

Regression Analysis

Last Updated:

October 3, 2023

Contents

1	Introduction	2
2	Simple Linear Regression	3
2.1	Estimating the Parameters with the Method of Least Squares . .	4
2.1.1	Estimation of β_0 and β_1	4
2.1.2	Properties of Fitted Regression Line	9
2.1.3	Estimation of σ^2	10
2.2	Hypothesis Testing on the Slope and Intercept	11
2.2.1	Using t -tests	12
2.2.2	Testing Significance	13
2.2.3	Analysis of Variance Tables (ANOVA)	13
2.3	Interval Estimation	15
2.3.1	Confidence Intervals on β_0 , β_1 , and σ^2	15
2.3.2	Interval Estimation of the mean Response	15
2.4	Coefficient of Determination	16
2.5	Correlation Coefficient	16
3	Multiple Linear Regression	18
3.1	Matrix Approach to Regression	18
3.1.1	Derivatives	19
3.2	Multiple Regression Models	19
3.2.1	Least Squares Estimation of Regression Coefficients . . .	19

Chapter 1

Introduction

The primary goal in regression is to develop a model that relates a set of explanatory variables X_1, \dots, X_p to a response variable Y , then test the model and use it for inference and prediction.

Given a set of n pairs of data Y_i and X_i , we attempt to fit a straight line to these points, using a simple regression model

$$Y_i = \beta_0 + \beta_1 X_i + \epsilon_i$$

Where ϵ_i represents an unobserved random error term, β_0 is the intercept and β_1 is the slope of the line. β_0 and β_1 are parameters that need to be estimated from observed data. The model can also be expressed in terms of $(X_i - \bar{X})$.

$$Y_i = (\beta_0 + \beta_1 \bar{X}) + \beta_1 (X_i - \bar{X}) + \epsilon_i$$

Where \bar{X} is the sample mean

$$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$$

This proposed model is linear in the parameters β_0 , β_1 , and would still be referred to as linear if we had X_i^2 instead of X_i . This model also makes the assumption that the random error terms ϵ_i are uncorrelated, have mean 0, and variance σ^2 . Under these assumptions, we have

$$E(Y_i) = \beta_0 + \beta_1 X_i$$

$$\text{Var}(Y_i) = \sigma^2$$

Chapter 2

Simple Linear Regression

The primary goal in regression is to develop a model that relates a set of explanatory variables X_1, \dots, X_p to a response variable Y , then test the model and use it for inference and prediction.

Given a set of n pairs of data Y_i and X_i , we attempt to fit a straight line to these points, using a simple regression model

$$Y_i = \beta_0 + \beta_1 X_i + \epsilon_i$$

Where ϵ_i represents an unobserved random error term, β_0 is the intercept and β_1 is the slope of the line. β_0 and β_1 are parameters that need to be estimated from observed data. The model can also be expressed in terms of $(X_i - \bar{X})$.

$$Y_i = (\beta_0 + \beta_1 \bar{X}) + \beta_1 (X_i - \bar{X}) + \epsilon_i$$

Where \bar{X} is the sample mean

$$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$$

This proposed model is linear in the parameters β_0 , β_1 , and would still be referred to as linear if we had X_i^2 instead of X_i . This model also makes the assumption that the random error terms ϵ_i are uncorrelated, have mean 0, and variance σ^2 . Under these assumptions, we have

$$E(Y_i) = \beta_0 + \beta_1 X_i$$

$$\text{Var}(Y_i) = \text{Var}(\beta_0 + \beta_1 X_i + \epsilon_i) = \sigma^2$$

Thus the mean of Y is a linear function of X however the variance of Y does not depend on a value of X .

The parameters β_0 and β_1 are called the regression coefficients. The slope β_1 is the change in the mean of the distribution of Y produced by a unit change

in Y . If the range of data on X includes $x = 0$, then the intercept β_0 is the mean of the distribution of the response Y when $x = 0$. If the range of x does not include zero, then β_0 has no practical interpretation.

2.1 Estimating the Parameters with the Method of Least Squares

The parameters β_0, β_1 are unknown and must be estimated from the data. Suppose we have n pairs of data $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$.

