

## Chapter 3 Summary: Linear Regression

Linear regression (yes, the one from middle or high school) may seem boring and not worth your time, but it is! It is widely used and many fancier models are generalizations or extensions of linear regression. So seriously, this is important.

### 3.1 Simple Linear Regression

**Simple Linear Regression** assumes that the response variable,  $Y$  has a linear relationship to a single predictor variable  $X$ .

$$Y \approx \beta_0 + \beta_1 X \quad (3.1)$$

#### Note

This is the same relationship as  $y = mx + b$  you likely remember from school.  $\beta_0$  represents the intercept,  $b$  while  $\beta_1$  represents the slope,  $m$ .

$\beta_0$  and  $\beta_1$  are referred to as the **coefficients** or **parameters**.

We will use our training data to estimate these parameters and then indicate and represent our **model** as

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 X \quad (3.2)$$

#### Note

We call this “hat notation.” In statistics, a hat over something indicates it’s either an **estimator** or an estimated value.

#### 3.1.1 Estimating the Coefficients

Since we do not know the true relationship between  $X$  and  $Y$ , we use the training data to estimate  $\beta_0$  and  $\beta_1$

#### Note

Just like you did in middle/high school when calculating the line of best fit

We will use **least squares** as the criteria to determine what values to use for our parameters,  $\beta_0$  and  $\beta_1$ .

The  $i$ th **residual** is the difference between the  $i$ th response variable and our prediction for that variable

$$e_i = y_i - \hat{y}_i$$

The **residual sum of squares** is defined

$$\begin{aligned}
\text{RSS} &= e_1^2 + e_2^2 + \cdots + e_n^2 \\
&= (y_1 - \hat{y}_1)^2 + (y_2 - \hat{y}_2)^2 + \cdots + (y_n - \hat{y}_n)^2 \\
&= (y_1 - \hat{\beta}_0 - \hat{\beta}_1 x_1)^2 + (y_2 - \hat{\beta}_0 - \hat{\beta}_1 x_2)^2 + \cdots + (y_n - \hat{\beta}_0 - \hat{\beta}_1 x_n)^2
\end{aligned}$$

The parameters can be calculated as:

$$\begin{aligned}
\hat{\beta}_1 &= \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} \\
\hat{\beta}_0 &= \bar{y} - \hat{\beta}_1 \bar{x}
\end{aligned} \tag{3.4}$$

where  $\bar{y}$  and  $\bar{x}$  are the sample mean, defined below

$$\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i \quad \text{and} \quad \bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$$

### 3.1.2 Assessing the Accuracy of the Coefficient Estimates

#### Standard Error

**Standard error** tells us the average amount that an estimate differs from the actual value.

$$\text{Var}(\hat{\mu}) = \text{SE}(\hat{\mu})^2 = \frac{\sigma^2}{n} \tag{3.7}$$

$$\text{SE}(\hat{\beta}_0)^2 = \sigma^2 \left[ \frac{1}{n} + \frac{\bar{x}^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \right] \tag{3.8}$$

$$\text{SE}(\hat{\beta}_1)^2 = \frac{\sigma^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \tag{also 3.8}$$

where  $\sigma^2 = \text{Var}(\epsilon)$  While we generally don't know  $\sigma^2$ , it can be estimated from the data.

The estimate of  $\sigma$  is called **residual standard error**

$$\sigma = \text{RSE} = \sqrt{\frac{\text{RSS}}{n-2}}$$

#### Confidence Intervals

Standard errors can be used to compute **confidence intervals**. The 95% confidence interval for  $\beta_0$  and  $\beta_1$  are approximately:

$$\begin{aligned}
&\hat{\beta}_0 \pm 2 \cdot \text{SE}(\hat{\beta}_0) \\
&\hat{\beta}_1 \pm 2 \cdot \text{SE}(\hat{\beta}_1)
\end{aligned} \tag{3.11}$$

This means there is approximately a 95% chance that the interval

$$\left[ \hat{\beta}_0 - 2 \cdot \text{SE}(\hat{\beta}_0), \hat{\beta}_0 + 2 \cdot \text{SE}(\hat{\beta}_0) \right] \tag{almost 3.10}$$

contains the true value for  $\beta_0$ . **NOTE: This is an incorrect interpretation of a confidence interval, but it is what the book writes. Please see the warning below**

**Interesting Tidbit!**

Did you know that the  $t$ -test was developed by in order to make better beer? William Sealy Gosset, while the head experimental brewer at Guinness, developed the  $t$ -test as a way to study the quality of the barley used in brewing Guinness.

**Note**

1.96 is closer to the correct value than 2. This value comes from the Z-score value for a 97.5% quantile of a  $t$ -distribution. You are likely to see these in any statistics class.

**Warning**

A 95% confidence interval does **NOT** mean we there is a 95% chance that the true parameter lies within the range.

What it really means is that if we sampled the data 100 times, each time calculating the parameter and confidence interval, 95% of those confidence intervals would contain the true value of the parameter. It's a small distinction, but I wanted to make it, even if the book did not.

We generally say that we are 95% confident that the true parameter lies in the range, not that the probability is 0.95.

