

STAT 331 - Applied Linear Models

Cameron Roopnarine

Last updated: November 2, 2020

1 Introduction to Regression Models

LECTURE 1 | 2020-09-08

DEFINITION 1.1: Response variable

A **response (dependent) variable** is the primary variable of interest, denoted by a capital roman letter Y .

DEFINITION 1.2: Explanatory Variable

An **explanatory (independent, predictor) variable** are variables that impact the response, denoted by x_i for $i = 1, \dots, p$.

DEFINITION 1.3: Regression Model

A **regression model** deals with modelling the functional relationship between a response variable and one or more explanatory variables.

EXAMPLE 1.4: Alligators in Florida

Let Y be the length in metres of an alligator and $x_1 := \{0, 1\}$ (male or female). The mass in an alligators stomach consists of fish (x_2), invertebrates (x_3), reptiles (x_4), birds (x_5), and other (x_6, \dots, x_p). We imagine we can explain Y in terms of (x_1, \dots, x_p) using some function such that $Y = f(x_1, \dots, x_p)$.

In this course, we will be looking at linear models.

DEFINITION 1.5: Linear model

A general **linear model** is defined as $Y = \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p + \varepsilon$ where Y is the response variable, (x_1, \dots, x_p) are the p explanatory variables, $(\beta_0, \beta_1, \dots, \beta_p)$ are the model parameters, and ε is the random error. We assume that (x_1, \dots, x_p) are fixed constants, β_0 is the intercept of Y , $(\beta_1, \dots, \beta_p)$ all quantify effect on x_j on Y , and $\varepsilon \sim N(0, \sigma^2)$.

REMARK 1.6

In general, the model will not perfectly explain the data.
“All models are wrong, but some are useful.”

$Y \sim N(\beta_0 + \beta_1 x_1 + \dots + \beta_p x_p, \sigma^2)$ since $\mathbb{E}[Y] = \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p$ and $\mathbb{V}(Y) = \mathbb{V}(\varepsilon) = \sigma^2$.

2 Simple Linear Regression

DEFINITION 2.1: Simple linear regression

A **simple linear regression** is a linear model that uses only one explanatory variable; that is, $Y = \beta_0 + \beta_1 x + \varepsilon$. The **data** in a simple linear regression consists of pairs (x_i, y_i) where $i = 1, \dots, n$.

REMARK 2.2

Before fitting any model, we might want to make a scatterplot to visualize if there is a linear relationship between x and y , or calculate the *correlation*.

DEFINITION 2.3: Correlation

The **correlation** of random variables X and Y is $\rho_{XY} = \frac{\text{Cov}(X, Y)}{\text{Sd}(X)\text{Sd}(Y)}$.

DEFINITION 2.4: Sample correlation

The **sample correlation** of all pairs (x_i, y_i) is

$$\begin{aligned} r &= \frac{\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2} \sqrt{\frac{1}{n-1} \sum_{i=1}^n (y_i - \bar{y})^2}} \\ &= \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}} \\ &= \frac{S_{xy}}{\sqrt{S_{xx} S_{yy}}} \end{aligned}$$

REMARK 2.5

The sample correlation measures the strength and direction of the linear relationship between X and Y . Note that $-1 \leq r \leq 1$. If $|r| \approx 1$, then there is a strong linear relationship, and if $|r| \approx 0$ then there is a lack of linear relationship. Also, if $r > 0$, then there is a positive relationship, and if $r < 0$ then there is a negative relationship. It does not tell us how to predict Y from X . To do so, we need to estimate β_0 and β_1 .

DEFINITION 2.6: Simple linear regression model

For data (x_i, y_i) for $i = 1, \dots, n$, the **simple linear regression model** is $Y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$ with the assumption that $\varepsilon_i \stackrel{\text{iid}}{\sim} N(0, \sigma^2)$. Therefore, $Y_i \sim N(\mu_i = \beta_0 + \beta_1 x_i, \sigma^2)$.

DEFINITION 2.7: Method of least squares

The method of estimating β_0 and β_1 by minimizing $S(\beta_0, \beta_1) = \sum_{i=1}^n (y_i - (\beta_0 + \beta_1 x_i))^2$ is referred to as the **method of least squares**.

REMARK 2.8

The least squares is equivalent to maximum likelihood estimate when $\varepsilon_i \stackrel{\text{iid}}{\sim} N(0, \sigma^2)$.

THEOREM 2.9: Least Square Estimates (LSEs) for SLR

Minimizing $S(\beta_0, \beta_1)$, gives the least square estimates

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} \quad \text{and} \quad \hat{\beta}_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} = \frac{S_{xy}}{S_{xx}}$$

Proof of: 2.9

$$\frac{\partial S}{\partial \beta_0} = 2 \sum_{i=1}^n [y_i - (\beta_0 + \beta_1 x_i)](-1) \quad \text{and} \quad \frac{\partial S}{\partial \beta_1} = 2 \sum_{i=1}^n [y_i - (\beta_0 + \beta_1 x_i)](-x_i).$$

Now,

$$\frac{dS}{d\beta_0} := 0 \iff \sum_{i=1}^n y_i - n\beta_0 - \beta_1 \sum_{i=1}^n x_i = 0 \iff \beta_0 = \bar{y} - \beta_1 \bar{x}$$

$$\begin{aligned} \frac{dS}{d\beta_1} &:= 0 \stackrel{\text{plug } \beta_0}{\iff} \sum_{i=1}^n [y_i - \bar{y} + \beta_1 \bar{x} - \beta_1 x_i] x_i = 0 \\ &\iff \sum_{i=1}^n x_i (y_i - \bar{y}) - \beta_1 \sum_{i=1}^n x_i (x_i - \bar{x}) = 0 \\ &\iff \beta_1 = \frac{\sum_{i=1}^n x_i (y_i - \bar{y})}{\sum_{i=1}^n x_i (x_i - \bar{x})} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} = \frac{S_{xy}}{S_{xx}} \end{aligned}$$

REMARK 2.10

We use a hat on the β 's to show that they are estimates.

DEFINITION 2.11: Fitted value, Residual

The expression $\hat{\mu}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i$ is called the **fitted value** that corresponds to the i th observation with x_i as the explanatory variable. The difference between y_i and $\hat{\mu}_i$, and $e_i = y_i - \hat{\mu}_i$ is referred to as the **residual**. It is the vertical distance between the observation y_i and the estimated line $\hat{\mu}_i$ evaluated at x_i .

LECTURE 3 | 2020-09-14

For $Y_i \sim N(\beta_0 + \beta_1 x_i, \sigma^2)$, the equation of fitted line is given by $y = \hat{\beta}_0 + \hat{\beta}_1 x$. Our interpretation of the parameters is as follows.

- $\hat{\beta}_0$ is the estimate of the expected response when $x = 0$ (but not always meaningful if outside range of x_i 's in data)
- $\hat{\beta}_1$ is the estimate of expected change in response for unit increase in x
- σ^2 is the “variability around the line” where $\sigma^2 = \mathbb{V}(\varepsilon_i) = \mathbb{V}(Y_i)$

Q: How should we estimate σ^2 ?

$$\varepsilon_i = Y_i - (\beta_0 + \beta_1 x_i) \quad \text{and} \quad e_i = y_i - (\hat{\beta}_0 + \hat{\beta}_1 x_i)$$

Our intuition tells us to use variability in the residuals to estimate σ^2 , so we use

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

where the first term looks like sample variance of e_i 's. The second equality follows since $\bar{e} = \bar{y} - (\hat{\beta}_0 + \hat{\beta}_1 \bar{x}) = 0$ by definition of our $\hat{\beta}_0$ estimate.

DEFINITION 2.12: Residual sum of squares

$SS(\text{Res}) = \sum_{i=1}^n (y_i - \hat{\mu}_i)^2 = \sum_{i=1}^n e_i^2$, is known as the **residual (error) sum of squares**.

REMARK 2.13

The $n-2$ will be looked at in more detail later, but for now it suffices to say that the degrees of freedom is $n-2$ or equivalently, n - number of parameters estimated. It allows $\hat{\sigma}^2$ to be an unbiased estimator for the true value of σ^2 ; that is, $\mathbb{E}[\hat{\sigma}^2] = \sigma^2$ whenever $\hat{\sigma}^2$ is viewed as a random variable.

THEOREM 2.14: Linear Combination of Independent Normal Random Variables

If $Y_i \sim N(\mu_i, \sigma^2)$, $i = 1, \dots, n$ independently, then

$$\sum_{i=1}^n a_i Y_i \sim N\left(\sum_{i=1}^n a_i \mu_i, \sum_{i=1}^n a_i^2 \sigma_i^2\right)$$

Proof of: 2.14

The proof is completed in STAT 330 with moment generating functions.

Viewing $\hat{\beta}_1$ as a random variable:

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(Y_i - \bar{Y})}{\sum_{i=1}^n (x_i - \bar{x})^2} = \frac{\sum_{i=1}^n (x_i - \bar{x})Y_i - \bar{Y} \overbrace{\sum_{i=1}^n (x_i - \bar{x})}^0}{\sum_{i=1}^n (x_i - \bar{x})x_i - \bar{x} \underbrace{\sum_{i=1}^n (x_i - \bar{x})}_0} = \frac{\sum_{i=1}^n (x_i - \bar{x})Y_i}{\sum_{i=1}^n (x_i - \bar{x})x_i} = \sum_{i=1}^n a_i Y_i$$

where $a_i = \frac{x_i - \bar{x}}{\sum_{i=1}^n x_i(x_i - \bar{x})}$. Therefore,

$$\mathbb{E}[\hat{\beta}_1] = \sum_{i=1}^n a_i \mathbb{E}[Y_i] = \frac{\sum_{i=1}^n (x_i - \bar{x})(\beta_0 + \beta_1 x_i)}{\sum_{i=1}^n x_i(x_i - \bar{x})} = \frac{\beta_0 \overbrace{\sum_{i=1}^n (x_i - \bar{x})}^0 + \beta_1 \sum_{i=1}^n x_i(x_i - \bar{x})}{\sum_{i=1}^n x_i(x_i - \bar{x})} = \beta_1$$

Now, we calculate the variance of $\hat{\beta}_1$:

$$\mathbb{V}(\hat{\beta}_1) = \sum_{i=1}^n a_i^2 \mathbb{V}(Y_i) = \sigma^2 \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{[\sum_{i=1}^n x_i(x_i - \bar{x})]^2} = \sigma^2 \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{[\sum_{i=1}^n (x_i - \bar{x})^2]^2} = \frac{\sigma^2}{S_{xx}}$$

Using our calculations from $\hat{\beta}_1$, and viewing $\hat{\beta}_0$ as a random variable:

$$\mathbb{E}[\hat{\beta}_0] = \mathbb{E}[\bar{Y}] - \bar{x} \mathbb{E}[\hat{\beta}_1] = \mathbb{E}\left[\frac{\sum_{i=1}^n Y_i}{n}\right] - \bar{x} \beta_1 = \frac{\sum_{i=1}^n (\beta_0 + \beta_1 x_i)}{n} - \beta_1 \bar{x} = \beta_0 + \beta_1 \bar{x} - \beta_1 \bar{x} = \beta_0$$

Now, we calculate the variance of $\hat{\beta}_0$:

$$\mathbb{V}(\hat{\beta}_1) = \mathbb{V}(\bar{Y} - \beta_1 \bar{x}) = \mathbb{V}(\bar{Y}) + (-\bar{x})^2 \mathbb{V}(\beta_1) = \mathbb{V}\left(\frac{\sum_{i=1}^n Y_i}{n}\right) + \bar{x}^2 \left(\frac{\sigma^2}{S_{xx}}\right) = \frac{n\sigma^2}{n^2} + \frac{\sigma^2 \bar{x}^2}{S_{xx}}$$

Also, since $\hat{\beta}_1$ and $\hat{\beta}_0$ are linear combination of Normal random variables, they follow a Normal distribution. Therefore, we get the following theorem.