2.1.1 Estimation of β_0 and β_1

The method of least squares is the most popular approach to fitting a regression model. Let Q be the sum of the error terms squared

$$Q = \sum_{i=1}^n \epsilon_i^2 = \sum_{i=1}^n (Y_i - \beta_0 - \beta_1 X_i)^2$$

Then we want to minimize Q with respect to the parameters β_1, β_2 ,

$$\frac{\partial Q}{\partial \beta_0} = -2 \sum_{i=1}^n (Y_i - \beta_0 - \beta_1 X_i) = 0$$

$$\frac{\partial Q}{\partial \beta_1} = -2 \sum_{i=1}^n (Y_i - \beta_0 - \beta_1 X_i) X_i = 0$$

We can rearrange these equations to get the following equations

$$\begin{aligned} & -2 \sum_{i=1}^n (Y_i - \beta_0 - \beta_1 X_i) = 0 \\ \implies & \sum_{i=1}^n Y_i - \sum_{i=1}^n \beta_0 - \sum_{i=1}^n \beta_1 X_i = 0 \\ \implies & \sum_{i=1}^n Y_i = n\beta_0 - \beta_1 \sum_{i=1}^n X_i \\ & -2 \sum_{i=1}^n (Y_i - \beta_0 - \beta_1 X_i) X_i = 0 \\ \implies & \sum_{i=1}^n Y_i X_i - \sum_{i=1}^n \beta_0 X_i - \sum_{i=1}^n \beta_1 X_i^2 = 0 \\ \implies & \sum_{i=1}^n Y_i X_i = \beta_0 \sum_{i=1}^n X_i - \beta_1 \sum_{i=1}^n X_i^2 \end{aligned}$$

These 2 equations are known as the normal equations and the solutions to them, call them b_0 , b_1 , are

$$b_0 = \bar{Y} - b_1 \bar{X}$$

$$b_1 = \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{\sum_{i=1}^n (X_i - \bar{X})^2} = \frac{\sum_{i=1}^n (X_i - \bar{X})Y_i}{\sum_{i=1}^n (X_i - \bar{X})^2} = \sum_{i=1}^n k_i Y_i$$

with

$$k_i = \frac{X_i - \bar{X}}{\sum_{i=1}^n (X_i - \bar{X})^2}$$

We also sometimes use a more compact notation, by denoting the corrected sum of squares for X and the sum of cross products of X_i Y_i as

$$S_{xx} = \sum_{i=1}^n x_i^2 - \frac{1}{n} \left(\sum_{i=1}^n x_i \right)^2 = \sum_{i=1}^n (x_i - \bar{x})^2$$

$$S_{xy} = \sum_{i=1}^n x_i y_i - \frac{1}{n} \left(\sum_{i=1}^n x_i \right) \left(\sum_{i=1}^n y_i \right) = \sum_{i=1}^n y_i (x_i - \bar{x})$$

So, we can write

$$b_1 = \frac{S_{xy}}{S_{xx}}$$

The observed difference between Y_i and the corresponding fitted value \hat{Y}_i is a residual. The i th residual is

$$e_i = y_i - \hat{y}_i = y_i - (b_0 + b_1 x_i)$$

Note that k_i has important properties, such as

$$\sum_{i=1}^n k_i = 0, \quad \sum_{i=1}^n k_i X_i = 1, \quad \sum_{i=1}^n k_i^2 = \frac{1}{\sum_{i=1}^n (X_i - \bar{X})^2}$$

$$\begin{aligned} \sum_{i=1}^n k_i &= \frac{\sum_{i=1}^n (X_i - \bar{X})}{\sum_{i=1}^n (X_i - \bar{X})^2} \\ &= \frac{n\bar{X} - n\bar{X}}{\sum_{i=1}^n (X_i - \bar{X})^2} = 0 \end{aligned}$$

$$\begin{aligned}
\sum_{i=1}^n k_i X_i &= \frac{\sum_{i=1}^n (X_i - \bar{X}) X_i}{\sum_{i=1}^n (X_i - \bar{X})^2} \\
&= \frac{\sum_{i=1}^n (X_i^2 - X_i \bar{X})}{\sum_{i=1}^n (X_i - \bar{X})^2} \\
&= \frac{\sum_{i=1}^n X_i^2 - \bar{X} \sum_{i=1}^n X_i}{\sum_{i=1}^n (X_i - \bar{X})^2} \\
&= \frac{\sum_{i=1}^n X_i^2 - n \bar{X}^2}{\sum_{i=1}^n (X_i^2 - 2X_i \bar{X} + \bar{X}^2)} \\
&= \frac{\sum_{i=1}^n X_i^2 - n \bar{X}^2}{\sum_{i=1}^n X_i^2 - 2n \bar{X}^2 + n \bar{X}^2} = 1
\end{aligned}$$

$$\begin{aligned}
\sum_{i=1}^n k_i^2 &= \sum_{i=1}^n \left(\frac{X_i - \bar{X}}{\sum_{i=1}^n (X_i - \bar{X})^2} \right)^2 \\
&= \frac{\sum_{i=1}^n (X_i - \bar{X})^2}{\sum_{i=1}^n (X_i - \bar{X})^4} \\
&= \frac{1}{\sum_{i=1}^n (X_i - \bar{X})^2}
\end{aligned}$$

