# Introduction to Least Squares Regression

*Chapter 6, Lab 2: Solutions*

*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 $\hat{y} = b_0 + b_1 x$, 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 $e_1^2 + e_2^2 + ... + e_n^2$ for the $n$ pairs of points in the dataset.

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

$$y = \beta_0 + \beta_1 x + \epsilon,$$

where $\epsilon$ is a normally distributed 'error term' with mean 0 and standard deviation $\sigma$.

The terms $\beta_0$ and $\beta_1$ are parameters with estimates $b_0$ and $b_1$. These estimates can be calculated from summary statistics: the sample means of $x$ and $y$ ($\overline{x}$ and $\overline{y}$), the sample standard deviations of $x$ and $y$ ($s_x, s_y$), and the correlation between $x$ and $y$ ($r$).

$$b_1 = r \frac{s_y}{s_x} \qquad b_0 = \overline{y} - b_1 \overline{x}$$

*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.samp` dataset.

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

2. From the data in `prevend.samp`, 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.

    The slope ($b_1$) is -1.26 and the intercept ($b_0$) is 137.55.

```
#use r as a calculator

#define constants
x.bar = mean(prevend.samp$Age)
y.bar = mean(prevend.samp$RFFT)
s.x = sd(prevend.samp$Age)
s.y = sd(prevend.samp$RFFT)
r = cor(prevend.samp$Age, prevend.samp$RFFT)

#calculate b_1
b_1 = r * (s.y/s.x)
b_1
```

```
## [1] -1.261359
```

```
#calculate b_0
b_0 = y.bar - (b_1 * x.bar)
b_0
```

```
## [1] 137.5497
```

    b) Verify your answer to part a) using the `lm( )` function.

    The results from the `lm( )` function equal the answer from part a).

```
#use lm( )
lm(prevend.samp$RFFT ~ prevend.samp$Age)
```

```
##
## Call:
## lm(formula = prevend.samp$RFFT ~ prevend.samp$Age)
##
## Coefficients:
##      (Intercept)   prevend.samp$Age
##          137.550            -1.261
```

    c) Write the equation of the least-squares line.

The equation of the least-squares line can be written as either $\hat{y} = 137.55 - 1.26x$ or $\widehat{RFFT} = 137.55 - 1.26(age)$.

**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.
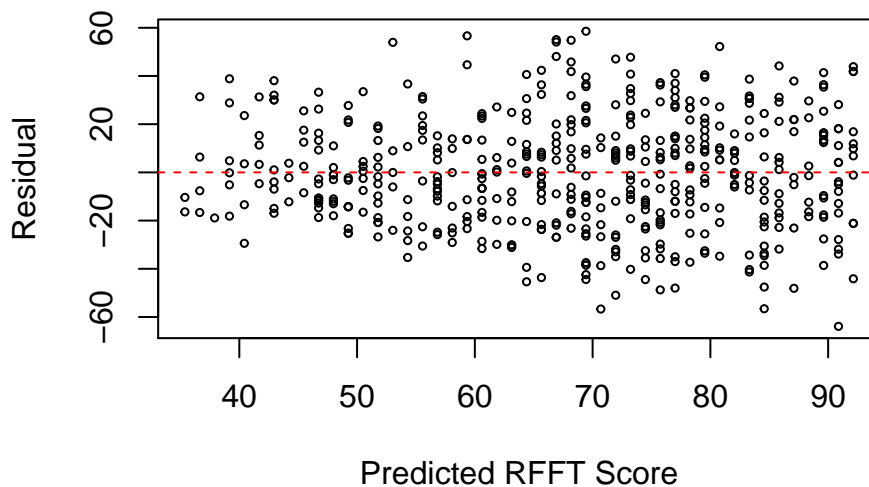
*RFFT and age in the* prevend *data*

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 = residuals(lm(RFFT ~ Age, data = prevend.samp))

