

Introduction to Least Squares Regression

Chapter 6, Lab 2

OpenIntro Biostatistics

Topics

- Calculating a least squares line
- Checking assumptions with residual plots

The previous lab introduced the mechanics of interpreting a line of best fit and the assumptions for linear regression. This lab formally introduces the statistical model for least squares regression and discusses the residual plots used to assess the assumptions for linear regression.

The material in this lab corresponds to Sections 6.2 and 6.3.1 of *OpenIntro Biostatistics*.

Introduction

Least squares regression

The vertical distance between a point in the scatterplot and the predicted value on the regression line is the **residual** for the point. For an observation (x_i, y_i) , where \hat{y}_i is the predicted value according to the line $\hat{y} = b_0 + b_1x$, the residual is the value $e_i = y_i - \hat{y}_i$.

The **least squares regression line** is the line which minimizes the sum of the squared residuals for all the points in the plot; i.e., the regression line is the line that minimizes $e_1^2 + e_2^2 + \dots + e_n^2$ for the n pairs of points in the dataset.

For a general population of ordered pairs (x, y) , the population regression model is

$$y = \beta_0 + \beta_1x + \epsilon,$$

where ϵ is a normally distributed ‘error term’ with mean 0 and standard deviation σ .

The terms β_0 and β_1 are parameters with estimates b_0 and b_1 . These estimates can be calculated from summary statistics: the sample means of x and y (\bar{x} and \bar{y}), the sample standard deviations of x and y (s_x, s_y), and the correlation between x and y (r).

$$b_1 = r \frac{s_y}{s_x} \quad b_0 = \bar{y} - b_1\bar{x}$$

Plots for checking assumptions

There are a variety of **residual plots** used to check the fit of a least squares line. The ones used in this textbook are scatterplots in which predicted values are on the x -axis and residual values on the y -axis. Residual plots are useful for checking the assumptions of linearity and constant variability.

To assess the normality of residuals, **normal probability plots** are used. These plots are also known as quantile-quantile plots, or Q-Q plots.

Calculating a least squares line

1. Run the following code chunk to load the prevend data and store a random sample of 500 individuals as prevend.sample.

```
#load the data
library(oibiostat)
data("prevend")

#take a random sample
set.seed(5011)
sample.size = 500
prevend.sample = prevend[sample(1:nrow(prevend), sample.size, replace = FALSE), ]
```

2. From the data in prevend.sample, calculate the least squares regression line for the relationship between RFFT score (RFFT) and age in years (Age).
 - a) Calculate b_1 and b_0 from summary statistics, using the formulas on the previous page.
 - b) Verify your answer to part a) using the `lm()` function.
 - c) Write the equation of the least-squares line.

Checking assumptions with residual plots

There are four assumptions that must be met for a linear model to be considered reasonable: linearity, constant variability, independent observations, and normally distributed residuals.

Even though linearity and constant variability can be assessed from the scatterplot of y versus x , it is helpful to consult residual plots for a clearer view. Normality of residuals is best assessed through a normal probability plot; although skew can be visible from a histogram of the residuals, deviations from normality are more obvious on a normal probability plot.

RFFT and age in the prevend data

3. Run the following chunk to create a residual plot where the residual values are plotted on the y -axis against predicted values from the model on the x -axis, using data in prevend.samp.

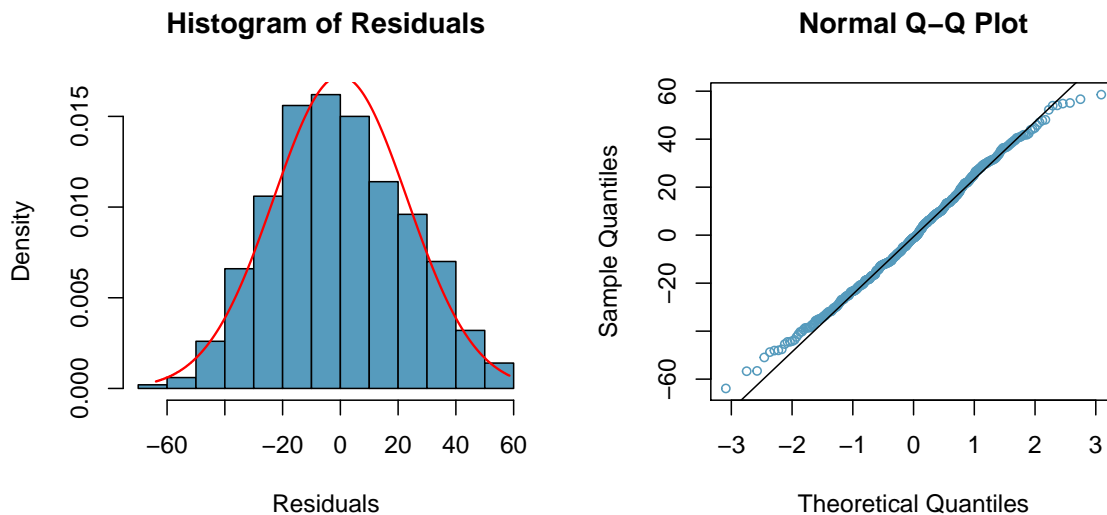
```
#store the residuals from the linear model
prevend.residuals = resid(lm(RFFT ~ Age, data = prevend.sample))

#store the predicted RFFT scores from the linear model
prevend.predicted = predict(lm(RFFT ~ Age, data = prevend.sample))

#create residual plot
plot(prevend.residuals ~ prevend.predicted)
abline(h = 0, col = "red", lty = 2)
```

- a) When a linear model is a good fit for the data, the residuals should scatter around the horizontal line $y = 0$ with no apparent pattern. Does a linear model seem appropriate for these data?
 - b) Does the variability of the residuals seem constant across the range of predicted RFFT scores?
4. Run the code chunk shown in the template to create a normal probability plot of the residuals. For comparison purposes, the following figure shows a histogram of the residual values overlaid with a normal curve and the normal probability plot.

Do the residuals appear to be normally distributed?



5. Overall, does it seem that a least squares regression line is an appropriate model for estimating the relationship between cognitive function (as measured by RFFT score) and age? Recall that as discussed in the previous lab, it is reasonable to assume the independence assumption holds for these data.

Clutch volume and body size in the frog data

The frog dataset in the `oibistat` package contains observations from a study conducted on a frog species endemic to the Tibetan Plateau. Researchers collected measurements on egg clutches and female frogs found at breeding ponds across five study sites.

Previous research suggests that larger body size allows females to produce egg clutches with larger volumes. Frog embryos are surrounded by a gelatinous matrix that may protect developing embryos from temperature fluctuation or ultraviolet radiation; a larger matrix volume provides added protection. In the data, clutch volume (`clutch.volume`) is recorded in cubic millimeters and female body size (`body.size`) is measured as length in centimeters.

The following questions step through examining whether a linear regression model is appropriate for the relationship between female body size and clutch volume.

6. Create a scatterplot of clutch volume versus female body size and plot the least squares line.

7. Create a residual plot where the residual values are plotted on the y -axis against predicted values from the model on the x -axis.
 - a) Does the linearity assumption seem to be satisfied?
 - b) Is the variability of the residuals constant across the range of predicted clutch volumes?
8. Assess whether it can be reasonably assumed that the observations are independent.
9. Create a Q-Q plot and assess whether the residuals appear to be normally distributed.
10. Evaluate whether a least squares regression line is an appropriate model for estimating the relationship between female body size and clutch volume.