The equation for the fitted line is then

$$\hat{Y} = b_0 + b_1 X$$

Or alternatively using $X - \bar{X}$,

$$\hat{Y} = (b_0 + b_1 \bar{X}) + b_1 (X - \bar{X})$$

Theorem 2.1.1 (Gauss Markov Theorem). *The least square estimators b_0 , b_1 are unbiased and have minimum variance among all unbiased linear estimators.*

Proof. Consider an unbiased linear estimator

$$\hat{\beta}_1 = \sum_{i=1}^n c_i Y_i$$

$\hat{\beta}_1$ must satisfy $E(\hat{\beta}_1) = \beta_1$.

$$\begin{aligned}
\beta_1 &= E(\hat{\beta}_1) \\
&= E\left(\sum_{i=1}^n c_i Y_i\right) \\
&= \sum_{i=1}^n c_i E(Y_i) \\
&= \sum_{i=1}^n c_i (\beta_0 + \beta_1 X_i) \\
&= \beta_0 \sum_{i=1}^n c_i + \beta_1 \sum_{i=1}^n c_i X_i
\end{aligned}$$

Therefore, $\sum_{i=1}^n c_i = 0$, and $\sum_{i=1}^n c_i X_i = 1$. We can also see that the variance is

$$\text{Var}(\hat{\beta}_1) = \sum_{i=1}^n c_i^2 \text{Var}(Y_i) = \sigma^2 \sum_{i=1}^n c_i^2$$

Now, set $c_i = k_i + d_i$ where k_i is as defined previously above and d_i are arbitrary constants. We want to show that the variance is minimized, so

$$\begin{aligned}
\text{Var}(\hat{\beta}_1) &= \sum_{i=1}^n c_i^2 \text{Var}(Y_i) \\
&= \sigma^2 \sum_{i=1}^n c_i^2 \\
&= \sigma^2 \sum_{i=1}^n (k_i + d_i)^2 \\
&= \sigma^2 \left(\sum_{i=1}^n k_i^2 + 2 \sum_{i=1}^n k_i d_i + \sum_{i=1}^n d_i^2 \right)
\end{aligned}$$

Note that the variance of b_1 is

$$\text{Var}(b_1) = \text{Var}\left(\sum_{i=1}^n k_i Y_i\right) = \sigma^2 \sum_{i=1}^n k_i^2 = \sigma^2 \frac{1}{\sum_{i=1}^n (X_i - \bar{X})^2}$$

Now notice that there is a relationship between the variance of $\hat{\beta}_1$ and b_1 , namely that the variance of $\hat{\beta}_1$ is the same as b_1 plus an additional constants but these

constants are indeed 0.

$$\begin{aligned}
\sum_{i=1}^n k_i d_i &= \sum_{i=1}^n k_i (c_i - k_i) \\
&= \sum_{i=1}^n k_i c_i - \sum_{i=1}^n k_i^2 \\
&= \sum_{i=1}^n c_i \frac{X_i - \bar{X}}{\sum_{i=1}^n (X_i - \bar{X})^2} - \frac{1}{\sum_{i=1}^n (X_i - \bar{X})^2} \\
&= \frac{\sum_{i=1}^n c_i X_i - \sum_{i=1}^n c_i \bar{X}}{\sum_{i=1}^n (X_i - \bar{X})^2} - \frac{1}{\sum_{i=1}^n (X_i - \bar{X})^2}
\end{aligned}$$

We know that $\sum_{i=1}^n c_i = 0$ and $\sum_{i=1}^n c_i X_i = 1$, so this becomes

$$\sum_{i=1}^n k_i d_i = \frac{1 - 0}{\sum_{i=1}^n (X_i - \bar{X})^2} - \frac{1}{\sum_{i=1}^n (X_i - \bar{X})^2} = 0$$

Therefore,

$$\text{Var}(\hat{\beta}_1) = \sigma^2 \left(\sum_{i=1}^n k_i^2 + \sum_{i=1}^n d_i^2 \right)$$

Clearly the variance is minimized when $d_i = 0$ for all i , thus

$$\text{Var}(\hat{\beta}_1) = \sigma^2 \sum_{i=1}^n k_i^2 = \text{Var}(b_1)$$

Thus the least squares estimator b_1 has minimum variance along all unbiased estimators. \square

We may write

$$\hat{Y} = b_0 + b_1 X$$

for the estimated or fitted line, and

$$e_i = Y_i - \hat{Y}_i$$

for the estimated i^{th} residual. The estimate for the variance σ^2 is then

$$\hat{\sigma}^2 = \frac{\sum_{i=1}^n e_i^2}{n - 2}$$

The estimate of the variance σ^2 is also known as the mean square error (MSE).