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

#create residual plot
plot(prevend.residuals ~ prevend.predicted,
 cex = 0.5,
 main = "Residual Plot for RFFT versus Age (n = 500)",
 xlab = "Predicted RFFT Score",
 ylab = "Residual")
abline(h = 0, col = "red", lty = 2)
```

3

## Residual Plot for RFFT versus Age (n = 500)



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?

Yes, a linear model seems appropriate; the residuals scatter about the horizontal line $y = 0$ with no apparent pattern. There is a roughly equal number of points above the line as below it, which indicates that the line goes through the center of the 'cloud' of data points.
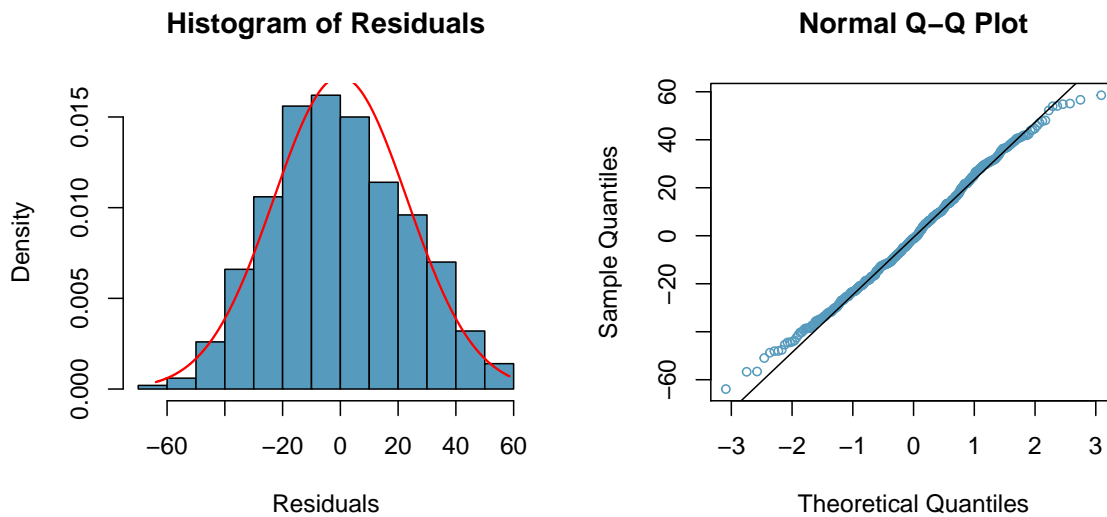
b) Does the variability of the residuals seem constant across the range of predicted RFFT scores?

The variability are generally constant across predicted RFFT scores between 60 to 90 points, but there seems to be less variability for lower predicted 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?

Yes, the residuals follow a straight line on the Q-Q plot, with only slight deviations from normality in the tails. This is also visible in the histogram of the residuals with an overlaid normal curve.

| Histogram of Residuals | Normal Q–Q Plot |
|---|---|



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.

Overall, a least squares regression line seems appropriate. There is some nonconstant variability for lower predicted RFFT scores that suggest caution should be exercised when using the model to conduct inference about older individuals (since older individuals are predicted to have lower scores).[1] At this point, it is more important to note that data almost never perfectly satisfy model assumptions; the learning goal here is to understand the mechanics of assumptions checking and to be able to identify especially severe violations.

*Clutch volume and body size in the* `frog` *data*

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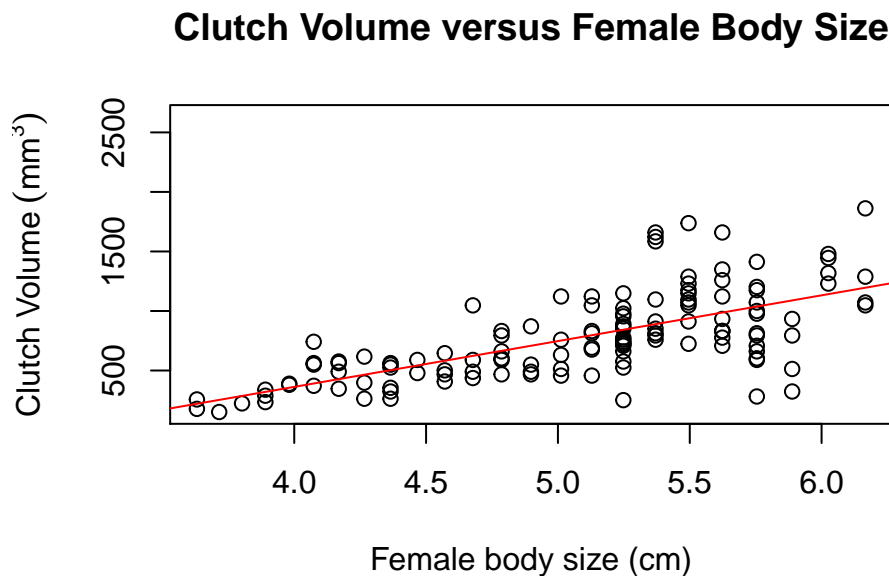
The following questions step through examining whether a linear regression model is appropriate for the relationship between female body size and clutch volume.

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

```
#load the data
data("frog")
```

---

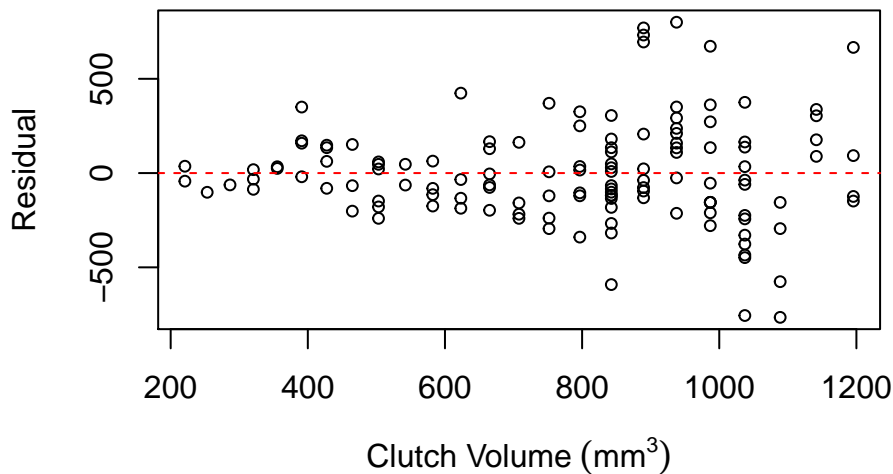[1] Inference for regression will be discussed in a later lab.