THEOREM 2.15: Distribution of LSEs

The distribution of the least square estimates are given by

$$\hat{\beta}_1 \sim N\left(\beta_1, \frac{\sigma^2}{S_{xx}}\right) \quad \text{and} \quad \hat{\beta}_0 \sim N\left(\beta_0, \sigma^2 \left(\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}}\right)\right)$$

Since $\mathbb{E}[\hat{\beta}_1] = \beta_1$, we say $\hat{\beta}_1$ is an unbiased estimator of β_1 . This implies that when the experiment is repeated a large number of times, the average of the estimates $\hat{\beta}_1$; that is, $\mathbb{E}[\hat{\beta}_1]$ coincides with the true value of β_1 . A similar argument can be made for β_0 .

Then, $\frac{\hat{\beta}_1 - \beta_1}{\sigma/\sqrt{S_{xx}}} \sim N(0, 1)$, but σ is unknown, so need to use $\hat{\sigma}$ to get $\frac{\hat{\beta}_1 - \beta_1}{\hat{\sigma}/\sqrt{S_{xx}}} \sim t(n-2)$.

DEFINITION 2.16: Standard deviation and standard error of $\hat{\beta}_1$

The **standard deviation** of $\hat{\beta}_1$ is defined as $\text{Sd}(\hat{\beta}_1) = \sigma/\sqrt{S_{xx}}$. The **estimated** standard deviation of $\hat{\beta}_1$ is also referred to as the **standard error** of the estimate $\hat{\beta}_1$, and we write $\text{Se}(\hat{\beta}_1) = \hat{\sigma}/\sqrt{S_{xx}}$.

DEFINITION 2.17: Student t distribution

Suppose $Z \sim N(0, 1)$ and $U \sim \chi^2(\nu)$, with Z and U independent. Then, $T = Z/\sqrt{U/\nu}$ has a **Student t distribution** with ν degrees of freedom.

THEOREM 2.18

For a simple linear regression model,

$$\frac{\hat{\sigma}^2(n-2)}{\sigma^2} = \frac{SS(\text{Res})}{\sigma^2} \sim \chi^2(n-2)$$

Proof of: 2.18

Too hard for sure.

Using the theorem stated, we justify the fact that replacing σ with $\hat{\sigma}$ gives us a $t(n-2)$ distribution.

$$\frac{\hat{\beta}_1 - \beta_1}{\hat{\sigma}/\sqrt{S_{xx}}} = \frac{\frac{\hat{\beta}_1 - \beta_1}{\sigma/\sqrt{S_{xx}}}}{\sqrt{\frac{\hat{\sigma}^2(n-2)}{\sigma^2} \left(\frac{1}{n-2}\right)}} = \frac{Z}{\sqrt{U/\nu}} = T \sim t(n-2)$$

where $\frac{\hat{\sigma}^2(n-2)}{\sigma^2} = U$, $\nu = n-2$, and $Z = \frac{\hat{\beta}_1 - \beta_1}{\hat{\sigma}/\sqrt{S_{xx}}}$. A $(1-\alpha)$ confidence interval for β_1 is

$$\hat{\beta}_1 \pm c \text{Se}(\hat{\beta}_1)$$

where c is the $1 - \frac{\alpha}{2}$ quantile of $t(n-2)$; that is, $P(|T| \leq c) = 1 - \alpha$ or $P(T \leq c) = 1 - \frac{\alpha}{2}$ where $T \sim t(n-2)$.

Hypothesis test: $H_0: \beta = 0$ versus $H_A: \beta_1 \neq 0$. If H_0 is true, then $\hat{\beta}_1/\text{Se}(\hat{\beta}_1) \sim t(n-2)$, so calculate the **t statistic** $t = \hat{\beta}_1/\text{Se}(\hat{\beta}_1)$, and reject H_0 at level α if $|t| > c$ where c is $1 - \frac{\alpha}{2}$ quantile of $t(n-2)$. Therefore, $p\text{-value} = P(|T| \geq |t|) = 2P(T \geq |t|)$.

LECTURE 4 | 2020-09-16

Suppose we want to predict the response y for a new value of x , say $x = x_0$. Then, SLR model says $Y_0 \sim N(\beta_0 + \beta_1 x_0, \sigma^2)$ where Y_0 is a random variable for response when $x = x_0$; that is, $\hat{Y}_0 = \hat{\beta}_0 + \hat{\beta}_1 x_0$. The fitted model predicts the value of y to be $\hat{y}_0 = \hat{\beta}_0 + \hat{\beta}_1 x_0$.

Also, $\mathbb{E}[\hat{Y}_0] = \mathbb{E}[\hat{\beta}_0] + x_0 \mathbb{E}[\hat{\beta}_1] = \beta_0 + \beta_1 x_0 = \mathbb{E}[Y_0]$, since $\hat{\beta}_i$ for $i = 0, 1$ are unbiased. Therefore, we can say that \hat{Y}_0 is an unbiased estimate of the random variable for the mean of Y_0 . For the variance of \hat{Y}_0 we write

$$\begin{aligned} \hat{Y}_0 &= \hat{\beta}_0 + \hat{\beta}_1 x_0 = \bar{Y} - \hat{\beta}_1 \bar{x} + \hat{\beta}_1 x_0 \\ &= \bar{Y} + \hat{\beta}_1 (x_0 - \bar{x}) \\ &= \sum_{i=1}^n \left[\frac{Y_i}{n} + (x_0 - \bar{x}) \left(\frac{(x_i - \bar{x})(Y_i - \bar{Y})}{S_{xx}} \right) \right] \\ &= \sum_{i=1}^n \left[\frac{Y_i}{n} + (x_0 - \bar{x}) \left(\frac{(x_i - \bar{x})Y_i}{S_{xx}} \right) \right] \\ &= \sum_{i=1}^n \left[\frac{1}{n} + \frac{(x_0 - \bar{x})(x_i - \bar{x})}{S_{xx}} \right] Y_i \\ &= \sum_{i=1}^n a_i Y_i \end{aligned}$$

where $a_i = \frac{1}{n} + \frac{(x_0 - \bar{x})(x_i - \bar{x})}{S_{xx}}$. Therefore,

$$\begin{aligned} \mathbb{V}(Y_0) &= \sum_{i=1}^n \left[\frac{1}{n} + \frac{(x_0 - \bar{x})(x_i - \bar{x})}{S_{xx}} \right]^2 \\ &= \sum_{i=1}^n \left[\frac{1}{n^2} + \frac{2(x_0 - \bar{x})(x_i - \bar{x})}{nS_{xx}} + \frac{(x_0 - \bar{x})^2(x_i - \bar{x})^2}{(S_{xx})^2} \right] \\ &= \sum_{i=1}^n \left[\frac{1}{n^2} \right] + \frac{2(x_0 - \bar{x})}{nS_{xx}} \sum_{i=1}^n (x_i - \bar{x}) + \frac{(x_0 - \bar{x})^2}{(S_{xx})^2} \sum_{i=1}^n (x_i - \bar{x})^2 \\ &= \frac{1}{n} + \frac{2(x_0 - \bar{x})}{S_{xx}} (0) + \frac{(x_0 - \bar{x})^2}{(S_{xx})^2} (S_{xx}) \\ &= \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{S_{xx}} \end{aligned}$$

We proved the following theorem.

THEOREM 2.19: Distribution of Prediction

The distribution of the prediction random variable is given by

$$\hat{Y}_0 \sim N \left(\beta_0 + \beta_1 x_0, \sigma^2 \left(\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{S_{xx}} \right) \right)$$

DEFINITION 2.20: Prediction error

The random variable for **prediction error** is defined as $Y_0 - \hat{Y}_0$ where Y_0 and \hat{Y}_0 are independent and \hat{Y}_0 is a function of Y_1, \dots, Y_n .

$$\begin{aligned} \mathbb{E}[Y_0 - \hat{Y}_0] &= \mathbb{E}[Y_0] - \mathbb{E}[\hat{Y}_0] = 0 \\ \mathbb{V}(Y_0 - \hat{Y}_0) &= \mathbb{V}(Y_0) + (-1)^2 \mathbb{V}(\hat{Y}_0) = \sigma^2 + \sigma^2 \left(\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{S_{xx}} \right) \end{aligned}$$

We proved the following theorem.

THEOREM 2.21: Distribution of Prediction Error

The distribution of the prediction error is given by

$$Y_0 - \hat{Y}_0 \sim N \left(0, \sigma^2 \left(1 + \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{S_{xx}} \right) \right)$$

Since σ is unknown, we use $\hat{\sigma}$ and get the following:

$$\frac{Y_0 - \hat{Y}_0}{\hat{\sigma} \sqrt{1 + \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{S_{xx}}}} \sim t(n-2)$$

Intuition for prediction error composed of 2 terms:

- $\mathbb{V}(Y_0)$: random error of new observation
- $\mathbb{V}(\hat{Y}_0)$ (predictor): estimating β_0 and β_1

Those are 2 sources of uncertainty.

REMARK 2.22

Be careful that the prediction may not make sense if x_0 is outside the range of the x_i 's in the data.

A $(1 - \alpha)$ prediction interval for the mean response $y_0 = \beta_0 + \beta_1 x_0$ at x_0 is

$$\hat{y}_0 \pm c \hat{\sigma} \sqrt{1 + \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{S_{xx}}}$$

where c is the $1 - \frac{\alpha}{2}$ quantile of $t(n-2)$.

EXAMPLE 2.23: Orange production 2018 in FL

We are given the following information.

- x : acres
- y : # boxes of oranges (thousands)
- (x_i, y_i) recorded for each of 25 FL counties
- $r = 0.964$
- $\bar{x} = 16133$
- $\bar{y} = 1798$
- $S_{xx} = 1.245 \times 10^{10}$
- $S_{xy} = 1.453 \times 10^9$

Now, $\hat{\beta}_1 = S_{xy}/S_{xx} = 0.1167$ has a positive slope, therefore x and y are positively correlated. The expected number of boxes produced is estimated to be about 117 higher per an additional acre.

Computing $\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} = -85.3$, we see that it is not meaningful to interpret, since it is the expected production if there were 0 acres (outside the range of x_i) as no county has $x = 0$.

Now suppose $SS(\text{Res}) = 1.31 \times 10^7$ the residuals are the differences between y_i and the fitted regression line.

- $\hat{\sigma}^2 = \frac{\sum_{i=1}^n e_i^2}{n-2} = \frac{1.31 \times 10^7}{25-2} = 5.7 \times 10^5$
- $Se(\hat{\beta}_1) = \frac{\hat{\sigma}}{\sqrt{S_{xx}}} = 0.00676$
- To test $H_0: \beta_1 = 0$, calculate $t = (\hat{\beta}_1 - 0)/Se(\hat{\beta}_1) = 0.1167/0.00676 \approx 17.3$, then elect the 0.975 quantile (for demonstration purposes) of $t(23)$ which is 2.07.
- Note that 17.3 is very unlikely to see in $t(23)$.

Since $17.3 \gg 2.07$, we reject H_0 at $\alpha = 0.05$ level, and conclude there's a significant linear relationship between acres and oranges produced.

The 95% confidence interval for β_1 is given by $0.1167 \pm 2.07(0.00676)$, which does not contain 0.

$$p\text{-value} = P(|t_{23}| \geq 17.3) = 2P(t_{23} \geq 17.3) \approx 1.2 \times 10^{-14}$$

Predict the # of boxes in thousands produced if we had 10000 acres to grow oranges.

$$\hat{\beta}_0 + \hat{\beta}_1 x_0 = -85.3 + (0.1167)(10000) \approx 1082$$

The 95% prediction interval is given by

$$1082 \pm 2.07 \sqrt{5.69 \times 10^5} \sqrt{1 + \frac{1}{25} + \frac{(6133)^2}{1.245 \times 10^{10}}} = [-512.0407, 2675.595]$$

REMARK 2.24

We are **not** trying to establish causation.

