# Unit 31
# A Hypothesis Test about Correlation and Slope in a Simple Linear Regression

Objectives:
- To perform a hypothesis test concerning the slope of a least squares line
- To recognize that testing for a statistically significant slope in a linear regression and testing for a statistically significant linear relationship (i.e., correlation) are actually the same test

Recall that we can use the Pearson correlation $r$ and the least squares line to describe a linear relationship between two quantitative variables $X$ and $Y$. We called a sample of observations on two such quantitative variables bivariate data, represented as $(x_1, y_1)$, $(x_2, y_2)$, ... , $(x_n, y_n)$. Typically, we let $Y$ represent the response variable and let $X$ represent the explanatory variable. From such a data set, we have previously seen how to obtain the Pearson correlation $r$ and the least squares line. We now want a hypothesis test to decide whether a linear relationship is statistically significant.

For example, return to Table 10-2 where data is recorded on $X$ = "the dosage of a certain drug (in grams)" and $Y$ = "the reaction time to a particular stimulus (in seconds)" for each subject in a study. We have reproduced this data here as

| Table 31-1 |
| --- |
| **Dosage and Reaction Time Data (from Table 10-2)** |
| The dosage of a stimulant drug and the reaction time to a stimulus are recorded for each of several subjects injected with the drug. |

| Dosage (grams) | 4 | 4 | 6 | 6 | 8 | 8 | 10 | 10 |
| --- | --- | --- | --- | --- | --- | --- | --- | --- |
| Reaction Time (seconds) | 7.5 | 6.8 | 4.0 | 4.4 | 3.9 | 3.1 | 1.4 | 1.7 |

Table 31-1. From the calculations in Table 10-3, you previously found the correlation between dosage and reaction time to be $r = -0.9628$, and from these same calculations you also found the least squares line to be $rct = 10.225 - 0.875(dsg)$, where $rct$ represents reaction time in seconds, and $dsg$ represents dosage in grams. However, these are just descriptive statistics; alone, they can not tell us whether there is sufficient evidence of a linear relationship between the two quantitative variables.

When we study the prediction of a quantitative response variable $Y$ from a quantitative explanatory variable $X$ based on the equation for a straight line, we say we are studying the *linear regression to predict Y from X* or the *linear regression of Y on X*. Deciding whether or not a linear relationship is significant is the same as deciding whether or not the slope in the linear regression to predict $Y$ from $X$ is different from zero. If the slope is zero, then $Y$ does not tend to change as $X$ changes; if the slope is not zero, then $Y$ will tend to change as $X$ changes.

We need a hypothesis test to decide if there is sufficient evidence that the slope in a simple linear regression is different from zero. Such a test could be one-sided or two-sided. If our focus was on finding evidence that the slope is greater than zero (i.e., finding evidence that the linear relationship was a positive one), then we would use a one-sided test. If our focus was on finding evidence that the slope is less than zero (i.e., finding evidence that the linear relationship was a negative one), then we would use a one-sided test. However, if our focus was on finding evidence that the slope is different from zero in either direction (i.e., finding evidence that the linear relationship was either positive or negative), then we would use a two-sided test.

With a simple random sample of bivariate data of size $n$, the following $t$ statistic with $n - 2$ degrees of freedom is available for the desired hypothesis test:

$$t_{n-2} = \frac{b - 0}{s_{Y|X} \Big/ \sqrt{\sum (x - \bar{x})^2}} \ ,$$

234

where $s_{Y|X}$ is called the *standard error of estimate*. The formula for this $t$ statistic has a format similar to other $t$ statistics we have encountered. The numerator is the difference between a sample estimate of the slope ($b$) and the hypothesized value of the slope (zero), and the denominator is an appropriate standard error. The calculation of this test statistic requires obtaining the slope $b$ of the least squares line, the sum of the squared deviations of the sample $x$ values from their mean $\sum (x - \bar{x})^2$, and the standard error of estimate $s_{Y|X}$. We have previously seen how to obtain $b$ and $\sum (x - \bar{x})^2$, but we have never previously encountered the standard error of estimate $s_{Y|X}$.

You can think of the standard error of estimate $s_{Y|X}$ for bivariate data similar to the way you think of the standard deviation $s$ for one sample of observations of a quantitative variable, which can be called *univariate data*. Recall that the standard deviation for a sample $x$ values is obtained by dividing $\sum (x - \bar{x})^2$ by one less than the sample size, and taking the square root of the result. The standard error of estimate $s_{Y|X}$ is a measure of the dispersion of data points around the least squares line similar to the way that the standard deviation $s$ is a measure of the dispersion of observations (around the sample mean) in a single sample. The standard error of estimate $s_{Y|X}$ is obtained by dividing the sum of the squared residuals by two less than the sample size, and taking the square root of the result. If you wish, you may do this calculation for the data of Table 31-1, after which you can calculate the test statistic $t_{n-2}$. In general, a substantial amount of calculation is needed to obtain this test statistic, and the use of appropriate statistical software or a programmable calculator is recommended. Consequently, shall not concern ourselves with these calculations in detail.