```
#create a scatterplot and add a regression line
plot(frog$clutch.volume ~ frog$body.size,
 main = "Clutch Volume versus Female Body Size",
 xlab = "Female body size (cm)", ylab = expression("Clutch Volume" ~ (mm^3)))
abline(lm(frog$clutch.volume ~ frog$body.size), col = "red")
```



**Clutch Volume versus Female Body Size**

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.

```
#create residual plot
frog.model = lm(clutch.volume ~ body.size, data = frog)
plot(resid(frog.model) ~ predict(frog.model),
 main = "Residual Plot for Clutch Volume vs. Body Size",
 ylab = "Residual",
 xlab = expression("Clutch Volume" ~ (mm^3)),
 cex = 0.75)
abline(h = 0, col = "red", lty = 2)
```

## Residual Plot for Clutch Volume vs. Body Size



a) Does the linearity assumption seem to be satisfied?

Linearity appears to be satisfied; there is no apparent pattern in the residuals and there is a roughly equal number of points above and below the $y = 0$ line.

b) Is the variability of the residuals constant across the range of predicted clutch volumes?

No; the residual plot makes it easier to see that the residuals are more variable for larger values of predicted clutch volume than smaller values.
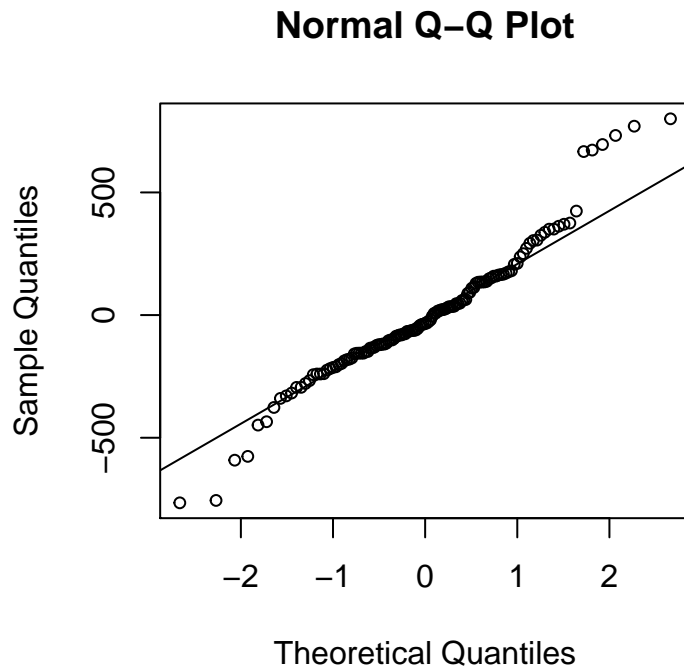
8. Assess whether it can be reasonably assumed that the observations are independent.

The data for body size and clutch volume are from frogs (and egg clutches) found at breeding ponds at five different study sites. In the absence of more specific biological information about the life cycle of this frog and its breeding habits, it seems reasonable to assume that the values of one pair (clutch volume and body size) provide no information about another pair. It would not be reasonable to assume independence if, for example, frogs found at a specific study site are more likely to be related to each other and that related frogs tend to have similar body sizes and/or clutch volumes.

9. Create a Q-Q plot and assess whether the residuals appear to be normally distributed.

The majority of the residuals closely follow a normal distribution, but there exists some deviation from normality in both tails.

```
#create q-q plot
qqnorm(resid(frog.model), cex = 0.75)
qqline(resid(frog.model))
```

## Normal Q–Q Plot



10. Evaluate whether a least squares regression line is an appropriate model for estimating the relationship between female body size and clutch volume.

The least squares line seems to be a reasonable fit for the data such that predictions about average clutch volume can be made based on female body size. However, since the constant variability assumption is violated, caution is required when using the model to conduct inference; this concern will be discussed further once inference for regression is covered.

The main point of showing the model with the `frog` data is to highlight the fact that data tend to be messy—the previous regression using the PREVEND data actually represent a case where data fit model assumptions surprisingly well. When assumptions are violated, the question of how to proceed is often a subtle one, requiring topics that would be covered in a specialized course on regression.