**Hypothesis Testing**

Standard errors can also be used to perform hypothesis testing. The most common of which is the *null hypothesis*. We will test to see if the data provides evidence to reject the null hypothesis (that two variables/phenomena/results have no relationship) in favor of the alternative hypothesis (that there is a relationship)

$$\underbrace{H_0}_{\text{Null hypothesis}} : \text{There is no relationship between } X \text{ and } Y \implies \beta_1 = 0 \quad (3.12)$$

$$\underbrace{H_1}_{\text{Alternative hypothesis}} : \text{There is a relationship between } X \text{ and } Y \implies \beta_1 \neq 0 \quad (3.13)$$

We compute the **t-statistic**:

$$t = \frac{\hat{\beta}_1 - 0}{\text{SE}(\hat{\beta}_1)} \quad (3.14)$$

If  $\beta_1 = 0$ , this will have a  $t$ -distribution with  $n - 2$  degrees of freedom. The probability of observing a value greater than or equal to  $|t|$  is called a **p-value**. If the p-value is small, there it is unlikely that the relationship between predictor and response is due to chance. We **reject** the null hypothesis (and accept the alternative hypothesis, if the p-value is "small enough")

**3.1.3 Assessing the Accuracy of the Model****Residual Standard Error**

Suppose we knew the exact true model,  $Y \approx \beta_0 + \beta_1 X$  (from equation 3.1), recall that there is the irreducible error,  $\epsilon$ , associated with every term. Residual Standard Error (RSE) will attempt to estimate the standard deviation of that irreducible error.

It does this by seeing how well your model fits the data! Convenient, let's calculate the standard deviation of our residuals!

$$\text{RSE} = \sqrt{\frac{1}{n-2} \text{RSS}} = \sqrt{\frac{1}{n-2} \sum_{i=1}^n (y_i - \hat{y}_i)^2} \quad (3.15)$$

RSE tells us about the spread of the observed values from the predicted ones, how far off the model's predictions are, on average.. A lower RSE means the predicted values are closer to the observed ones.

RSE is in whatever units the  $Y$  variable is in, so it can be hard to understand what it really means. For that, let's look at...

### $R^2$ Statistic

$R^2$  also measures the accuracy of the model, but does it as a proportion so the values are always between 0 and 1.

$$R^2 = \frac{\text{TSS} - \text{RSS}}{\text{TSS}} = 1 - \frac{\text{RSS}}{\text{TSS}} \quad (3.17)$$

where TSS, the **total sum of squares** is defined

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

- TSS is the total variance in the response,  $Y$ . Think of it as the amount of variability in the response (not dependent on the model)
- RSS measures the amount of variability that is left unexplained after performing the regression
- TSS - RSS then measures the amount of variability in the response that is explained by performing the regression.
- $R^2$  measures the proportion of variability in  $Y$  that can be explained using  $X$

An  $R^2$  value of 1 means the model perfectly explains all variability. An  $R^2$  of 0 means the model explains none of it. It does just as well as predicting the mean.

**Should I add correlation?**

## 3.2 Multiple Linear Regression

What happens if we have more than one dependent variable? If we have  $p$  input variables, a multiple linear regression model takes the form:

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \cdots + \beta_p X_p + \epsilon \quad (3.19)$$

$X_j$  is the  $j$ th predictor,  $\beta_j$  is the corresponding coefficient. It can be interpreted as the average effect on  $Y$  of a one unit increase in  $X_j$ , given that no other predictor variables change.

\*3.2.1 Estimating the Regression Coefficients Our multiple linear regression model will be

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x_1 + \hat{\beta}_2 x_2 + \cdots + \hat{\beta}_p x_p. \quad (3.21)$$

Like in simple linear regression, we will choose our  $\beta$ s in order to minimize the residual sum of squares.

$$\begin{aligned} \text{RSS} &= \sum_{i=1}^n (y_i - \hat{y}_i)^2 \\ &= \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_{i1} - \hat{\beta}_2 x_{i2} - \cdots - \hat{\beta}_p x_{ip})^2 \end{aligned}$$

You'll never compute the coefficients by hand, so we omit the actual calculation. It's fine to know that they are chosen to minimize RSS.

**Note**

While not covered by ISLP, in order to make the book more approachable (and not dependent on linear algebra), you will often see multiple linear regression represented in matrix form. If we let

$$\beta = \begin{bmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \\ \hat{\beta}_2 \\ \vdots \\ \hat{\beta}_n \end{bmatrix} \quad \text{and} \quad \mathbf{X}_i = [1 \quad x_{i1} \quad x_{i2} \quad \cdots \quad x_{in}]$$

then

$$\hat{y}_i = \mathbf{X}_i \beta$$

Note that you will likely see this represented as multiple rows in  $\mathbf{X}$  and  $\mathbf{Y}$  representing different observations. I've omitted that to make it easier to understand.