To illustrate the $t$ test about the slope in a simple linear regression, let us consider using a 0.05 significance level to perform a hypothesis test to see if the data of Table 31-1 (Table 10-2) provide evidence that the linear relationship between drug dosage and reaction time is significant, or in other words, evidence that the slope in the linear regression to predict reaction time from drug dosage is different from zero.

The first step is to state the null and alternative hypotheses, and choose a significance level. This test is two-sided, since we are looking for a relationship which could be positive or negative (i.e., a slope which could be different from zero in either direction). We can complete the first step of the hypothesis test as follows:

$H_0$: slope $= 0$ vs. $H_1$: slope $\neq 0$ ($\alpha = 0.05$, two-sided test) .

We could alternatively choose to write the hypotheses as follows:

$H_0$: The linear relationship between drug dosage and reaction time is not significant
  vs.                                                                ($\alpha = 0.05$)
$H_1$: The linear relationship between drug dosage and reaction time is significant

The second step is to collect data and calculate the value of the test statistic. Whether you choose to do all the calculations yourself or to use the appropriate statistical software or programmable calculator, you should find the test statistic to be $t_6 = -8.723$.

The third step is to define the rejection region, decide whether or not to reject the null hypothesis, and obtain the $p$-value of the test. Since we have chosen $\alpha = 0.05$ for a two-sided test, the rejection region is defined by the $t$-score with $df = 6$ above which 0.025 of the area lies; below the negative of this $t$-score will lie 0.025 of the area, making a total area of 0.05. From Table A.3, we find that $t_{6;0.025} = 2.447$. We can then define our rejection region algebraically as

$t_6 \leq -2.447$ or $t_6 \geq +2.447$ (i.e., $|t_6| \geq 2.447$) .

Since our test statistic, calculated in the second step, was found to be $t_6 = -8.723$, which is in the rejection region, our decision is to reject $H_0$: slope $= 0$; in other words, our data provides sufficient evidence to believe that $H_1$: slope $\neq 0$ is true. From Table A.3, we find that $t_{6;0.0005} = 5.595$, and since our observed test statistic is $t_6 = -8.723$, we see that $p$-value $< 0.0005$.

To complete the fourth step of the hypothesis test, we can summarize the results of the hypothesis test as follows:

Since $t_6 = -8.723$ and $t_{6;0.025} = 2.447$, we have sufficient evidence to reject $H_0$. We conclude that the slope in the linear regression to predict reaction time from drug dosage is different from zero, or in other words, the linear relationship between dosage and reaction time is significant ($p$-value $< 0.0005$). The data suggest that the slope is less than zero, or in other words, that the linear relationship is negative.

If we do not reject the null hypothesis, then no further analysis is necessary, since we are concluding that the linear relationship is not significant. On the other hand, if we do reject the null hypothesis, then we would want to describe the linear relationship, which of course is easily done by finding the least squares line. This was already done previously for the data of Table 31-1 in Unit 10 where we were working with the same data in Table 10-2. There are also further analyses that are possible once the least squares line has been obtained for a significant regression. These analyses include hypothesis testing and confidence intervals concerning the slope, the intercept, and various predicted values; there are also prediction intervals which can be obtained. These analyses are beyond the scope of our discussion here but are generally available from the appropriate statistical software or programmable calculator.

Finally, it is important to keep in mind that the $t$ test concerning whether or not a linear relationship is significant (i.e., the test concerning a slope) is based on the assumptions that the relationship between two quantitative variables is indeed linear and that the variable $Y$ has a normal distribution for each given value of $X$. Previously, we have discussed some descriptive methods for verifying the linearity assumption. Hypothesis tests are available concerning both the linearity assumption and the normality assumption. These tests are beyond the scope of our discussion here but are generally available from the appropriate statistical software or programmable calculator.

---

### Table 31-2
## Age and Grip Strength Data (from Table 10-4)

The age (years) and right-hand grip strength (pounds of force) are recorded for each of several right-handed males.

| Age | 15 | 17 | 19 | 11 | 16 | 22 | 17 | 25 | 12 | 14 | 25 | 23 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Grip Strength | 50 | 54 | 66 | 46 | 58 | 54 | 64 | 80 | 46 | 70 | 76 | 80 |

---

**Self-Test Problem 31-1.** In Self-Test Problems 10-1 and 11-1, the between age and grip strength among right-handed males is being investigated, with regard to the prediction of grip strength from age. The Age and Grip Strength Data, displayed in Table 10-4, was recorded for several right-handed males. For convenience, we have redisplayed the data here in Table 31-2.

   (a)  Explain how the data in Table 31-2 is appropriate for a hypothesis test concerning whether or not a linear relationship is significant.