2.1.2 Properties of Fitted Regression Line

- (i) $\sum_{i=1}^n e_i = 0$. Recall that $\hat{Y} = b_0 + b_1 X = (b_0 + b_1 \bar{X}) + b_1(X - \bar{X})$, and

$$\bar{Y} = b_0 + b_1 \bar{X}$$

So $\hat{Y} = \bar{Y} + b_1(X - \bar{X})$, then

$$\begin{aligned} \sum_{i=1}^n e_i &= \sum_{i=1}^n (Y_i - \hat{Y}_i) \\ &= \sum_{i=1}^n Y_i - \sum_{i=1}^n \hat{Y}_i \\ &= \sum_{i=1}^n Y_i - \sum_{i=1}^n (\bar{Y} + b_1(X_i - \bar{X})) \\ &= n\bar{Y} - n\bar{Y} + b_1 \sum_{i=1}^n (X_i - \bar{X}) \\ &= n\bar{Y} - n\bar{Y} + b_1(n\bar{X} - n\bar{X}) = 0 \end{aligned}$$

- (ii) $\sum_{i=1}^n Y_i = \sum_{i=1}^n \hat{Y}_i$. This follows from the previous property since

$$\sum_{i=1}^n e_i = \sum_{i=1}^n Y_i - \sum_{i=1}^n \hat{Y}_i = 0 \implies \sum_{i=1}^n Y_i = \sum_{i=1}^n \hat{Y}_i$$

- (iii) $\sum_{i=1}^n X_i e_i = 0$. This can be shown from the definition

$$\begin{aligned} \sum_{i=1}^n X_i e_i &= \sum_{i=1}^n X_i (Y_i - \hat{Y}_i) \\ &= \sum_{i=1}^n X_i (Y_i - b_0 - b_1 X_i) \\ &= \sum_{i=1}^n X_i Y_i - b_0 \sum_{i=1}^n X_i - b_1 \sum_{i=1}^n X_i^2 \\ &= b_0 \sum_{i=1}^n X_i + b_1 \sum_{i=1}^n X_i - b_0 \sum_{i=1}^n X_i - b_1 \sum_{i=1}^n X_i \\ &= 0 \end{aligned}$$

This is significant because it tells us that the dot product between the vector of explanatory variables $\vec{X} = (X_1, \dots, X_n)^T$ is orthogonal to the

vector of error terms $\vec{e} = (e_1, \dots, e_n)^T$, and from the previous property we get that

$$\vec{e} \cdot \mathbf{1}_n = \sum_{i=1}^n e_i = 0$$

Hence the vectors $\{1_n, X - \bar{X}1_n\}$ are linearly independent and form a basis of the estimation space.

(iv) By applying the Pythagorean Theorem to the previous property we get

$$\begin{aligned} \|Y\|^2 &= \|\hat{Y}\|^2 + \|Y - \hat{Y}\|^2 \\ \sum_{i=1}^n Y_i^2 &= \sum_{i=1}^n \hat{Y}_i^2 + \sum_{i=1}^n e_i^2 \\ &= \sum_{i=1}^n \bar{Y}^2 + b_1^2 \sum_{i=1}^n (X_i - \bar{X})^2 + \sum_{i=1}^n e_i^2 \\ \implies \sum_{i=1}^n Y_i^2 - n\bar{Y}^2 &= b_1^2 \sum_{i=1}^n (X_i - \bar{X})^2 + \sum_{i=1}^n e_i^2 \\ \sum_{i=1}^n (Y_i - \bar{Y})^2 &= b_1^2 \sum_{i=1}^n (X_i - \bar{X})^2 + \sum_{i=1}^n (Y_i - \hat{Y}_i)^2 \end{aligned}$$

This shows us the the total sum of squares is equal to the regression sum of squares plus the error sum of squares.

(v) The point (\bar{X}, \bar{Y}) is on the fitted line.