**Note: Interpreting Regression Coefficients**

- Ideally, want predictors to be uncorrelated (though rarely happens)
  - Lets each coefficient be estimated and tested separately
  - Can make interpretations such as "a unit change in  $X_j$  is associated with a  $\beta_j$  change in  $Y$ , while all the other variables stay fixed"
- Correlations amongst predictors cause problems
  - Variance of all coefficients tends to increase
  - Interpretations becomes difficult, when  $X_j$  changes, so too do others
    - \* Suppose a company has a fixed advertising budget, if you increase TV marketing, you must decrease internet marketing.
- Claims of causality should be avoided for observational data
- It's possible (and likely) that a predictor that might have a big effect when the sole predictor has little to no effect, if it's correlated with another predictor.

**3.2.2 Some Important Questions**

1. Is at least one of the predictors  $X_1, X_2, \dots, X_p$  useful in predicting the response?
2. Do all the predictors help to explain  $Y$ , or is only a subset of the predictors useful?
3. How well does the model fit the data?
4. Given a set of predictor values, what response value should we predict, and how accurate is our prediction?

**Is there a Relationship between the Response and Predictors?**

We use the F-Statistic

$$F = \frac{\text{TSS} - \text{RSS}/p}{\text{RSS}/(n - p - 1)} \quad (3.23)$$

- Tests that all coefficients are 0.
- When there is no relationship between the response and predictors, the  $F$ -statistics takes on a value close to 1. If there is a relationship (and  $H_a$  is true), the  $F$  statistic is greater than 1.
- How large the  $F$ -statistic must be depends on the values of  $n$  and  $p$ .
- For large  $n$ , smaller  $F$ -statistic can provide evidence against  $H_0$ . If  $n$  is small, need a larger  $F$ -statistic.

If the null hypothesis is that a specific subset of coefficients are zero (for convenience, the ones testing are at the end of the list):

$$H_0 : \beta_{p-q+1} = \beta_{p-q+2} = \cdots = \beta_p = 0$$

we fit a second model that uses all the variables **except** the subset we are testing. Call  $RSS_0$  the residual sum of squares of this second model. then the  $F$ -statistic is

$$F = \frac{RSS_0 - RSS/q}{RSS/(n - p - 1)} \quad (3.24)$$

- The  $p$ -values of individual regressors (which provided information on if that variable was related to the response), is equivalent to the  $F$ -statistic omitting that single variable from the model, leaving the rest in.

#### Warning: Why bother with F statistics?

Suppose we have 100 predictors  $p = 100$  and none are related to the response, meaning  $H_0 : \beta_1 = \beta_2 = \cdots = \beta_p = 0$  is true. There is a 5% chance that any given  $p$ -value is below 0.05, totally by chance! So we would expect 5 of them to have small  $p$ -values, despite there being no true association between the predictor and response! This means, **If we use the individual  $t$ -statistics and associated  $p$ -values in order to decide whether or not there is any association between the variables and the response, there is a very high chance that we will incorrectly conclude that there is a relationship.** The  $F$ -statistic does not suffer this problem because it adjusts for the number of predictors. If  $H_0$  is true (and there is no association between any predictor and the response), there is only a 5% chance that the  $F$ -statistic will result in a  $p$  value below 0.05.

- Using an  $F$ -statistic to test for any association works when  $p$  is relatively small compared to  $n$ .
  - If  $p > n$ , then there are more coefficients  $\beta_j$  to estimate than observations from which to estimate. Then, we cannot use least-squares to fit the multiple linear regression, so the  $F$ -statistic cannot be used
  - When  $p$  is large, some approaches discussed soon, like forward selection can be used. More detail about handling high-dimensional settings is in Chapter 6.

### Deciding on Important Variables

- Most direct approach is called **all subsets** or **best subsets** regression. Compute the least squares fit for all possible subsets and then choose between them based on some criterion that balances training error with model size.
- Can't do when large number of regressors, since there are  $2^p$  possible subsets. Even  $p = 40$  results in over a billion models!

### Three Main Approaches

- **Forward Selection** - Begin with **null model** (has intercept, but no predictors). Then fit  $p$  simple linear regressions and add to the null model, the variable that results in the lowest RSS. Repeat for the new one-variable, two-variable, etc model.
  - Greedy approach - might include variables early on that later become redundant
- **Backward Selection** - start with all variables in model, remove the variable with largest  $p$  value. Repeat process
  - Cannot be used if  $p > n$

- **Mixed Selection** - Start with null model, add variables one-by-one. If as add new variables, the  $p$ -value for a model reaches a threshold, remove it. Continue going forward and backwards until all variables have low  $p$  value and all variables not in the model have large  $p$ -value.

### Model Fit

- **RSE**
- $R^2$  - Recall, in simple linear regression, is correlation of response and variable. In multiple,  $\text{Cor}(Y, \hat{Y})^2$ 
  - Note: Property of fitted linear model is maximizes the correlation among all possible linear models
  - $R^2$  close to 1 means model explains a large portion of variance in the response variable.
  - $R^2$  always increases as more variables are added to the model, even if those variables only weakly associated with response, since get better fit to training data. Small increase probably means omit variable.

enditemize

### Predictions

## 3.3 Other Considerations in the Regression Model

## 3.4 The Marketing Plan

## 3.5 Comparison of Linear Regression with K-Nearest Neighbors