   (b)  A 0.05 significance level is chosen for a hypothesis test to see if there is any evidence that the linear relationship between age and grip strength among right-handed males is positive, or in other words, evidence that the slope is positive in the linear regression to predict grip strength from age. Complete the four steps of the hypothesis test by completing the table titled *Hypothesis Test for Self-Test Problem 31-1*. You should find that $t_{10} = +3.814$.

   (c)  Verify that the least squares line (found in Self-Test Problem 11-1) is $grp = 26 + 2(age)$, and write a one sentence interpretation of the slope of the least squares line.

   (d)  Considering the results of the hypothesis test, decide which of the Type I or Type II errors is possible, and describe this error.

   (e)  Decide whether $H_0$ would have been rejected or would not have been rejected with each of the following significance levels: (i) $\alpha = 0.01$, (ii) $\alpha = 0.10$.

# *Hypothesis Test for Self Test Problem 31-1*

Step 1    $H_0$:

          $H_1$:                                                                    $\alpha =$

Step 2

Step 3

Step 4

---

**Answers to Self-Test Problems**

**31-1**    (a) The data consists of a random sample of observations of the two quantitative variables age and grip strength.

(b)    Step 1: $H_0$: slope = 0  vs.  $H_1$: slope > 0   ($\alpha = 0.05$, one-sided test)
          OR

$H_0$: The linear relationship between age and grip strength is not significant

   vs.                                                                              $(\alpha = 0.05)$

$H_1$: The linear relationship between age and grip strength is significant

          Step 2: $t_{10} = +3.814$

          Step 3: The rejection is $t_{10} \geq +1.812$. $H_0$ is rejected; $0.0005 < p$-value $< 0.005$.

          Step 4: Since $t_{10} = 3.814$ and $t_{10;\ 0.05} = 1.812$, we have sufficient evidence to reject $H_0$. We conclude that the slope is positive in the regression to predict grip strength from age among right-handed males, or in other words, that the linear relationship between age and grip strength is positive ($0.0005 < p$-value $< 0.005$).

(c) The squares line (found in Self-Test Problem 11-1) is $grp = 26 + 2(age)$. With each increase of one year in age, grip strength increases on average by about 2 lbs.

(d) Since $H_0$ is rejected, the Type I error is possible, which is concluding that slope > 0 (i.e., the positive linear relationship is significant) when in reality slope = 0 (i.e., the linear relationship is not significant).

(e) $H_0$ would have been rejected with both $\alpha = 0.01$ and $\alpha = 0.10$.

**Summary**

  When studying a possible linear relationship between two quantitative variables $X$ and $Y$, we can obtain the Pearson correlation $r$ and the least squares line from a sample consisting of *bivariate* data, represented as $(x_1, y_1), (x_2, y_2), \ldots, (x_n, y_n)$. With $Y$ representing the response variable and $X$ representing the explanatory variable, we say that we are studying the *linear regression to predict Y from X* or the *linear regression of Y on X*. Deciding whether or not the linear relationship between $X$ and $Y$ is significant is the same as deciding whether or not the slope in the linear regression of $Y$ on $X$ is different from zero. To make this decision, the $t$ statistic with $n - 2$ degrees of freedom which can be used is

$$t_{n-2} = \frac{b - 0}{s_{Y|X} \Big/ \sqrt{\sum (x - \bar{x})^2}} \quad ,$$

where $s_{Y|X}$ is called the *standard error of estimate*, calculated by dividing the sum of the squared residuals by two less than the sample size, and taking the square root of this result. The numerator of this $t$ statistic is the difference between a sample estimate of the slope ($b$) and the hypothesized value of the slope (zero), and the denominator is an appropriate standard error. These calculations can be greatly simplified by the use of appropriate statistical software or a programmable calculator.

  The null hypothesis in this hypothesis test could be stated as either $H_0$: slope $= 0$ or $H_0$: The linear relationship between $X$ and $Y$ is not significant, and this test and the test could be one-sided or two-sided. If our focus was on finding evidence that the slope is greater than zero (i.e., finding evidence that the linear relationship was a positive one), then we would use a one-sided test. If our focus was on finding evidence that the slope is less than zero (i.e., finding evidence that the linear relationship was a negative one), then we would use a one-sided test. However, if our focus was on finding evidence that the slope is different from zero in either direction (i.e., finding evidence that the linear relationship was either positive or negative), then we would use a two-sided test.

  If we do not reject the null hypothesis, then no further analysis is necessary, since we are concluding that the linear relationship is not significant. On the other hand, if we do reject the null hypothesis, then we would want to describe the linear relationship, which of course is easily done by finding the least squares line. Further analyses that are possible once the least squares line has been obtained for a significant regression include hypothesis testing and confidence intervals concerning the slope, the intercept, and various predicted values; there are also prediction intervals which can be obtained. The $t$ test concerning whether or not a linear relationship is significant (i.e., the test concerning a slope) is based on the assumptions that the relationship between two quantitative variables is indeed linear and that the variable $Y$ has a normal distribution for each given value of $X$. Hypothesis tests are available concerning both the linearity assumption and the normality assumption.