The example done in R is included in the next page.


```

# Read data from florange.csv and input it into the dat vector.
dat <- read.csv("florange.csv")
# Done to make the predict function work well.
x <- dat$acres
y <- dat$boxes
# Output the first 6 rows in dat.
head(dat)

```

```

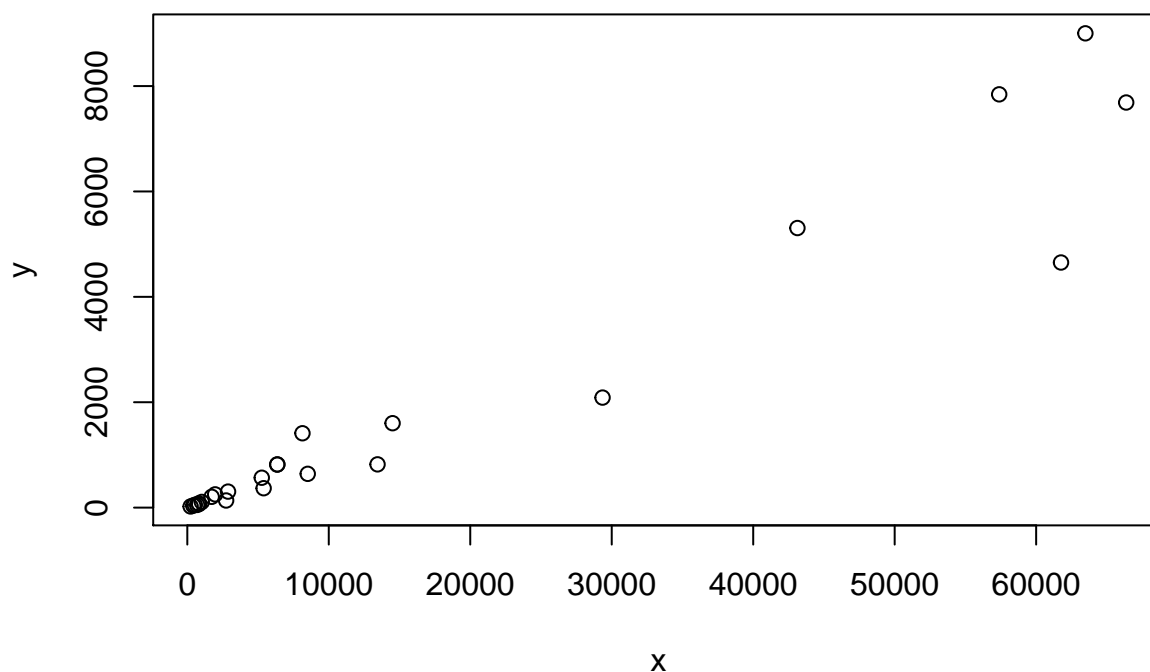
##      county boxes acres
## 1  Brevard    51   696
## 2 Charlotte  821 13447
## 3  Collier  2088 29351
## 4   DeSoto  7688 66365
## 5   Glades   368  5396
## 6   Hardee  5306 43126

```

```

# Draw a scatterplot with x-axis as `acres` and y-axis as `boxes`.
plot(x,y)

```



```

# Compute some common variables with common functions.
r <- cor(x,y)
xbar <- mean(x)
ybar <- mean(y)
cat("r:", r, "xbar:", xbar, "ybar:", ybar)

```

```
## r: 0.9635098 xbar: 16132.64 ybar: 1797.56
```

Therefore, $r = 0.9635098$, $\bar{x} = 16132.64$, and $\bar{y} = 1797.56$.

```

# Compute some common variables manually.
Sxx <- sum( (x - xbar)^2 )
Sxy <- sum( (x - xbar) * (y - ybar) )
cat("Sxx: ", Sxx, "Sxy: ", Sxy)

```

```
## Sxx: 12450023404 Sxy: 1453128337
```

Therefore, $S_{xx} = 12450023404 = 1.245 \times 10^{10}$ and $S_{xy} = 1453128337 = 1.453 \times 10^9$.

```
# R's lm function fits linear models
```

```
lm.1 <- lm(y~x)
```

```
summary(lm.1)
```

```
##
```

```
## Call:
```

```
## lm(formula = y ~ x)
```

```
##
```

```
## Residuals:
```

```
##      Min       1Q   Median       3Q      Max
## -2470.81    -6.17    71.72   106.46  1677.32
```

```
##
```

```
## Coefficients:
```

```
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept) -85.391989  186.178031  -0.459    0.651
## x              0.116717   0.006761  17.263 1.16e-14 ***
```

```
## ---
```

```
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
##
```

```
## Residual standard error: 754.4 on 23 degrees of freedom
```

```
## Multiple R-squared:  0.9284, Adjusted R-squared:  0.9252
```

```
## F-statistic:    298 on 1 and 23 DF,  p-value: 1.164e-14
```

From the summary, we can see that $\hat{\beta}_0 = -85.391989$, $\hat{\beta}_1 = 0.116717$, $\text{Se}(\hat{\beta}_1) = 0.006761$, $t = 17.263$, $p\text{-value} = 1.64 \times 10^{-14}$, and $\hat{\sigma} = 754.4$.

```
# Sum Squared Fitted Values
```

```
sum(lm.1$fitted.values^2)
```

```
## [1] 250385207
```

```
# Sum Squared Residuals
```

```
sum(lm.1$residuals^2)
```

```
## [1] 13089860
```

Therefore, $SS(\text{Res}) = \sum_{i=1}^n e_i^2 = 13089860 = 1.31 \times 10^7$.

```
# Manual calculation of sigma^2 estimate
```

```
sum(lm.1$residuals^2) / 23
```

```
## [1] 569124.3
```

Therefore, $\hat{\sigma}^2 = 569124.3 = 5.7 \times 10^5$.

```
# Manual calculation of sigma estimate
```

```
sqrt(sum(lm.1$residuals^2) / 23)
```

```
## [1] 754.4033
```

Therefore, $\hat{\sigma} = 754.4$.

```
# t distribution values
```

```
qt(0.975,23)
```

```
## [1] 2.068658
```

Therefore, $c = 2.07$.

```
# 95% confidence interval
```

```
confint(lm.1)
```

```
##              2.5 %      97.5 %  
## (Intercept) -470.5305905 299.7466119  
## x           0.1027305   0.1307034
```

```
# 95% prediction interval with predicted boxes if we had 10000 acres
```

```
predict(lm.1, data.frame(x=10000), interval="prediction")
```

```
##      fit      lwr      upr  
## 1 1081.777 -512.0407 2675.595
```

Q: Is σ the same for all values of y ?

A: It appears to not in the sense that the variance appears to be higher with respect to higher acres. Sigma will be smaller when there's less acres. Later, this will be testing equal variance or homoscedastic assumption. Later, when we talk about variable transformations we can consider taking the logarithm.

Q: Are the error terms plausibly independent? In other words, does knowing one e_i (residual) help predict e_j (another residual) for a different county?

A: There's diagnostics for checking this. However, intuitively there could be some common factors at play when two counties are geographically close.

3 Multiple Linear Regression

DEFINITION 3.1: Multiple linear regression

A **multiple linear regression** (MLR) model is defined as

$$Y = \beta_0 + \beta_1 x_1 + \cdots + \beta_p x_p + \varepsilon$$

which links a response variable y to several independent explanatory variables x_1, x_2, \dots, x_p .

EXAMPLE 3.2: Rocket MLR

- x_1 : nozzle area (large or small, 0 or 1)
- x_2 : mixture in propellant, ratio oxidized fuel
- Y : thrust

Want to develop linear relationship between response y and x_1, x_2 ; that is, we want to develop a linear relationship between thrust and both nozzle area and mixture in propellant.

In a MLR, there are n observations, where each consists of p response variables (y_i), and p explanatory variables ($x_{i1}, x_{i2}, \dots, x_{ip}$). Then,

$$Y_i \sim N(\underbrace{\beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \cdots + \beta_p x_{ip}}_{\mathbb{E}[Y_i] = \mu_i}, \sigma^2)$$

or $Y_i = \mu_i + \varepsilon_i$ where $\varepsilon_i \stackrel{\text{iid}}{\sim} N(0, \sigma^2)$. We can write in vector/matrix form

$$\begin{bmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{bmatrix}_{n \times 1} = \begin{bmatrix} \beta_0 + \beta_1 x_{11} + \cdots + \beta_p x_{1p} \\ \beta_0 + \beta_1 x_{21} + \cdots + \beta_p x_{2p} \\ \vdots \\ \beta_0 + \beta_1 x_{n1} + \cdots + \beta_p x_{np} \end{bmatrix} + \begin{bmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_n \end{bmatrix}_{n \times 1}$$

Which we can more commonly write as $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$ where

$$\mathbf{Y} = \begin{bmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{bmatrix}_{n \times 1} \quad \mathbf{X} = \begin{bmatrix} 1 & x_{11} & x_{12} & \cdots & x_{1(p-1)} & x_{1p} \\ 1 & x_{21} & x_{22} & \cdots & x_{2(p-1)} & x_{2p} \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 1 & x_{(n-1)1} & x_{(n-1)2} & \cdots & x_{(n-1)(p-1)} & x_{(n-1)p} \\ 1 & x_{n1} & x_{n2} & \cdots & x_{n(p-1)} & x_{np} \end{bmatrix}_{n \times (p+1)} \quad \boldsymbol{\beta} = \begin{bmatrix} \beta_0 \\ \vdots \\ \beta_p \end{bmatrix}_{(p+1) \times 1} \quad \boldsymbol{\varepsilon} = \begin{bmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_n \end{bmatrix}_{n \times 1}$$

DEFINITION 3.3: Random vector

We call $\mathbf{Y} = (Y_1, Y_2, \dots, Y_n)^\top$ a **random vector**.

DEFINITION 3.4: Mean vector

The **mean vector** of \mathbf{Y} is defined as $\mathbb{E}[\mathbf{Y}] = (\mathbb{E}[Y_1], \mathbb{E}[Y_2], \dots, \mathbb{E}[Y_n])^\top$.

DEFINITION 3.5: Covariance matrix

The **covariance matrix** (or **variance-covariance matrix**) of \mathbf{Y} is defined as

$$\mathbb{V}(\mathbf{Y}) = \begin{bmatrix} \mathbb{V}(Y_1) & \text{Cov}(Y_1, Y_2) & \cdots & \text{Cov}(Y_1, Y_{n-1}) & \text{Cov}(Y_1, Y_n) \\ \text{Cov}(Y_2, Y_1) & \mathbb{V}(Y_2) & \cdots & \text{Cov}(Y_2, Y_{n-1}) & \text{Cov}(Y_2, Y_n) \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ \text{Cov}(Y_{n-1}, Y_1) & \text{Cov}(Y_{n-1}, Y_2) & \cdots & \mathbb{V}(Y_{n-1}) & \text{Cov}(Y_{n-1}, Y_n) \\ \text{Cov}(Y_n, Y_1) & \text{Cov}(Y_n, Y_2) & \cdots & \text{Cov}(Y_n, Y_{n-1}) & \mathbb{V}(Y_n) \end{bmatrix}_{n \times n}$$

PROPOSITION 3.6: Properties of Covariance Matrix

Let \mathbf{Y} be a random vector and $\mathbf{a} \in \mathbb{R}^n$, then the covariance matrix has the following properties.

- (1) Symmetric since $\text{Cov}(Y_i, Y_j) = \text{Cov}(Y_j, Y_i)$; that is $\mathbb{V}(\mathbf{Y})^\top = \mathbb{V}(\mathbf{Y})$.
- (2) Positive semi-definite since $\mathbf{a}^\top \mathbb{V}(\mathbf{Y}) \mathbf{a} \geq 0$ for all $\mathbf{a} \in \mathbb{R}^n$.
- (3) $\mathbb{V}(\mathbf{Y}) = \mathbb{E}[(\mathbf{Y} - \mathbb{E}[\mathbf{Y}]) (\mathbf{Y} - \mathbb{E}[\mathbf{Y}])^\top]$

Proof of: 3.6

Trivial.

PROPOSITION 3.7: Properties of Random Vector

Let \mathbf{a} be a $1 \times n$ matrix (row vector) of constants and A be an $n \times n$ matrix of constants, then the random vector has the following properties.

