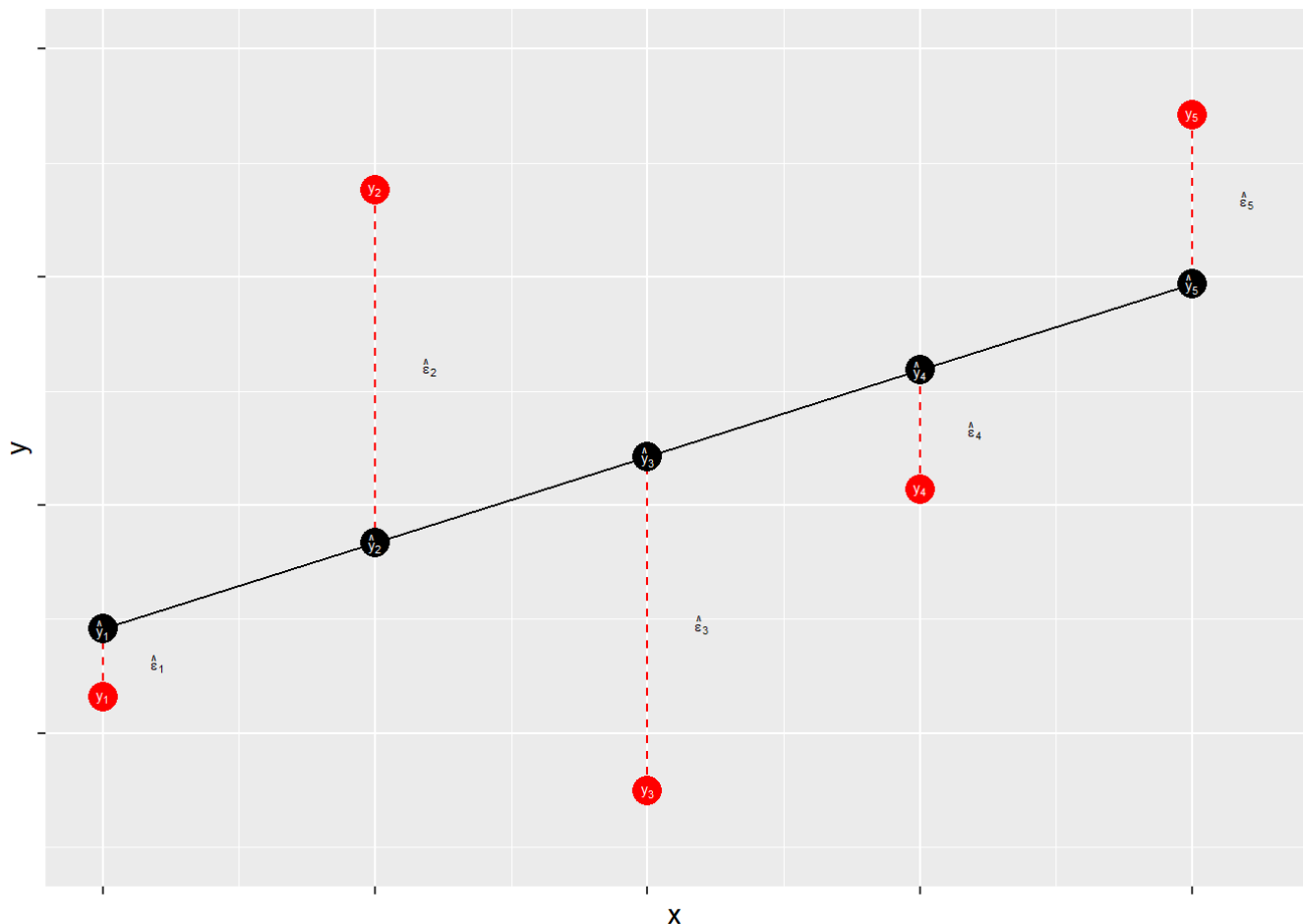


Figure 7.2.3: Residuals: Statistical Forensics

The residual term contains the accumulation (sum) of errors that can result from measurement issues, modeling problems, and irreducible randomness. Ideally, the residual term contains lots of small and independent influences that result in an overall random quality of the distribution of the errors. When that distribution is not random – that is, when the distribution of error has some systematic quality – the estimates of $\hat{\alpha}$ and $\hat{\beta}$ may be biased. Thus, when we evaluate our models we will focus on the shape of the distribution of our errors.

What's in ϵ ?

Measurement Error

- *Imperfect operationalizations*
- *Imperfect measure application*

Modeling Error

- *Modeling error/mis-specification*
- *Missing model explanation*
- *Incorrect assumptions about associations*
- *Incorrect assumptions about distributions*

Stochastic “noise”

- *Unpredictable variability in the dependent variable*

The goal of regression analysis is to minimize the error associated with the model estimates. As noted, the residual term is the estimated error, or overall miss" (e.g., $Y_i - \hat{Y}_i$). Specifically, the goal is to minimize the sum of the squared errors, $\sum_{i=1}^n e_i^2$. Therefore, we need to find the values of $\hat{\alpha}$ and $\hat{\beta}$ that minimize $\sum_{i=1}^n e_i^2$.

Note that for a fixed set of data $\{Y_i, X_i\}$, each possible choice of values for $\hat{\alpha}$ and $\hat{\beta}$ corresponds to a specific residual sum of squares, $\sum_{i=1}^n e_i^2$. This can be expressed by the following functional form:

$$S(\hat{\alpha}, \hat{\beta}) = \sum_{i=1}^n e_i^2 = \sum_{i=1}^n (Y_i - \hat{Y}_i)^2 = \sum_{i=1}^n (Y_i - \hat{\alpha} - \hat{\beta}X_i)^2 \quad (7.5)$$

Minimizing this function requires specifying estimators for $\hat{\alpha}$ and $\hat{\beta}$ such that $S(\hat{\alpha}, \hat{\beta}) = \sum_{i=1}^n e_i^2$ is at the lowest possible value. Finding this minimum value requires the use of calculus, which will be discussed in the next chapter. Before that, we walk through a quick example of simple regression

This page titled [7.2: Estimating Linear Models](#) is shared under a [CC BY 4.0](#) license and was authored, remixed, and/or curated by [Jenkins-Smith et al. \(University of Oklahoma Libraries\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.