(vi) The sum of residuals weighted by their corresponding fitted value is 0, that is

$$\sum_{i=1}^n y_i e_i = 0$$

(vii) Under the normality assumption, $e_i \stackrel{\text{iid}}{\sim} N(0, \sigma^2)$. The method of maximum likelihood leads to the method of least squares.

$$L(\beta_0, \beta_1, \sigma^2) = \left(\frac{1}{\sqrt{2\pi}\sigma} \right)^n \exp \left(-\frac{1}{2\sigma^2} \sum_{i=1}^n \epsilon_i^2 \right)$$

So maximizing $L(\beta_0, \beta_1, \sigma^2)$ is equivalent to minimizing $\sum \epsilon_i^2$.

2.1.3 Estimation of σ^2

We need to estimate σ^2 to test hypotheses and construct interval estimates pertinent to the regression model. Ideally we would like this estimate not to depend on the adequacy of the fitted model. This is only possible when there are several observations on y for at least one value of x , or when prior information

concerning σ^2 is available. When this approach cannot be used, the estimate of σ^2 is obtained from the residual or error sum of squares.

$$SSE = \sum_{i=1}^n e_i^2 = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

We can substitute \hat{y}_i for $b_0 + b_1 x_i$ and simplify to get

$$SSE = \sum_{i=1}^n y_i^2 - n\bar{y}^2 - b_1 S_{xy}$$

Moreover, the correct sum of squares of the response variable is

$$SST = \sum_{i=1}^n y_i^2 - n\bar{y}^2 = \sum_{i=1}^n (y_i - \bar{y})^2$$

Thus,

$$SSE = SST - b_1 S_{xy}$$

The residual sum of squares has $n - 2$ degrees of freedom, because we reserve 2 degrees of freedom for the estimators b_0 , b_1 . We will later show that the expected value for SSE is

$$E(SSE) = (n - 2)\sigma^2$$

So an unbiased estimator of σ^2 is

$$\hat{\sigma}^2 = \frac{SSE}{n - 2} = MSR$$

The quantity MSR is known as the **residual mean square**. The root of $\hat{\sigma}^2$ is known as the **standard error of regression**.

2.2 Hypothesis Testing on the Slope and Intercept

To perform hypotheses tests and construct confidence intervals, we require that we make the additional assumption that the model errors ϵ_i are normally distributed. Thus, the complete assumptions are that the errors are normally and independently distributed with mean 0 and variance σ^2 , written as $\{\epsilon_i\} \stackrel{\text{iid}}{\sim} N(0, \sigma^2)$. We will discuss how these assumptions can be checked through residual analysis later.

Suppose that we have the model $Y_i = \beta_0 + \beta_1 X_i + \epsilon_i$, where $\{\epsilon_i\} \stackrel{\text{iid}}{\sim} N(0, \sigma^2)$. Then

$$(a) \quad \frac{b_1 - \beta_1}{se(b_1)} \sim t_{n-2} \text{ where } se^2(b_1) = \frac{MSE}{\sum_{i=1}^n (X_i - \bar{X})^2}$$

(b) $\frac{b_0 - \beta_0}{se(b_0)} \sim t_{n-2}$ where

$$se^2(b_0) = MSE \left(\frac{1}{n} + \frac{\bar{X}^2}{\sum_{i=1}^n (X_i - \bar{X})^2} \right)$$

(c) MSE is an unbiased estimate of σ^2 and is independent of b_0, b_1 . Furthermore

$$\frac{(n-2)MSE}{\sigma^2} \sim \chi_{n-2}^2$$

Proof. Proof will be shown when we generalize this using matrices in later sections. \square

2.2.1 Using t -tests

Suppose we want to test that the slope is equal to a constant, β , we have the hypotheses

$$H_0 : \beta_1 = \beta, H_1 : \beta_1 \neq \beta$$

Since $\{\epsilon_i\} \stackrel{\text{iid}}{\sim} N(0, \sigma^2)$, the observations y_i are normally distributed with $\beta_0 + \beta_1 x_i$ and variance σ^2 . Then, b_1 is a linear combination of the observations, so it is normally distributed with mean β_1 and variance σ^2/S_{xx} . Therefore, our test statistic becomes

$$Z_0 = \frac{b_1 - \beta}{\sqrt{\sigma^2/S_{xx}}}$$

If the null hypothesis is true, then $Z_0 \sim N(0, 1)$. If σ^2 is known then we would use Z_0 to test our hypotheses. However, σ^2 is typically unknown. We've seen that MSE is an unbiased estimator for σ^2 , and we've established that $(n-2)MSE/\sigma^2 \sim \chi_{n-2}^2$.