- (1) $\mathbb{E}[\mathbf{aY}] = \mathbf{aY}$
- (2) $\mathbb{E}[A\mathbf{Y}] = A\mathbb{E}[\mathbf{Y}]$
- (3) $\mathbb{V}(\mathbf{aY}) = \mathbf{a}\mathbb{V}(\mathbf{Y})\mathbf{a}^\top$
- (4) $\mathbb{V}(A\mathbf{Y}) = A\mathbb{V}(\mathbf{Y})A^\top$

Proof of: 3.7

We prove property (4) only.

$$\begin{aligned} \mathbb{V}(A\mathbf{Y}) &= \mathbb{E}[(A\mathbf{Y} - \mathbb{E}[A\mathbf{Y}]) (A\mathbf{Y} - \mathbb{E}[A\mathbf{Y}])^\top] \\ &= \mathbb{E}[(A\mathbf{Y} - A\mathbb{E}[\mathbf{Y}]) (A\mathbf{Y} - A\mathbb{E}[\mathbf{Y}])^\top] \\ &= \mathbb{E}[A(\mathbf{Y} - \mathbb{E}[\mathbf{Y}]) (A(\mathbf{Y} - \mathbb{E}[\mathbf{Y}]))^\top] \\ &= \mathbb{E}[A(\mathbf{Y} - \mathbb{E}[\mathbf{Y}]) (\mathbf{Y} - \mathbb{E}[\mathbf{Y}])^\top A^\top] \\ &= A\mathbb{E}[(\mathbf{Y} - \mathbb{E}[\mathbf{Y}]) (\mathbf{Y} - \mathbb{E}[\mathbf{Y}])^\top] A^\top \\ &= A\mathbb{V}(\mathbf{Y})A^\top \end{aligned}$$

EXAMPLE 3.8: Calculations with MLR Variables

Let $\mathbf{Y} = (Y_1, Y_2, Y_3)^\top$. Suppose $\mathbb{E}[\mathbf{Y}] = (3, 1, 2)^\top$. Let $\mathbb{V}(\mathbf{Y}) = \begin{bmatrix} 4 & 1/2 & -2 \\ 1/2 & 1 & 0 \\ -2 & 0 & 3 \end{bmatrix}$ and $\mathbf{a} = (1, -1, 2)$

and $A = \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 7 & 8 & 9 \end{bmatrix}$. Note that \mathbf{a} is a 1×3 row vector. Compute the following.

- (i) $\mathbb{E}[\mathbf{aY}]$
- (ii) $\mathbb{V}(\mathbf{aY})$

- (iii) $\mathbb{E}[AY]$
(iv) $\mathbb{V}(AY)$

Solution. We do the first two and leave the rest as an exercise.

(i) $\mathbb{E}[aY] = a\mathbb{E}[Y] = [1 \quad -1 \quad 2] \begin{bmatrix} 3 \\ 1 \\ 2 \end{bmatrix} = 1(3) - 1(1) + 2(2) = 6.$

(ii)

$$\begin{aligned} \mathbb{V}(aY) &= a\mathbb{V}(Y)a^\top \\ &= [1 \quad -1 \quad 2] \begin{bmatrix} 4 & 1/2 & -2 \\ 1/2 & 1 & 0 \\ -2 & 0 & 3 \end{bmatrix} \begin{bmatrix} 1 \\ -1 \\ 2 \end{bmatrix} \\ &= [1 \quad -1 \quad 2] \begin{bmatrix} 4(1) + (1/2)(-1) - 2(2) \\ (1/2)(1) + 1(-1) + 0(2) \\ -2(1) + 0(-1) + 3(2) \end{bmatrix} \\ &= [1 \quad -1 \quad 2] \begin{bmatrix} -1/2 \\ -1/2 \\ 4 \end{bmatrix} \\ &= 1(-1/2) - 1(-1/2) + 2(4) \\ &= 8 \end{aligned}$$

DEFINITION 3.9: Multivariate normal distribution

Let $Y = (Y_1, \dots, Y_n)^\top$ be a random vector. We say that $Y \sim \text{MVN}(\mu, \Sigma)$; that is, Y follows a **multivariate normal distribution** (MVN) when

$$f(y; \mu, \Sigma) = \frac{1}{(2\pi)^{n/2} |\Sigma|^{1/2}} \exp \left\{ -\frac{1}{2} (y - \mu)^\top \Sigma^{-1} (y - \mu) \right\}$$

where μ is defined as the **mean vector**, and Σ is defined as the **covariance matrix**. Note that Σ^{-1} is the inverse of the covariance matrix and $|\Sigma|$ is the determinant of Σ .

THEOREM 3.10: Properties of Multivariate Normal Distribution

Let $Y = (Y_1, \dots, Y_n)^\top \sim \text{MVN}(\mu, \Sigma)$ and a be a $1 \times n$ row vector of constants and A be an $n \times n$ matrix of constants.

(1) Linear transformations of MVN is MVN, so

$$aY \sim \text{MVN}(a\mu, a\Sigma a^\top)$$

$$AY \sim \text{MVN}(A\mu, A\Sigma A^\top)$$

(2) Marginal distribution of Y_i is Normal,

$$Y_i \sim N(\mu_i, \Sigma_{ii})$$

In fact, any subset of Y_i 's is MVN

(3) Conditional MVN is MVN, e.g. $Y_1 \mid Y_2, \dots, Y_n$

(4) Another property:

$$\text{Cov}(Y_i, Y_j) = 0 \iff Y_i, Y_j \text{ independent}$$

that is, Y_i and Y_j are uncorrelated.

$$\Sigma_{ij} = 0$$

Recall that last lecture, for a MLR, we have $\mathbf{Y} = \mathbf{X}\mathbf{B} + \boldsymbol{\varepsilon}$ with the assumption that $\varepsilon \stackrel{\text{iid}}{\sim} N(0, \sigma^2)$. Therefore, for a random vector $\boldsymbol{\varepsilon}$, we have

$$\boldsymbol{\varepsilon} \sim \text{MVN} \left(\begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} \sigma^2 & 0 & \cdots & 0 & 0 \\ 0 & \sigma^2 & \cdots & 0 & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & \cdots & \sigma^2 & 0 \\ 0 & 0 & \cdots & 0 & \sigma^2 \end{bmatrix} \right) = (\mathbf{0}_{n \times 1}, \sigma^2 \mathbf{I}_{n \times n})$$

since $\text{Cov}(\varepsilon_1, \varepsilon_2) = 0$ due to independence.

Thus, $\mathbf{Y} \sim \text{MVN}(\mathbf{X}\mathbf{B}, \sigma^2 \mathbf{I})$.

DEFINITION 3.11: Least squares for MLR

We define the **least squares for a multiple linear regression model** as

$$S(\beta_0, \beta_1, \dots, \beta_p) = \sum_{i=1}^n (y_i - \underbrace{(\beta_0 + \beta_1 x_{i1} + \cdots + \beta_p x_{ip})}_{\mathbb{E}[Y_i] = \mu_i})^2$$

THEOREM 3.12: Least Square Estimates (LSEs) for MLR

Minimizing $S(\beta_0, \beta_1, \dots, \beta_p)$, gives the least squares estimate $\hat{\boldsymbol{\beta}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$.

Proof of: 3.12

The first partial is $\frac{\partial S}{\partial \beta_0} = \sum_{i=1}^n 2(y_i - \mu_i)(-1)$, and all other partials for $j = 1, \dots, p$ are

$$\frac{\partial S}{\partial \beta_j} = \sum_{i=1}^n 2(y_i - \mu_i)(-x_{ij})$$

Set $\frac{\partial S}{\partial \beta_0} = 0$ and $\frac{\partial S}{\partial \beta_j} = 0$ for $j = 1, \dots, p$ to get

$$\begin{cases} \sum_{i=1}^n (y_i - \mu_i) = 0 \iff \mathbf{1}_{n \times n}^\top (\mathbf{y} - \boldsymbol{\mu}) = 0 \\ \sum_{i=1}^n (y_i - \mu_i) x_{ij} = 0 \iff \mathbf{x}_j^\top (\mathbf{y} - \boldsymbol{\mu}) = 0 \quad j = 1, \dots, p \end{cases}$$

since we recall that

$$\mathbf{X} = \begin{bmatrix} 1 & x_{11} & x_{12} & \cdots & x_{1p} \\ \vdots & & & & \\ 1 & x_{n1} & x_{n2} & \cdots & x_{np} \end{bmatrix} = [\mathbf{1}_{n \times 1} \quad \mathbf{x}_1 \quad \cdots \quad \mathbf{x}_{p-1} \quad \mathbf{x}_p]$$

Therefore,

$$\mathbf{X}^\top (\mathbf{y} - \mathbf{X}\mathbf{B}) = 0 \iff \mathbf{X}^\top \mathbf{y} - \mathbf{X}^\top \mathbf{X}\mathbf{B} = 0 \iff \mathbf{X}^\top \mathbf{X}\mathbf{B} = \mathbf{X}^\top \mathbf{y} \iff \mathbf{B} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$$

assuming $\mathbf{X}^\top \mathbf{X}$ is invertible (full rank of $p + 1$ or linearly independent columns). So, the LS solution for \mathbf{B} is given by $\hat{\mathbf{B}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$.

DEFINITION 3.13: Residuals for MLR

The **residuals** for a multiple linear regression model is defined as

$$e_i = y_i - \underbrace{(\hat{\beta}_0 + \hat{\beta}_1 x_{i1} + \cdots \hat{\beta}_p x_{ip})}_{\text{fitted value } \mu_i}$$

or equivalently, $\hat{\mu} = X\hat{B}$ and $e = y - \hat{\mu}$.

The estimate σ^2 based on e_i 's is

$$\hat{\sigma}^2 = \frac{SS(\text{Res})}{n - (p + 1)} = \frac{\sum_{i=1}^n e_i^2}{n - p - 1} = \frac{e^\top e}{n - p - 1}$$

since d.f. is $n - (\text{no. estimated parameters})$. When viewed as a random variable,

$$\frac{(n - p - 1)\hat{\sigma}^2}{\sigma^2} \sim \chi^2(n - p - 1)$$

Inference for $\hat{\beta} = (\hat{\beta}_0, \dots, \hat{\beta}_p)^\top = (X^\top X)^{-1} X^\top Y$.

Note that $\hat{\beta}$ is a matrix of constants and Y is a random vector, and $Y \sim \text{MVN}(X\beta, \sigma^2 I)$, so

$$\begin{aligned} \mathbb{E}[\hat{\beta}] &= \mathbb{E}[(X^\top X)^{-1} X^\top Y] \\ &= (X^\top X)^{-1} X^\top \mathbb{E}[Y] \\ &= (X^\top X)^{-1} (X^\top X) \beta \\ &= \beta \end{aligned}$$

That is, $\mathbb{E}[\hat{\beta}_0], \dots, \mathbb{E}[\hat{\beta}_p] = \beta_p$ all unbiased.

$$\begin{aligned} \mathbb{V}((X^\top X)^{-1} X^\top Y) &= (X^\top X)^{-1} X^\top \mathbb{V}(Y) [(X^\top X)^{-1} X^\top]^\top \\ &= (X^\top X)^{-1} X^\top \sigma^2 I (X^\top)^\top [(X^\top X)^{-1}]^\top && X^\top X \text{ symmetric} \\ &= \sigma^2 (X^\top X)^{-1} (X^\top X) (X^\top X)^{-1} \end{aligned}$$

Since $\hat{\beta}$ is a linear transformation of Y we have $\hat{\beta} \sim \text{MVN}(\beta, \sigma^2 \underbrace{(X^\top X)^{-1}}_V)$. We proved the following theorem.

THEOREM 3.14: Distribution of $\hat{\beta}_j$

The distribution of a given $\hat{\beta}_j$ is

$$\hat{\beta}_j \sim N(\beta_j, \sigma^2 V_{jj})$$

from marginal property of MVN.

$$\begin{aligned} \frac{\hat{\beta}_j - \beta_j}{\sigma \sqrt{V_{jj}}} &\sim N(0, 1) \\ \frac{\hat{\beta}_j - \beta_j}{\hat{\sigma} \sqrt{V_{jj}}} &\sim t(n - p - 1) \end{aligned}$$

DEFINITION 3.15: Standard error for $\hat{\beta}_j$

We define the **standard error** of $\hat{\beta}_j$ as

$$\text{Se}(\hat{\beta}_j) = \hat{\sigma} \sqrt{V_{jj}}$$

So, a $(1 - \alpha)$ confidence interval for β_j is

$$\hat{\beta}_j \pm c \text{Se}(\hat{\beta}_j)$$

where c is $(1 - (\alpha/2))$ quantile of $t(n - p - 1)$.

