# Simple Linear Regression

*Chapter 6, Lab 1*

*OpenIntro Biostatistics*

**Topics**

- Examine visual relationships between numerical variables
- Estimate a regression line using least squares
- Interpret a least squares regression line
- Evaluate model fit
- Perform statistical inference with regression
- Prediction with the least squares regression line
- Check the necessary conditions for regression

This lab introduces simple linear regression.

The material in this lab corresponds to Sections 6.1 - 6.4 of *OpenIntro Biostatistics*.

## Background

Triglycerides levels provide strong indication of heart health. Triglycerides are a type of lipid stored in your blood derived from unused calories. When your body needs energy between meals hormones release triglycerides. In some cases, people who eat excess high carbohydrate foods (more than the body requires) have high triglyceride levels. This condition is known as hypertriglyceridemia.

A lipid panel or blood test for cholesterol can easily determine a person's triglyceride level. The following thresholds show generally accepted standards for certain triglyceride levels (mg/dL).

- Normal: < 150 mg/dL
- Moderate: 150 to 199 mg/dL
- High: 200 to 499 mg/dL
- Very High: > 500 mg/dL

A patient with high and very high triglyceride levels are at extreme risk for arteriosclerosis (hardening of the arteries or thickening of the artery walls). This condition increases the risk of stroke and heart and vascular failure. In addtion to heart disease, high triglycerides contribute to obesity, pancreatitis, hypothyroidism, and diabetes.

Prevention or at least mitigation of hypertriglyceridemia can be handled with multiple methods. Maintaining a health lifestyle, including regular exercise, proper diet, and limited alcohol consumption helps ensure triglyceride levels are managed appropriately. In some cases living a healthy lifestyle is not enough and medication, such as Lipitor or Crestor (statins) must also be taken to manage triglyceride levels.

The `tri` dataset provides a list of de-identified patients who visited their primary care physician in the past month. For each patient, data elements including their triglyceride level, age, waist circumference measurement, gender, and patient disease indicators were recorded.

The `tri` dataset can be found in the `oibiostat` package. For more information, use the `help(tri)` command.

```
library(oibiostat)
data("tri")
```

## Exploratory Data Analysis

1. Explore the `tri` dataset. What is the dimensions, data types, etc.?

2. Visually inspect the relationship between a patient's waist circumference measurement, `WAIST_CIR` and their associated triglyceride level, `TEST_RESULT_NUMERIC`. What type of relationship exists?

3. What is the strength of the linear relationship between a patient's waist circumference measurement, `WAIST_CIR` and their associated triglyceride level, `TEST_RESULT_NUMERIC`. Does this indicate larger waist circumference measurements cause higher triglyceride levels? Explain.

4. Explore the relationships between the remaining variables, `AGE`, `SEX_NAME`, and other disease indicators in relationship to the target variable, `TEST_RESULT_NUMERIC`. Is a scatterplot appropriate to evaluate the relationship between `SEX_NAME` and `TEST_RESULT_NUMERIC`? If not, use the appropriate visual tool to assess if there is any difference in triglyceride levels across genders.

5. Based on your exploratory data analysis, which of the variables examined best explains the variability in `TEST_RESULT_NUMERIC`?

## Least Squares Estimation

In this section we will fit two regression lines using the least squares approach. First, we are interested in examining the relationship between a patient's waist circumference measurement, `WAIST_CIR` and their associated triglyceride level, `TEST_RESULT_NUMERIC`. Secondly, we are interested in examining the relationship between a patient's gender, `SEX_NAME` and their associated triglyceride level, `TEST_RESULT_NUMERIC`.

Determine the best fitting linear regression equations using the following procedures.

### Triglyceride Level ~ Waist Circumference

1. Calculate the slope, $b_0$ and y-intercept, $b_1$ using the following expressions.

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

2. Now, use R to calculate the coefficients of the least squares line using the `lm()` function. Compare your results.

**Triglyceride Level ~ Gender**

Refer to *Question 4* in the Exploratory Data Analysis section above, specifically the summary statistics and visualizations created to infer the relationship between triglyceride levels, `TEST_RESULT_NUMERIC` and gender, `SEX_NAME`.

3. Does the observed means for triglyceride levels by gender appear to be widely different? Calculate the two sample means of triglyceride level for each gender (if not determined earlier). Calculate the difference in both sample means.

4. Now, build a linear regression model using R's `lm()` function to explain triglyceride levels, `TEST_RESULT_NUMERIC` using gender, `SEX_NAME` as your *only* explanatory variable.

5. Comment on the relationship between the estimated coefficients for $\beta_0$ and $\beta_1$ and the sample means of triglyceride levels calculated for each gender in Question 1. What does the difference in sample means of triglyceride levels for each gender equal in terms of estimated coefficients?

## Interval Estimation

Just as we constructed confidence intervals for parameters such as $\mu$, we can also calculate confidence intervals for regression parameters, $\beta$. We will discuss two approaches to calculating confidence intervals for regression parameters in the model, *Triglyceride Level ~ Waist Circumference*.

1. Calculate the 95% confidence interval for the slope, $\beta_1$ using the following expression. Recall the expression for mean squared error, $MSE = \frac{1}{n-2} \sum (y_i - \hat{y}_i)^2$

$$b_1 \pm t_{\frac{\alpha}{2}, n-2} \sqrt{\frac{MSE}{\sum (x_i - \bar{x})^2}}$$

2. Now, use R to calculate the 95% confidence interval for the slope, $\beta_1$ using `confint()` function. Compare your results.