$$t_0 = \frac{b_1 - \beta}{\sqrt{MSE/S_{xx}}}$$

If the null hypothesis is true, $t_0 \sim t_{n-2}$. We compare the observed value t_0 with the upper $\alpha/2$, of the t_{n-2} distribution. So we reject the null hypothesis

$$|t_0| > t_{\alpha/2, n-2}$$

We can also test with the p -value. From the equation for t_0 , the denominator is called the **estimated standard error** of the slope.

$$se(b_1) = \sqrt{\frac{MSE}{S_{xx}}}$$

So, we often write t_0 is

$$t_0 = \frac{b_1 - \beta}{se(\beta_1)}$$

We test the intercept in a similar manner,

$$H_0 : \beta_0 = \beta, \quad H_1 : \beta_0 \neq \beta$$

We use a similar test statistic,

$$t_0 = \frac{b_0 - \beta}{se(b_0)}$$

and we reject the null hypothesis when $|t_0| > t_{\alpha/2, n-2}$.

2.2.2 Testing Significance

A special case for hypotheses is

$$H_0 : B_1 = 0, \quad H_1 : B_1 \neq 0$$

These hypotheses relate to the **significance of regression**. Failing to reject the null hypothesis means there is no linear relationship between x and y , we would reject the null hypothesis when $|t_0| > t_{\alpha/2, n-2}$.

2.2.3 Analysis of Variance Tables (ANOVA)

Analysis of variance can be used to test significance of regression. The analysis of variance of variance is based on a partitioning of the total variability of the response variable y , given by

$$y_i - \bar{y} = (\hat{y}_i - \bar{y}) + (y_i - \hat{y}_i)$$

Then, taking the sum of the square of both sides

$$\sum_{i=1}^n (y_i - \bar{y})^2 = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 + \sum_{i=1}^n (y_i - \hat{y}_i)^2 + 2 \sum_{i=1}^n (\hat{y}_i - \bar{y})(y_i - \hat{y}_i)$$

Notice that

$$\begin{aligned} 2 \sum_{i=1}^n (\hat{y}_i - \bar{y})(y_i - \hat{y}_i) &= 2 \sum_{i=1}^n \hat{y}_i (y_i - \hat{y}_i) - 2\bar{y} \sum_{i=1}^n (y_i - \hat{y}_i) \\ &= 2 \sum_{i=1}^n \hat{y}_i e_i - 2\bar{y} \sum_{i=1}^n e_i = 0 \end{aligned}$$

Therefore,

$$\sum_{i=1}^n (y_i - \bar{y})^2 = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 + \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

The left side is the corrected sum of squares of the observations, which we denote by SST or $SSTO$. Notice that $y_i - \hat{y}_i = e_i$, so that term is the sum of residuals squared SSE . We call $\sum_{i=1}^n (\hat{y}_i - \bar{y})^2$ the **regression sum of squares**. So we have

$$SST = SSR + SSE$$

The regression sum of squares can also be computed by

$$SSR = b_1^2 S_{xx}$$

The **degrees of freedom** for each sum of squares is as follows.

- The total sum of squares SST has $df_T = n - 1$ since we lose a degree of freedom for the constraint

$$\sum_{i=1}^n (y_i - \bar{y}) = 0$$

- The regression sum of squares SSR has $df_R = p - 1$ where p is the number of variables (including y).
- The residual sum of squares SSE has $df_E = n - 2$ degrees of freedom since 2 constraints are placed on $e_i = y_i - \hat{y}_i$ with the estimation for β_0 and β_1 .

We create a table to summarize our results from statistical analysis.

Source	SS	DF	MS=SS/df	E(MS)
Regression	$SSR = b_1^2 \sum (X_i - \bar{X})^2$	$p - 1$	MSR	$\sigma^2 + \beta_1^2 \sum (X_i - \bar{X})^2$
Error	$SSE = \sum (Y_i - \hat{Y}_i)^2$	$n - p$	MSE	σ^2
Total	$SSTO = \sum (Y_i - \bar{Y})^2$	$n - 1$		

Each of the sums of squares is a quadratic form where the rank of the corresponding matrix is the degrees of freedom indicated. Chochran's theorem applies and we conclude that the quadratic forms are independent and have chi-squared distributions. Note that

$$\frac{SSR}{\sigma^2} = \frac{b_1^2 \sum (X_i - \bar{X})^2}{\sigma^2} \sim \chi_{p-1}^2$$

$$\frac{SSE}{\sigma^2} = \frac{\sum (Y_i - \hat{Y}_i)^2}{\sigma^2} \sim \chi_{n-p}^2$$

Then, the ratio between 2 chi-squared distributions divided by their degrees of freedom has a F-distribution with their respective degrees of freedom.

$$F = \frac{SSR/\sigma^2(p-1)}{SSE/\sigma^2(n-p)} = \frac{SSR/(p-1)}{SSE/(n-p)} = \frac{MSR}{MSE} \sim F_{p-1, n-p}$$

The degrees of freedom are determined by the amount of data required to calculate each expression.