To test $H_0: \beta_j = 0$ vs $H_A: \beta_j \neq 0$, calculate t -statistic $t = \frac{\hat{\beta}_j}{\text{Se}(\hat{\beta}_j)}$ reject at level α if $|t| > c$ and p -value is $2P(T \geq |t|)$ where $T \sim t(n - p - 1)$.

Interpretation of $\hat{\beta}$: fitted linear regression model says $\widehat{\mathbb{E}[Y]}$ (estimate of the expected response) is $\hat{\beta}_0 + \hat{\beta}_1 x_1 + \dots + \hat{\beta}_p x_p$.

- $\hat{\beta}_0$ is the estimate of expected response when all explanatory variables are equal to 0.
- $\hat{\beta}_j$ is the estimated change in expected response for a unit increase in x_j , when holding all other explanatory variables constant, e.g.

$$\hat{\beta}_0 + \hat{\beta}_1(x_1 + 1) + \dots + \hat{\beta}_p x_p - (\hat{\beta}_0 + \hat{\beta}_1 x_1 + \dots + \hat{\beta}_p x_p) = \hat{\beta}_1$$

REMARK 3.16

When it's written V_{jj} , that means the $j + 1^{\text{th}}$ column and $j + 1^{\text{th}}$ row since we start from index 0 for these matrices. Some unfortunate events may have happened on the quiz to me due to this.

EXAMPLE 3.17: Rocket MLR

Let $n = 12$, $\hat{\beta} = (473.6, 16.7, -1.09)^{\top} = (\hat{\beta}_0, \hat{\beta}_1, \hat{\beta}_2)^{\top}$.

- x_1 : nozzle area (1 = L, 0 = S)
- x_2 : propellant ratio
- Y : thrust

$$\hat{\sigma} = \sqrt{\frac{\sum_{i=1}^{12} e_i^2}{12 - 1 - 2}} = \sqrt{\frac{\mathbf{e}^{\top} \mathbf{e}}{9}} = 2.655$$

Interpretation of $\hat{\beta}$:

- $\hat{\beta}_1$ estimated change in expected thrust is 16.7 when changing small to large nozzle while holding other variables (propellant ratio) constant.
- $\hat{\beta}_2$ estimated thrust to decrease by 1.09 on average for a unit increase in propellant ratio while holding other variables (nozzle area) constant.

Given $\text{Se}(\hat{\beta}_2) = 0.94$, we compute the t -statistic for $H_0: \beta_2 = 0$ vs $H_A: \beta_2 \neq 0$ which is $t = -1.09/0.94 = -1.16$.

$$p\text{-value} = 2P(T \geq 1.16) = 0.275 \text{ from R where } T \sim t(9)$$

Do not reject H_0 (e.g. $\alpha = 0.05$), therefore propellant ratio does not significantly influence thrust.

LECTURE 7 | 2020-09-28

Recall that $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon} \sim \text{MVN}(\mathbf{X}\boldsymbol{\beta}, \sigma^2 \mathbf{I})$, and

- Estimates: $\hat{\boldsymbol{\beta}} = (\mathbf{X}^{\top} \mathbf{X})^{-1} \mathbf{X}^{\top} \mathbf{Y}$
- Fitted values: $\hat{\boldsymbol{\mu}} = \mathbf{X} \hat{\boldsymbol{\beta}}$
- Residuals: $\mathbf{e} = \mathbf{y} - \hat{\boldsymbol{\mu}}$
- Constants: $\mathbf{X} = [\mathbf{1} \quad \mathbf{x}_1 \quad \dots \quad \mathbf{x}_p]_{n \times (p+1)}$

- Values of responses: $\mathbf{y} = (y_1, y_2, \dots, y_n)^\top \in \mathbb{R}^n$

Author's Note: Geometric interpretation of data is omitted in these notes because I'm simply too lazy.

The span of X is $\text{Span}(X) = \{b_0 \mathbf{1} + b_1 \mathbf{x}_1 + \dots + b_p \mathbf{x}_p : b_0, \dots, b_p \in \mathbb{R}\} \subset \mathbb{R}^n$ which is all linear combinations of columns of X which is a subspace of \mathbb{R}^n , and by assumption we know $\text{rank}(X) = p + 1$.

We can say $\text{Span}(X)$ represents all possible vector values $X\mathbf{b}$ where $\mathbf{b} = (b_0, b_1, \dots, b_p)^\top$.

Generally, $\mathbf{y} \notin \text{Span}(X)$, so since the linear model is an approximation, ϵ variability not explained by model.

Intuitively, it makes sense to choose an estimate $\hat{\beta}$ so that $X\hat{\beta}$ is as close to \mathbf{y} as possible. Therefore, \mathbf{e} must be orthogonal to $\text{Span}(X) \iff \mathbf{e}$ is orthogonal to all columns of X .

$$\begin{aligned} \mathbf{1}^\top \cdot (\mathbf{y} - \hat{\boldsymbol{\mu}}) &= 0 \\ \mathbf{x}_1^\top \cdot (\mathbf{y} - \hat{\boldsymbol{\mu}}) &= 0 \\ &\vdots \\ \mathbf{x}_p^\top \cdot (\mathbf{y} - \hat{\boldsymbol{\mu}}) &= 0 \end{aligned}$$

which is the same as LS estimates. We also know $\hat{\boldsymbol{\mu}} = X\hat{\beta}$ and $\mathbf{e} = \mathbf{y} - \hat{\boldsymbol{\mu}}$.

DEFINITION 3.18: Hat matrix

The **hat matrix** is defined as $H = X(X^\top X)^{-1}X^\top$.

PROPOSITION 3.19: Properties of Hat Matrix

Let H be a hat matrix, then H has the following properties.

- (1) H is symmetric; that is, $H = H^\top$.
- (2) H is idempotent; that is, $H^2 = HH = H$.
- (3) $I - H$ is symmetric idempotent; that is, $(I - H)^2 = (I - H)(I - H) = I - H$.

Proof of: 3.19

We prove all three because it's easy.

- (1) $H^\top = [X(X^\top X)^{-1}X^\top]^\top = X(X^\top X)^{-1}X^\top = H$.
- (2) $HH = X(X^\top X)^{-1}(X^\top X)(X^\top X)^{-1}X^\top = H$.
- (3) $(I - H)(I - H) = I(I - H) - H(I - H) = II - IH - HI + HH = I - 2H + HH = I - 2H + H = I - H$.

Let's view $\hat{\boldsymbol{\mu}}$ and \mathbf{e} as random vectors

$$\begin{aligned} \hat{\boldsymbol{\mu}} &= X\hat{\beta} = X(X^\top X)^{-1}X^\top \mathbf{Y} = H\mathbf{Y} \\ \mathbf{e} &= \mathbf{Y} - \hat{\boldsymbol{\mu}} = I\mathbf{Y} - H\mathbf{Y} = (I - H)\mathbf{Y} \\ \mathbb{E}[\hat{\boldsymbol{\mu}}] &= \mathbb{E}[H\mathbf{Y}] = H\mathbb{E}[\mathbf{Y}] = X(X^\top X)^{-1}X^\top \underbrace{\mathbb{E}[\mathbf{Y}]}_{X\boldsymbol{\beta}} = X\boldsymbol{\beta} \\ \mathbb{V}(\hat{\boldsymbol{\mu}}) &= \mathbb{V}(H\mathbf{Y}) = H\mathbb{V}(\mathbf{Y})H^\top = H\sigma^2 I H^\top = \sigma^2(HH^\top) = \sigma^2 H \\ \mathbb{E}[\mathbf{e}] &= \mathbb{E}[(I - H)\mathbf{Y}] = \mathbb{E}[\mathbf{Y}] - \mathbb{E}[H\mathbf{Y}] = X\boldsymbol{\beta} - X\boldsymbol{\beta} = 0 \\ \mathbb{V}(\mathbf{e}) &= (I - H)\mathbb{V}(\mathbf{Y})(I - H)^\top = \sigma^2(I - H)(I - H)^\top = \sigma^2(I - H) \end{aligned}$$

So since $\hat{\boldsymbol{\mu}}$ and \mathbf{e} are linear transformations of \mathbf{Y} we have proved the following theorem.

THEOREM 3.20: Distribution of $\hat{\mu}$ and e

$\hat{\mu}$ and \hat{e} have the following distribution.

$$\begin{aligned}\hat{\mu} &\sim MVN(X\beta, \sigma^2 H) \\ \hat{e} &\sim MVN(0, \sigma^2(I - H))\end{aligned}$$

Suppose we want to predict response for \mathbf{x}_0 where the first 1 represents the intercept in the row vector.

$$\mathbf{x}_0 = [1 \quad x_{01} \quad x_{02} \quad \cdots \quad x_{0p}]_{1 \times (p+1)}$$

Let Y_0 random variable representing the response associated with \mathbf{x}_0 . The MLR says

$$Y_0 \sim N(\beta_0 + \beta_1 x_{01} + \cdots + \beta_p x_{0p}, \sigma^2)$$

So we predict the value

$$\hat{y}_0 = \hat{\beta}_0 + \hat{\beta}_1 x_{01} + \cdots + \hat{\beta}_p x_{0p} = \mathbf{x}_0 \hat{\beta}$$

which represents the estimated mean response given $x_{01}, x_{02}, \dots, x_{0p}$. Corresponding distribution has

$$\begin{aligned}\mathbb{E}[\hat{Y}_0] &= \mathbf{x}_0 \mathbb{E}[\hat{\beta}] = \mathbf{x}_0 \beta = \mathbb{E}[Y_0] \\ \mathbb{V}(\hat{Y}_0) &= \mathbf{x}_0 \mathbb{V}(\hat{\beta}) \mathbf{x}_0^\top = \mathbf{x}_0 \sigma^2 (X^\top X)^{-1} \mathbf{x}_0^\top\end{aligned}$$

We have proved the following theorem.

THEOREM 3.21: Distribution of Predictor

The distribution of \hat{Y}_0 which is a function of Y_1, \dots, Y_n is

$$\hat{Y}_0 \sim N(\mathbf{x}_0 \beta, \sigma^2 \mathbf{x}_0 (X^\top X)^{-1} \mathbf{x}_0^\top)$$

$$\frac{\hat{Y}_0 - \mathbf{x}_0 \beta}{\sigma \sqrt{\mathbf{x}_0 (X^\top X)^{-1} \mathbf{x}_0^\top}} \sim N(0, 1)$$

$$\frac{\hat{Y}_0 - \mathbf{x}_0 \beta}{\hat{\sigma} \sqrt{\mathbf{x}_0 (X^\top X)^{-1} \mathbf{x}_0^\top}} \sim t(n - (p + 1)) = t(n - p - 1)$$

A $(1 - \alpha)$ confidence interval for the mean response $y_0 = \mathbf{x}_0 \hat{\beta}$ given \mathbf{x}_0 is

$$\hat{y}_0 \pm c \hat{\sigma} \sqrt{\mathbf{x}_0 (X^\top X)^{-1} \mathbf{x}_0^\top}$$

where c is the $1 - \alpha/2$ quantile of $t(n - p - 1)$.

Prediction error: $Y_0 - \hat{Y}_0$ which are independent since Y_0 is a random variable with variance σ^2 and \hat{Y}_0 is a function of Y_1, \dots, Y_n . Therefore,

$$\begin{aligned}\mathbb{E}[Y_0 - \hat{Y}_0] &= \mathbf{x}_0 \beta - \mathbf{x}_0 \beta = 0 \\ \mathbb{V}(Y_0 - \hat{Y}_0) &= \mathbb{V}(Y_0) + (-1)^2 \mathbb{V}(\hat{Y}_0) = \sigma^2 + \sigma^2 (\mathbf{x}_0 (X^\top X)^{-1} \mathbf{x}_0^\top)\end{aligned}$$

We have proved the following theorem.