## Interpretation

**Triglyceride Level ~ Waist Circumference**

1. Interpret $b_1$, $b_0$, and the 95% confidence interval for $\beta_1$ in the context of the research question. Does, $b_0$ have any contextual meaning?

**Triglyceride Level ~ Gender**

2. Interpret $b_1$ and $b_0$ in the context of the research question. Does, $b_0$ have any contextual meaning?

## Model Fit

In this section we will evaluate the fit of both models estimated previously.

### Triglyceride Level ~ Waist Circumference

1. Calculate the strength of the association, $r$ between a patient's waist circumference measurement, `WAIST_CIR` and their triglyceride level, `TEST_RESULT_NUMERIC`.

2. What is the proportion of variability in the response, triglyceride levels, explained by the model? Calculate using the following expression. Does the model, including only the explanatory variable, waist circumference, do a fair job at explaining the variability in triglyceride levels according to $R^2$?

$$R^2 = \frac{s_y^2 - s_{residuals}^2}{s_y^2}$$

### Triglyceride Level ~ Gender

3. What is the proportion of variability in the response, triglyceride levels, explained by the model? Calculate using the expression offered in *Question 2*. Does the model, including only the explanatory variable, gender, do a fair job at explaining the variability in triglyceride levels according to $R^2$?

## Statistical Inference

Since the least squares method offers a way to estimate population parameters, we are essentially performing statistical inference, where the slope of the least squares line, $b_1$ is an estimate of the population slope, $\beta_1$ and the intercept of the least squares line, $b_0$ is an estimate of the population intercept, $\beta_0$.

In regresssion the typical hypothesis for $\beta_1$ is as follows:

$$H_0 : \beta_1 = 0$$

$$H_A : \beta_1 \neq 0$$

where under the null hypothesis, $H_0$ the explanatory and response variables are not associated and under the alternative hypothesis the explanatory and response variables are associated. The test statistic, $t - statistic$ is provided below:

$$t = \frac{b_1}{s.e.(b_1)}$$

where the $t-statistic$ follows a $t$ distribution with $df = n - 2$.

### Triglyceride Level ~ Waist Circumference

1. Setup the appropriate hypothesis test for $\beta_{Waist-Circumference}$ and provide the test statistic and $p-value$ for the testing outcome. Rule on the test and provide your interpretation in the context of the research hypothesis.

2. Use the R function, `summary(model)` to confirm your $t-statistic$ and $p-value$ for $\beta_{Waist-Circumference}$.

### Triglyceride Level ~ Gender

3. Use the R function, `summary(model)` to determine the $t-statistic$ and $p-value$. Also, use the R function `t.test()` to evaluate the difference between triglyceride levels according to gender. Do you see any similiarities between the test statistics in the regression output and $t$ test output?

## Prediction

In addition to performing inference about the association between the response and explanatory variables, we can use the least squares line to predict responses according to values of the explanatory variable that fall within the scope of the model. For example one might like to know the predicted triglyceride level for a patient with a waist circumference, $x_{Waist-Circumference}$ or the predicted triglyceride level for a patient that is $x_{Male}$.

### Triglyceride Level ~ Waist Circumference

1. Use your estimated simple linear regression model for *Triglyceride Level ~ Waist Circumference* to predict the triglyceride level of a patient whos waist circumference measures 32 inches.

2. Now use the R function `predict()` to confirm your answer in part 1 and provide a *prediction* interval.

### Triglyceride Level ~ Gender

3. Use your estimated simple linear regression model *Triglyceride Level ~ Gender* to predict the triglyceride level of a patient whos is 'Female'.

4. Now use the R function `predict()` to confirm your answer in part 3 and provide a *prediction* interval.

## Checking Model Assumptions for Regression

In order to use linear regression to better understand the relationship between two variables, the following conditions need to validated. You will need to verify each is met or call out concerns based on your observations of the residual plots, normal probability plots, and histogram of your residuals. In this section we will only use the *Triglyceride Level ~ Waist Circumference* simple linear regression model.

- **Linearity**: data show a linear trend in change for response $y$ as a function of predictor $x$.
- **Constant variability**: the variability of the response variable about the line remains roughly constant as the predictor variable changes.
- **Independent observations**: the $(x, y)$ pairs are independent; i.e., values of one pair provide no information about other pairs.
- **Approximate normality of residuals**: $e_i = y_i - \hat{y}_i$ where $e_i$ are normally distributed under $N(0, \sigma)$.

1. Create a residual plot where the predicted values, $\hat{y}$ are on the *x-axis* and the residuals are on the *y-axis*. Use the following R function to add a horizontal line to your residual plot, *abline(h = 0, col = "red")*. What do you observe regarding conditions 1 and 2?

2. Create a normal probability plot where the theoretical quantiles for a normal distribution are on the *x-axis* and the observed quantiles from the data are on the *y-axis*, and a histogram of the residuals. Observe the tails of the distribution in the normal probability to monitor for departures from normality. Is the histogram of the residuals nearly normal? Does them confirm or deny assumption 4?

3. Use a common sense approach paired with your understanding of the study to validate assumption 3. Explain your reasoning.

4. Does your investigation of the conditions for linear regression provide pause for concern? If so, what might you suggest to correct such issues and continue using linear regression as your choice for model building? Hint: *Transformations, weighted regression, leverage values/outlier identification, etc.* Explain.