To summarize, the ANOVA table indicates how one can test the null hypothesis

$$H_0 : \beta_1 = 0$$

$$H_1 : \beta_1 \neq 0$$

The null Hypothesis is that the slope of the line is equal to 0. Under the null hypothesis, the expected mean square for regression and the expected mean square error are separate independent estimates of the variance σ^2 .

2.3 Interval Estimation

2.3.1 Confidence Intervals on β_0 , β_1 , and σ^2 .

The width of the confidence intervals are a measure of the quality of the regression line. If the error is normally and independently distributed by our assumption, then $(b_1 - \beta_1)/se(b_1)$ and $(b_0 - \beta_0)/se(b_0)$ follow a t distribution with $n - 2$ degrees of freedom. So, a $100(1 - \alpha)$ percent. The confidence interval for the slope β_1 is

$$b_1 - t_{\alpha/2, n-2} se(b_1) \leq b_1 \leq b_1 + t_{\alpha/2, n-2} se(b_1)$$

and for the intercept β_0 ,

$$b_0 - t_{\alpha/2, n-2} se(b_0) \leq b_0 \leq b_0 + t_{\alpha/2, n-2} se(b_0)$$

The interpretation for these intervals is, if we were to take repeated samples of the same size at the same x levels and construct 95% CIs on the slope for each sample, then 95% of those intervals will contain the true value of β_1 .

As we've seen earlier, the sampling distribution of $(n - 2)MSE/\sigma^2$ follows a chi-square distribution with $n - 2$ degrees of freedom, so

$$P \left\{ \chi_{1-\alpha/2, n-2}^2 \leq \frac{(n - 2)MSE}{\sigma^2} \leq \chi_{\alpha/2, n-2}^2 \right\}$$

Thus the $100(1 - \alpha)$ percent CI on σ^2 is

$$\frac{(n - 2)MSE}{\chi_{\alpha/2, n-2}^2} \leq \sigma^2 \leq \frac{(n - 2)MSE}{\chi_{1-\alpha/2, n-2}^2}$$

2.3.2 Interval Estimation of the mean Response

Another important part of the regression model is estimating the mean response $E(y)$ for a particular regressor variable x . Assuming that x_0 is any value of the regressor variable within the range of the original data on x that we used to

create the model. Then, an unbiased estimator for $E(y|x_0)$ can be found from the fitting model

$$\widehat{E(y|x_0)} = \hat{\mu}_{y|x_0} = b_0 + b_1x_0$$

Note that $\hat{\mu}_{y|x_0}$ follows a normal distribution since it is a linear combination of the observations y_i . The variance is

$$\text{Var}(\hat{\mu}_{y|x_0}) = \text{Var}(b_0 + b_1x_0) = \text{Var}(\bar{y} - b_1(x_0 - \bar{x})) = \frac{\sigma^2}{n} + \frac{\sigma^2(x_0 - \bar{x})^2}{S_{xx}}$$

The sampling distribution for

$$\frac{\hat{\mu}_{y|x_0} - E(y|x_0)}{\sqrt{MSE(1/n + (x_0 - \bar{x})^2/S_{xx})}}$$

is a t distribution with $n - 2$ degrees of freedom. Then the CI is given as

$$\left[\hat{\mu}_{y|x_0} \pm t_{\alpha/2, n-2} \sqrt{MSE \left(\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{S_{xx}} \right)} \right]$$

2.4 Coefficient of Determination

The quantity

$$R^2 = \frac{SSR}{SST} = 1 - \frac{SSE}{SST}$$

is called the **coefficient of determination**. R^2 is also called the proportion of variation explained by the regressor x since SST is a measure of variability in y without considering the effect of x , and SSE is the variability in y after considering x . Since $0 \leq SSE \leq SST$, then $0 \leq R^2 \leq 1$. An R^2 value close to 1 means **most of the variability of y is explained by x** .