THEOREM 3.22: Distribution of Prediction Error

The distribution of the prediction error is

$$Y_0 - \hat{Y}_0 \sim N(0, \sigma^2(1 + \mathbf{x}_0(X^\top X)^{-1}\mathbf{x}_0^\top))$$

A $(1 - \alpha)$ prediction interval for the mean response $y_0 = \mathbf{x}_0^\top \hat{\boldsymbol{\beta}}$ given \mathbf{x}_0 is

$$\hat{y}_0 \pm c\hat{\sigma}\sqrt{1 + \mathbf{x}_0(X^\top X)^{-1}\mathbf{x}_0^\top}$$

where c is the $1 - \alpha/2$ quantile of $t(n - p - 1)$.

REMARK 3.23

Our intuition tells us that the prediction interval is wider than the confidence interval for mean. In other words, estimating an average is “easier” than an individual response.

The example done in R is included in the next page.

```
## NASA rocket data example
```

```
## From: R.S. Jankovsky, T.D. Smith, A.J. Pavli (1999). "High-Area-Ratio Rocket  
## Nozzle at High Combustion Chamber Pressure-Experimental and Analytical  
## Validation".
```

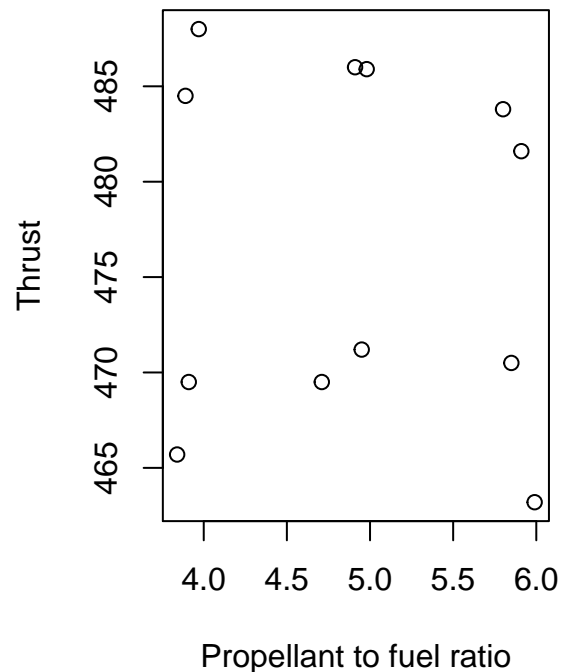
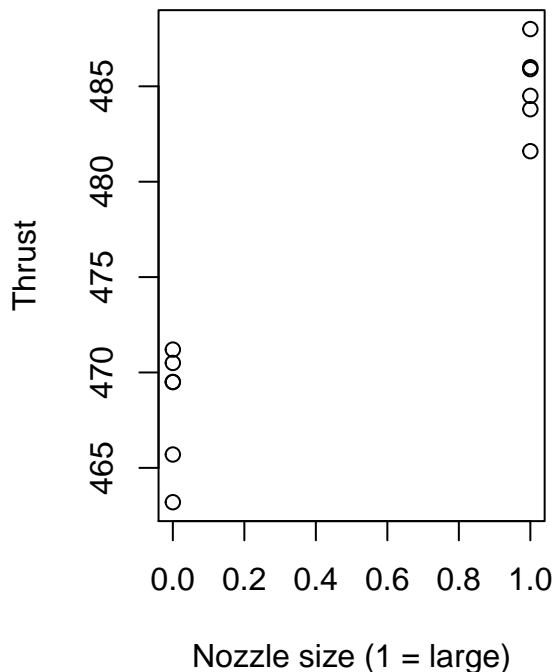
```
# setwd(...) first if your CSV file is somewhere else  
rocket <- read.csv(file="rocket.csv")  
# output all data in rocket vector  
rocket
```

```
##      thrust nozzle propratio  
## 1    488.0      1      3.97  
## 2    481.6      1      5.91  
## 3    485.9      1      4.98  
## 4    486.0      1      4.91  
## 5    484.5      1      3.89  
## 6    483.8      1      5.80  
## 7    463.2      0      5.99  
## 8    471.2      0      4.95  
## 9    469.5      0      3.91  
## 10   470.5      0      5.85  
## 11   469.5      0      4.71  
## 12   465.7      0      3.84
```

Y (thrust) is the response variable, and there are two explanatory variables x_1, x_2 (nozzle, propratio) where nozzle is coded as 1 if it's large.

```
# Scatter plots where mfrow is used to put multiple  
# plots on one image
```

```
par(mfrow = c(1,2))  
plot(rocket$nozzle, rocket$thrust, ylab="Thrust", xlab="Nozzle size (1 = large)")  
plot(rocket$propratio, rocket$thrust, ylab="Thrust", xlab="Propellant to fuel ratio")
```



Left is

nozzle size vs thrust. Right is propellant relationship vs thrust.

```
# Fit MLR using lm
m1 <- lm(thrust ~ nozzle + propratio, data = rocket)
summary(m1)

##
## Call:
## lm(formula = thrust ~ nozzle + propratio, data = rocket)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -3.8459 -1.7555  0.5934  1.2906  3.3008
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  473.6039     4.7158  100.430 4.88e-15 ***
## nozzle       16.7383     1.5329   10.919 1.71e-06 ***
## propratio   -1.0948     0.9414   -1.163  0.275
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 2.655 on 9 degrees of freedom
## Multiple R-squared:  0.9303, Adjusted R-squared:  0.9148
## F-statistic: 60.05 on 2 and 9 DF,  p-value: 6.238e-06
```

On the left it's Y (response variable) and on the right it's x_1, x_2 (explanatory variables). From summary, we get the estimate vector $\hat{\beta} = (473.6039, 16.7383, -1.0948)^\top$.

```
# Manual beta estimates where rep is used to make the columns of 1s
X <- cbind(rep(1, 12), rocket$nozzle, rocket$propratio) # X matrix
y <- matrix(rocket$thrust, ncol = 1) # response vector
beta_hat <- solve(t(X) %*% X) %*% t(X) %*% y
beta_hat
```

```
##           [,1]
## [1,] 473.603924
## [2,] 16.738319
## [3,] -1.094822
```

`solve` is used for the inverse. `%*%` is used for matrix-matrix multiplication, and `t(X)` is used for transposing X .

```
# Manual sigma estimate
mu_hat <- X %*% beta_hat # fitted values
e <- y - mu_hat # residuals
sigma_hat <- sqrt((t(e) %*% e) / 9) # Note n-p-1 = 12-2-1 = 9
sigma_hat
```

```
##           [,1]
## [1,] 2.6545
```

```
sigma_hat <- sqrt( sum(e^2) / 9) # equivalent
sigma_hat
```

```
## [1] 2.6545
```

- $\hat{\mu} = X\hat{\beta}$

- $e = y - \hat{\mu}$
- $\hat{\sigma} = \sqrt{\left(\sum_{i=1}^n e_i^2\right)/9} = 2.6545$, or
- $\hat{\sigma} = \sqrt{(e^\top e)/9} = 2.6545$

```
# Covariance matrix of beta_hat
vcov(m1)
```

```
##           (Intercept)      nozzle  propratio
## (Intercept)  22.238325 -1.02316688 -4.32080608
## nozzle      -1.023167   2.34987593 -0.03102117
## propratio   -4.320806 -0.03102117  0.88631920
```

```
sqrt(diag(vcov(m1))) # SEs of individual betas
```

```
## (Intercept)      nozzle  propratio
##  4.7157528   1.5329305   0.9414453
```

```
# Manual
se_beta <- sigma_hat * sqrt(diag(solve(t(X) %*% X)))
se_beta
```

```
## [1] 4.7157528 1.5329305 0.9414453
```

- $Se(\hat{\beta}) = \hat{\sigma} \sqrt{(X^\top X)^{-1}} = (4.71, 1.53, 0.94)^\top$

```
# Estimate the mean response for units with small nozzle and propellant ratio 5.5
# include a 95% CI
predict(object = m1, newdata = data.frame(nozzle = 0, propratio = 5.5),
        interval = "confidence", level = 0.95)
```

```
##           fit      lwr      upr
## 1 467.5824 464.7929 470.3719
```

Therefore, $\hat{y}_0 = 467.58$. The 95% confidence interval for the mean response given \mathbf{x}_0 is $[464.7929, 470.3719]$.

```
# Manual calculation
x0 <- matrix(c(1, 0, 5.5), nrow = 1)
y0_hat <- x0 %*% beta_hat
y0_hat
```

```
##           [,1]
## [1,] 467.5824
```

```
# mu0 is also known as \hat{Y}_0
se_mu0 <- sigma_hat * sqrt(x0 %*% solve(t(X) %*% X) %*% t(x0))
se_mu0
```

```
##           [,1]
## [1,] 1.233132
```

```
crit_val <- qt(0.975, 9)
ci_lo <- y0_hat - crit_val*se_mu0
ci_hi <- y0_hat + crit_val*se_mu0
c(y0_hat, ci_lo, ci_hi)
```

```
## [1] 467.5824 464.7929 470.3719
```

- $\mathbf{x}_0 = [1 \ 0 \ 5.5]$

- $\hat{y}_0 = \mathbf{x}_0 \hat{\beta} = 467.5824$
- $Se(\hat{Y}_0) = \hat{\sigma} \sqrt{\mathbf{x}_0 (X^\top X)^{-1} \mathbf{x}_0^\top} = 1.233132$

Therefore, $\hat{y}_0 = 467.58$. The 95% confidence interval for the mean response given \mathbf{x}_0 is [464.7929, 470.3719].

```
# Predict the value of the response for a unit with small nozzle and propellant ratio 5.5
# include a 95% PI
predict(object = m1, newdata = data.frame(nozzle = 0, propratio = 5.5),
        interval = "prediction", level = 0.95)
```

```
##           fit           lwr           upr
## 1 467.5824 460.9612 474.2036
```

Therefore, $y_0 = 467.5824$. The 95% prediction interval for the response (y_0) given \mathbf{x}_0 is [460.9612474.2036].

```
# Manual calculation for an individual
x0 <- matrix(c(1, 0, 5.5), nrow = 1)
y0_hat <- x0 %*% beta_hat
se_y0 <- sigma_hat * sqrt(1+ x0 %*% solve(t(X) %*% X) %*% t(x0))
se_y0
```

```
##           [,1]
## [1,] 2.926941
```

```
crit_val <- qt(0.975,9)
pi_lo <- y0_hat - crit_val*se_y0
pi_hi <- y0_hat + crit_val*se_y0
c(y0_hat, pi_lo, pi_hi)
```

```
## [1] 467.5824 460.9612 474.2036
```

- $Se(Y_0 - \hat{Y}_0) = \hat{\sigma} \sqrt{1 + \mathbf{x}_0 (X^\top X)^{-1} \mathbf{x}_0^\top} = 2.926941$

Handling categorical variables: when there are explanatory variables with values that fall into one of several categories.

- e.g. nozzle large/small, if just binary, code as 1 and 0
- ordered small, medium, large or not red, blue green

Approach: can convert to indicator variables or treat as numerical if it makes sense to do so.

Example: CQI (2018)

Extract a few variables:

	Acidity	Method
1	8.7	Washed-wet
2	8.3	Washed-wet
3	8.2	Natural-dry
4	8.4	Semi-washed/pulped

Flavour (response)

How to set up X ? For example,

$$x_{i2} = \begin{cases} 0 & \text{dry} \\ 1 & \text{semi} \\ 2 & \text{wet} \end{cases}$$

Not generally appropriate unless we think a response is linear according to this scheme.

More flexible approach: indicator/dummy variables

$$x_{i2} = \begin{cases} 1 & \text{semi} \\ 0 & \text{otherwise} \end{cases}, \quad x_{i3} = \begin{cases} 1 & \text{wet} \\ 0 & \text{otherwise} \end{cases}$$

Therefore,

$$X = \begin{bmatrix} 1 & 8.7 & 0 & 1 \\ 1 & 8.3 & 0 & 1 \\ 1 & 8.2 & 0 & 0 \\ 1 & 8.4 & 1 & 0 \end{bmatrix}$$

Why not $x_{i4} = \begin{cases} 1 & \text{dry} \\ 0 & \text{otherwise} \end{cases}$? If we did that, we would have

$$X = \begin{bmatrix} 1 & 8.7 & 0 & 1 & 0 \\ 1 & 8.3 & 0 & 1 & 0 \\ 1 & 8.2 & 0 & 0 & 1 \\ 1 & 8.4 & 1 & 0 & 0 \end{bmatrix}$$

This has linearly dependent columns since $x_4 = 1 - x_2 - x_3$. There is no new information and X would not have full rank.

Model: $Y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \beta_3 x_{i3} + \varepsilon_i$.

Interpretation:

- mean flavour if acidity = x_{01} and method dry is $\beta_0 + \beta_1 x_{01}$.
- mean flavour if acidity = x_{01} and method wet is $\beta_0 + \beta_1 x_{01} + \beta_3$.
- mean flavour if acidity = x_{01} and method semi is $\beta_0 + \beta_1 x_{01} + \beta_2$.

- β_2 is the difference between semi and dry in expected response (holding acidity constant)
- β_3 is the difference between wet and dry in expected response (holding acidity constant)
- $\beta_2 - \beta_3$ is the difference between semi and wet (holding other variables constant)

$\hat{\beta} \sim \text{MVN}(\beta, \sigma^2 V)$ where $V = (X^\top X)^{-1}$.

- We know $\hat{\beta}_j \sim N(\beta_j, \sigma^2 V_{jj})$ with $\text{Se}(\hat{\beta}_j) = \hat{\sigma} \sqrt{V_{jj}}$ where $j = 0, \dots, p$.
- What about $\beta_2 - \beta_3$?

$$\mathbb{V}(\hat{\beta}_2 - \hat{\beta}_3) = \mathbb{V}(\hat{\beta}_2) - \mathbb{V}(\hat{\beta}_3) - 2\text{Cov}(\hat{\beta}_2, \hat{\beta}_3) = \sigma^2 V_{22} + \sigma^2 V_{33} - 2\sigma^2 V_{23}$$

Therefore,

$$\text{Se}(\hat{\beta}_2 - \hat{\beta}_3) = \hat{\sigma} \sqrt{V_{22} + V_{33} - 2V_{23}}$$

Now, we can construct a CI for $\beta_2 - \beta_3$.

In general, for an explanatory variable with k categories. We need $k - 1$ indicator variables.

LECTURE 9 | 2020-10-05

Analysis of variance (ANOVA): how well does our regression model fit our response variable?

Variability in response can be measured by “total sum of squares:”

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

as seen in HW1, it's closely related to sample variance of y_1, \dots, y_n , which is $\text{SS}(\text{Total})/(n - 1)$.

ANOVA decomposes $\text{SS}(\text{Total}) = \text{SS}(\text{Reg}) + \text{SS}(\text{Res})$ where $\text{SS}(\text{Reg})$ is the regression sum of squares and $\text{SS}(\text{Res})$ is the residual sum of squares.

The regression sum of squares is variation explained by the model and the residual sum of squares is the variation not explained by the regression model.

Using the fact that

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

When regression fits data well, the observations y_i tend to be much closer to $\hat{\mu}_i$. Note that \bar{y} is line a regression line with $\beta_1 = 0$.

Mathematically,

$$\underbrace{\sum_{i=1}^n (y_i - \bar{y})^2}_{\text{SS}(\text{Total})} = \underbrace{\sum_{i=1}^n (\hat{\mu}_i - \bar{y})^2}_{\text{SS}(\text{Reg})} + \underbrace{\sum_{i=1}^n (y_i - \hat{\mu}_i)^2}_{\text{SS}(\text{Res})}$$

since we showed that $\sum_{i=1}^n (\hat{\mu}_i - \bar{y}) \underbrace{(y_i - \hat{\mu}_i)}_{e_i} = 0$ in HW1 for SLR. It's also true for MLR since

$$\sum_{i=1}^n (\hat{\mu}_i - \bar{y}) e_i = \sum_{i=1}^n (e_i \hat{\mu}_i) - \bar{y} \sum_{i=1}^n e_i = \hat{\mu}^\top e - \bar{y} \mathbf{1}^\top e = 0$$

Recall: $\mathbf{1}^\top e = 0$ is one of LS equations, and $\hat{\mu} = X\hat{\beta}$ is in $\text{Span}(X)$, so e is orthogonal to $\text{Span}(X)$, so $\hat{\mu}^\top e = 0$.

F is used to test the overall significance of regression (later).

Table 1: ANOVA Table

Source	d.f.	SS	Mean Square	F
Regression	p	SS(Reg)	SS(Reg)/ p	MS(Reg)/MS(Res)
Residual	$n - p - 1$	SS(Res)	SS(Res)/($n - p - 1$) = $\hat{\sigma}^2$	
Total	$n - 1$	SS(Total)		

We call the **coefficient of determination** $R^2 = \text{SS(Reg)}/\text{SS(Total)} = 1 - \text{SS(Res)}/\text{SS(Total)}$. clearly, $0 \leq R^2 \leq 1$. It is the proportion of variation (in our response variable) that is explained by the regression model. Larger R^2 means the fitted values are closer to the observations y_i , which means the residuals are small; that is, smaller SS(Res). Note that (HW1) in SLR, R^2 is equivalent to the square of the sample correlation between x and y based on $(x_1, y_1), \dots, (x_n, y_n)$.

Table 2: Rocket ANOVA Table

Source	d.f.	SS	Mean Square	F
Regression	2	846.2	423.1	60
Residual	9	63.42	7.05	
Total	11	909.62		

Response thrust $R^2 = 846.2/909.62 \approx 0.93$. R^2 interpretation: regression model with nozzle size and propellant ratio explains 93% of variation in thrust (response).

LECTURE 10 | 2020-10-07

Hypothesis testing based on F distribution

So far we've tested $H_0: \beta_j = 0$ vs $H_A: \beta_j \neq 0$ involving individual parameters, using t distribution.

Now consider hypothesis test of the form $H_0: A\beta = \mathbf{0}$ where A is a matrix of constraints specifying linear combinations of parameters.

EXAMPLE 3.24: Coffee Continued

The full model is:

$$Y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \beta_3 x_{i3} + \varepsilon_i$$

- Y_i is the flavour
- x_{i1} is acidity
- x_{i2} is 1 if semi, and 0 otherwise.
- x_{i3} is 1 if wet, and 0 otherwise.

Example 1.

- $H_0: \beta_1 = \beta_2 = \beta_3 = 0$ versus
- H_A : at least one of $\beta_1, \beta_2, \beta_3$ not 0.
- If H_0 is true, the model reduces to $Y_i = \beta_0 + \varepsilon_i$.
- This tests overall significance of regression (whether any of predictors impact response)
- $A = \begin{bmatrix} 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}$. Note that row i considers the constraint of $\beta_i = 0$ for $i = 1, 2, 3$ in this example.

Example 2.

- $H_0: \beta_2 = \beta_3 = 0$

- If H_0 is true, $Y_i = \beta_0 + \beta_1 x_{i1} + \varepsilon_i$
- Q: Is reduced model with only acidity plausible?
- $A = \begin{bmatrix} 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}$. Note that $A\beta = \mathbf{0}_{1 \times 2}$

Example 3.

- $H_0: \beta_2 - \beta_3 = 0$
- $H_A: \beta_2 \neq \beta_3$
- $Y_i = \beta_0 + \beta_1 x_{i1} + \beta_2(x_{i2} + x_{i3}) + \varepsilon_i$ where $(x_{i2} + x_{i3})$ is 1 if semi/wet and 0 if dry.
- Do the wet and semi methods have the same impact on the response (holding acidity constant)?
- $A = \begin{bmatrix} 0 & 0 & 1 & -1 \end{bmatrix}$

In general, with ℓ constraints. A is an $\ell \times (p + 1)$ matrix with rank ℓ . Recall that

$$\text{Span}(X) = \{\beta_0 \mathbf{1} + \beta_1 \mathbf{x}_1 + \cdots + \beta_p \mathbf{x}_p\}$$

Let

$$\text{Span}(X)_A = \{\beta_0 \mathbf{1} + \beta_1 \mathbf{x}_1 + \cdots + \beta_p \mathbf{x}_p : A\beta = \mathbf{0}\}$$

which is a subspace of $\text{Span}(X)$ since any vector in $\text{Span}(X)_A$ is also in $\text{Span}(X)$. We call $\text{Span}(X)_A$ the $\text{Span}(X)$ with constraint A on β .

Let $\hat{\mu}_A$ denote the fitted values from fitting the reduced model. The residual if we fit the model with $A\beta = \mathbf{0}$ is $e_A = \mathbf{y} - \hat{\mu}_A$.

If $H_0: A\beta = \mathbf{0}$ is true, then $\hat{\mu}$ and $\hat{\mu}_A$ should be close; that is, the model makes similar predictions whether we set $A\beta = \mathbf{0}$ or not when fitting the model.

So to assess whether H_0 is plausible, look at $\|\hat{\mu} - \hat{\mu}_A\|$ where $\|\cdot\|$ is Euclidean or L_2 norm. That is,

$$\|\hat{\mu} - \hat{\mu}_A\| = \sqrt{(\hat{\mu} - \hat{\mu}_A)^\top (\hat{\mu} - \hat{\mu}_A)}$$

If it's "large" or "small" (close to 0) where large gives evidence against H_0 and small gives evidence for H_0 .

By Pythagoras,

$$\|\mathbf{y} - \hat{\mu}_A\|^2 = \|\mathbf{y} - \hat{\mu}\|^2 + \|\hat{\mu} - \hat{\mu}_A\|^2 \quad \text{or} \quad \|e_A\|^2 = \|e\|^2 + \|\hat{\mu} - \hat{\mu}_A\|^2$$

or equivalently $e_A^\top e_A = e^\top e + \|\hat{\mu} - \hat{\mu}_A\|^2$ where $e_A^\top e_A$ is the sum of squares residual in the reduced model and $e^\top e$ is the sum of squares residual in the full model.

We define $e_A^\top e_A = \text{SS}(\text{Res})_A$ and $e^\top e = \text{SS}(\text{Res})$.

Thus, $\|\hat{\mu} - \hat{\mu}_A\|^2 = \text{SS}(\text{Res})_A - \text{SS}(\text{Res}) \geq 0$ additional sum of squares explained by full model vs reduced one with constraints A .

Practical implications:

- $\text{SS}(\text{Res})$ cannot decrease when constraints applied.
- Equivalently, full model always has small (or equal) $\text{SS}(\text{Res})$ for a fixed $\text{SS}(\text{Tot})$ and thus higher R^2 compared to a reduced model.

Define test statistic:

$$F = \frac{(\text{SS}(\text{Res})_A - \text{SS}(\text{Res}))/\ell}{\text{SS}(\text{Res})/(n - p - 1)} = \frac{(\text{SS}(\text{Res})_A - \text{SS}(\text{Res}))/\ell}{\hat{\sigma}^2}$$

DEFINITION 3.25: F distribution

If $U \sim \chi^2(a)$ and $V \sim \chi^2(b)$ are independent. We say F follows an F **distribution** if

$$F = \frac{U/a}{V/b}$$

and write $F \sim F(a, b)$.

Here, we have these facts when H_0 is true

$$V = \frac{\hat{\sigma}^2(n-p-1)}{\sigma^2} \sim \chi^2(n-p-1)$$

$$U = \frac{\|\hat{\boldsymbol{\mu}} - \hat{\boldsymbol{\mu}}_A\|^2}{\sigma^2} \sim \chi^2(\ell)$$

where U and V are independent. Therefore,

$$F = \frac{\frac{\|\hat{\boldsymbol{\mu}} - \hat{\boldsymbol{\mu}}_A\|^2}{\sigma^2} \left(\frac{1}{\ell}\right)}{\frac{\hat{\sigma}^2(n-p-1)}{\sigma^2} \left(\frac{1}{n-p-1}\right)} \sim F(\ell, n-p-1)$$

when H_0 is true. Reject H_0 : $A\boldsymbol{\beta} = \mathbf{0}$ at level α if F is greater than $(1-\alpha)$ quantile of $F(\ell, n-p-1)$ and p -value is $P(Y \geq F)$ where $Y \sim F(\ell, n-p-1)$.

Relation to T distribution: Say $Y \sim t(a)$

$$Y = \frac{Z}{\sqrt{U/a}}$$

where $Z \sim N(0, 1)$ and $U \sim \chi^2(a)$ are independent. Squaring everything,

$$Y^2 = \frac{Z^2}{U/a}$$

and we know $Z^2 \sim \chi^2(1)$. Therefore, $Y^2 \sim F(1, a)$ (we divide by 1 in the numerator).