2.5 Correlation Coefficient

The pearson correlation coefficient, denoted by ρ , related to b_1 is given as

$$\rho = \frac{\text{Cov}(X, Y)}{\sqrt{\text{Var}(X)}\sqrt{\text{Var}(Y)}}$$

This measures the linear correlation between 2 variables. When applied to a sample,

$$r = b_1 \left(\frac{S_{xx}}{SST} \right)^{\frac{1}{2}} = \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{\sqrt{\sum_{i=1}^n (X_i - \bar{X})^2} \sqrt{\sum_{i=1}^n (Y_i - \bar{Y})^2}} = \frac{S_{xy}}{(S_{xx}SST)^{1/2}}$$

Note that $-1 \leq r \leq 1$. To test hypotheses on ρ , we have 2 cases. The hypotheses for testing if the correlation is 0 is as follows

$$H_0 : \rho = 0, H_1 : \rho \neq 0$$

When testing the null hypothesis $\rho = 0$, we use a t statistic given as

$$t = \frac{r\sqrt{n-2}}{\sqrt{1-r^2}} \sim t_{n-2}$$

When testing

$$H_0 : \rho = \rho_0, \quad H_1 : \rho \neq \rho_0$$

We use a Z statistic,

$$Z = \operatorname{arctanh} r = \frac{1}{2} \ln \left(\frac{1+r}{1-r} \right) \sim N \left(\mu_z, \frac{1}{n-3} \right)$$

where

$$\mu_z = \frac{1}{2} \ln \left(\frac{1+\rho}{1-\rho} \right)$$

Now we can standardize our statistic to obtain a standard normal test statistic

$$Z_0 = (\operatorname{arctanh}(r) - \operatorname{arctanh}(\rho_0))\sqrt{n-3}$$

We can obtain our confidence interval with

Chapter 3

Multiple Linear Regression

We call a regression model with more than one regressor variable a **multiple regression model**.

3.1 Matrix Approach to Regression

We will first cover simple linear regression in matrix form.

Let $Y = [Y_1, \dots, Y_n]^T$ be a column data vector, and we'll define the expected value as

$$E(Y) = \begin{bmatrix} E(Y_1) \\ \vdots \\ E(Y_n) \end{bmatrix}$$

Proposition 3.1.1. *If $Z = AY + B$ for a matrix of constants A , and B , then*

$$E(Z) = AE(Y) + B$$

Proof. Simply from the definition of expectations on vectors,

$$E(Z_i) = E\left(\left[\sum_j a_{ij}Y_j\right] + b_i\right) = \sum_j a_{ij}E(Y_j) + b_i$$

So

$$E(Z) = AE(Y) + B$$

□

Definition 3.1.1. *The covariance of a vector of data*

$$Y = \begin{bmatrix} Y_1 \\ \vdots \\ Y_n \end{bmatrix}$$

is

$$\text{Cov}(Y) = E([Y - E(Y)][Y - E(Y)]^T) = \Sigma$$

Proposition 3.1.2. $\text{Cov}(AY) = A\Sigma A^T$.

Definition 3.1.2. A random vector Y has a multivariate normal distribution if its density is given by

$$f(y_1, \dots, y_n) = \frac{|\Sigma|^{-1/2}}{\exp\left(-\frac{1}{2}(Y - \mu)^T \Sigma^{-1}(Y - \mu)\right)}$$

where

$$Y^T = (y_1, \dots, y_n), \quad \mu^T = (\mu_1, \dots, \mu_n)$$

we denote this by

$$Y \sim N_n(\mu, \Sigma)$$

Theorem 3.1.1. Let $Y \sim N_n(\mu, \Sigma)$. Let A be an arbitrary $p \times n$ matrix of constants. Then

$$Z = AY + B \sim N_p(A\mu + B, A\Sigma A^T)$$

This theorem implies that any linear combination of normal variates has a normal distribution. This theorem won't be proved here.

3.1.1 Derivatives

- $z = a'y \rightarrow \frac{\partial z}{\partial y} = a$
- $z = y'y \rightarrow \frac{\partial z}{\partial y} = 2y$
- $z = a' Ay \rightarrow \frac{\partial z}{\partial y} = A'a$
- $z = y' Ay \rightarrow \frac{\partial z}{\partial y} = A'y + Ay$
- If A is symmetric, then $z = y' Ay \rightarrow \frac{\partial z}{\partial y} = 2A'y$

3.2 Multiple Regression Models

Suppose we have 2 regressor variables, a multiple regression model that may describe a relationship with our data is

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \epsilon$$

The parameter β_1 indicates the expected change in response per unit change in x_1 when x_2 is held constant. Similarly β_2 measures the change in y per unit change in x_2 when x_1 is held constant.

3.2.1 Least Squares Estimation of Regression Coefficients

The method of **least squares** can be used to estimate the regression coefficients. Suppose $n > k$ observations are available, and let y_i denote the i th observed response x_{ij} denote the i th observation on level