Thus, if our hypothesis test has one constraint, then F test is equal to t test of same hypothesis; for example, H_0 : $\beta_1 = 0$ versus H_A : $\beta_1 \neq 0$.

LECTURE 11 | 2020-10-19

Recall the general linear hypothesis: H_0 : $A\boldsymbol{\beta} = \mathbf{0}$ vs H_A : $A\boldsymbol{\beta} \neq \mathbf{0}$ where A gives ℓ constraints.

$$F \text{ statistic} = \frac{(\text{SS}(\text{Res})_A - \text{SS}(\text{Res}))/\ell}{\text{SS}(\text{Res})/(n-p-1)} = \frac{(\text{SS}(\text{Res})_A - \text{SS}(\text{Res}))/\ell}{\hat{\sigma}^2}$$

compare to $F(\ell, n-p-1)$.

Special case: overall test of significance

“are any predictors related to response?”

- H_0 : $\beta_1 = \beta_2 = \dots = \beta_p = 0$
- H_A : $\beta_j \neq 0$ for at least one j

$$A = \begin{bmatrix} 0 & 1 & 0 & 0 & \cdots & 0 & 0 \\ 0 & 0 & 1 & 0 & \cdots & 0 & 0 \\ \vdots & & & \ddots & & & \\ 0 & 0 & 0 & 0 & \cdots & 1 & 0 \\ 0 & 0 & 0 & 0 & \cdots & 0 & 1 \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_{p-1} \\ \beta_p \end{bmatrix}$$

If H_0 is true: $Y_i = \beta_0 + \varepsilon_i$ where $Y_i \sim N(\beta_0, \sigma^2)$.

Fit reduced model; that is, in this case estimate β_0 using least squares, minimize $\sum_{i=1}^n (y_i - \beta_0)^2$, which can be shown $\hat{\beta}_0 = \bar{y}$. So,

$$SS(\text{Res})_A = \sum_{i=1}^n e_i^2 = \sum_{i=1}^n (y_i - \hat{\mu}_i)^2 = \sum_{i=1}^n (y_i - \bar{y})^2 = SS(\text{Total})$$

Then,

$$F = \frac{(SS(\text{Total}) - SS(\text{Res}))/p}{SS(\text{Res})/(n - p - 1)} = \frac{SS(\text{Reg})/p}{SS(\text{Res})/(n - p - 1)} = \frac{MS(\text{Reg})}{MS(\text{Res})} \leftarrow F \text{ statistic on ANOVA table}$$

LECTURE 12 | 2020-10-21

Multicollinearity: occurs when some explanatory variables have a **strong linear** relationship amongst themselves. For example, this might occur exactly

$$\mathbf{x}_3 = \alpha_0 \mathbf{1} + \alpha_1 \mathbf{x}_1 + \alpha_2 \mathbf{x}_2$$

in which case the columns of X would be **linearly dependent** and $X^\top X$ does not have an inverse. Practically, there is no new info including \mathbf{x}_3 when $\mathbf{x}_1, \mathbf{x}_2$ are in the model. **Approximately**,

$$\mathbf{x}_3 \approx \alpha_0 \mathbf{1} + \alpha_1 \mathbf{x}_1 + \alpha_2 \mathbf{x}_2$$

in which case, the columns of X are close to being linearly dependent which causes the variances $\mathbb{V}(\hat{\beta}_j)$ to be **inflated**, in turn leads to inaccurate confidence intervals and conclusions of hypothesis tests for the regression parameters, in practice. $\text{Se}(\hat{\beta}_j)$ when fitting models can change drastically when adding/moving variables from the model.

EXAMPLE 3.26: Hockey (NHL)

Suppose Goals + Assists = Points and we want to predict a forward's salary.

- x_1 = Goals
- x_2 = Assists
- x_3 = Points

However, $x_3 = x_1 + x_2$ so we have exact multicollinearity.

EXAMPLE 3.27: Burmese Pythons in Florida (2017)

- y = fat content
- x_1 = mass
- x_2 = overall length
- x_3 = snout-to-vent length

It turns out that x_2 and x_3 are highly correlated. Including all variables in regression lead to inflated $\text{Se}(\hat{\beta}_2)$ and $\text{Se}(\hat{\beta}_3)$.

3.1 Detection of Multicollinearity

If two predictors are related

- Scatterplot matrix [all possible pairs of scatterplots b/w y, x_1, x_2, \dots, x_p]
- Correlation matrix (all pairwise correlations)

DEFINITION 3.28: Variance inflation factor

For multicollinearity between more than two predictors, we can define the **variance inflation factor** (VIF).

$$\text{VIF}_j = \frac{\mathbb{V}(\hat{\beta}_j)}{\mathbb{V}(\hat{\beta}_j^*)} \geq 1$$

for $j = 1, \dots, p$, where $\hat{\beta}_j$ is the estimate of β_j with all predictors in the model, and $\hat{\beta}_j^*$ estimate of β_j based on regression with x_j only. It can be shown that $\text{VIF}_j \geq 1$.

Fit MLR of x_j in terms of other predictors; that is,

$$x_{ij} = \alpha_0 + \alpha_1 x_{i1} + \dots + \alpha_{j-1} x_{i(j-1)} + \alpha_{j+1} x_{i(j+1)} + \dots + \alpha_p x_{ip} + \varepsilon_{ij}$$

and compute R^2 for this model, call it R_j^2 .

Intuition: if R_j^2 is close to 1, x_j is strongly related linearly to other predictors. It can be shown that

$$\text{VIF}_j = \frac{1}{1 - R_j^2}$$

Values of VIF larger than 10 are taken as solid evidence of multicollinearity; that is, $R_j^2 > 0.9$.

LECTURE 13 | 2020-10-26

Model selection: Given p explanatory variables, find the subset $k \leq p$ of explanatory variables (“reduced model”) that gives us the “best” model: goodness of fit, interpretability, predictive performance.

Some related concepts:

1. F tests compare between 2 specific models where test adequacy of a “reduced” model (subset, “nested”) relative to full model.

Quiz 4: $\beta_1 = \beta_2$ in part d-f

2. Multicollinearity: can affect interpretability of $\hat{\beta}_j$ usual interpretation “holding other variables constant” doesn’t really work when x_j is strongly correlated with other predictors.
3. R^2 is the proportion of variability in the response explained by the regression model. It always increases when adding variables.
4. $\hat{\sigma}^2$ is estimated residual variable, used for prediction, want $\hat{\sigma}^2$ small

Two key ingredients:

- Metric (or criterion) for comparing different models with potentially different number of predictors
- selection/search strategy (which models should we fit and test?)

Examples of metrics for model selection:

Adjusted R^2

$$R_{\text{adj}}^2 = 1 - \frac{\text{SS(Res)}/(n - k - 1)}{\text{SS(Total)}/(n - 1)}$$

for model with k predictors.

- $SS(\text{Res})/(n - k - 1)$ estimated $\hat{\sigma}^2$ for model with k predictors
- $SS(\text{Total})/(n - 1)$ is the sample variance of responses y_i .

$$\begin{aligned} R_{\text{adj}}^2 &= 1 - \frac{n-1}{n-k-1} = 1 - \left(1 + \frac{k}{n-k-1}\right) (1 - R^2) \\ &= R^2 - (1 - R^2) \frac{k}{n-k-1} \end{aligned}$$

Intuition: R_{adj}^2 accounts for number variables in model, *penalizes* inclusion of unimportant predictors; that is, $SS(\text{Res})$ has little decrease when adding such variables. Meanwhile, R^2 always increases with more predictors, but R_{adj}^2 can decrease if $SS(\text{Res})$ change is small.

While R_{adj}^2 loses its usual interpretation of R^2 , but can be used as a measure of “goodness of fit” and model selection criterion (e.g. pick subset of predictors that gives the highest R_{adj}^2).

EXAMPLE 3.29

$n = 25$ observations, $SS(\text{Total}) = 20$, in total $p = 6$ predictors. Suppose we’re considering on a subset of $k = 4$ predictors, and find:

	red	full
	$k = 4$	$p = 6$
$SS(\text{Total})$	20	20
$SS(\text{Res})$	10	9.8
R^2	0.5	$10.2/20 = 0.51$
R_{adj}^2	$1 - \frac{10/20}{20/24} \approx 0.4$	$1 - \frac{9.8/18}{20/24} \approx 0.347$
$\hat{\sigma}^2$	$10/20 = 0.5$	$9.8/18 \approx 0.544$

- $n - k - 1$ d.f. Res in reduced
- $n - p - 1$ d.f. Res in full

Remarks:

- $R_{\text{adj}}^2 < R^2$, but as $n \rightarrow \infty$, $R_{\text{adj}}^2 \rightarrow R^2$.
- model with higher R_{adj}^2 has lower $\hat{\sigma}^2$, thus is a reasonable metric for model selection

Akaike Information Criterion (AIC)

Let n be sample size, q is number of parameters [in MLR k predictors + 1 (intercept) + 1 (σ^2)]

$$AIC = 2q - \ln[L(\hat{\theta})]$$

where $L(\hat{\theta})$ is the likelihood function evaluated at $\hat{\theta}$ (parameter estimates). Note that LS estimates of β are equivalent to MLE under the usual normal assumptions on ϵ . Also, $2q$ is the penalty for including more predictors. With more parameters, $L(\hat{\theta})$ increases, offset by penalty $2q$.

Therefore, model with lower AIC is preferred; that is, differences in AIC matter not the value itself.

Bayesian Information Criterion (BIC)

Similar to AIC, but more strongly penalizes inclusion of more variables.

$$BIC = q \ln(n) - 2 \ln[L(\hat{\theta})]$$

where $q \ln(n)$ depends on sample size.

Recap:

- R^2 , AIC, BIC are all based on comparing the fitted models. In other words, they look at the explanatory power of the model.
- They all have penalties to try to prevent “overfitting.” That is, having too many variables might end up modelling spurious relationships that are actually noise.

Mean Square Prediction Error (MSPE)

Consider predictive performance of model on *new* data; that is, data *not* used in fitting of models. “Is model generalizable to new data?” Overfitted models tend to have high prediction error.

For example, via cross-validation schemes. We’ve given 4 examples of metrics/criteria for comparing models. Imagine we have p predictors:

- $\binom{p}{1}$ 1 predictors
- $\binom{p}{2}$ 2 predictors
- \vdots
- $\binom{p}{p}$ p predictors

$$\sum_{j=0}^p \binom{p}{j} = 2^p$$

Occam’s Razor: “The simplest explanation is usually the best one.”—William Ockham

LECTURE 14 | 2020-10-28

Model Selection

- Criteria: R_{adj}^2 , AIC, BIC, MSPE, etc. explicitly penalizes unnecessarily complex models.
 - Search strategies (use with chosen criterion)
- (i) Brute force: fit all possible regressions. With p predictors, we have $\sum_{j=0}^p \binom{p}{j} = 2^p$ possible models to fit.
- Finds optimal model that may be computationally intensive (or infeasible) if p is large.
- Idea: Find a “good” (useful) model in reasonable computational time (not necessarily optimal). Many strategies focus on adding/removing variables one at a time.
- (ii) Forward selection: add one variable at a time to model.
- Start with model that only has intercept β_0 .
 - Fit p simple linear regression models

$$\mathbf{y} = \beta_0 \mathbf{1} + \beta_1 \mathbf{x}_j + \varepsilon \quad j = 1, \dots, p$$

- Pick the best of p models (with 1 predictor) according to chosen criterion, and add that variable x_j to model.
- Fit $(p - 1)$ models containing x_j and one other variable.
 - If none of $(p - 1)$ models improves criterion, stop.
 - Pick the best of $(p - 1)$ models according to criterion, so now we have 2 variables in the model.

Continue adding 1 variable at a time in this way until we can no more variables improve the criterion. The final model is one with the best criterion after we *stop*; that is, no further improvement is possible.

Note: Much faster than brute force as the maximum number of models to fit is:

$$p + (p - 1) + \cdots + 2 + 1 = \sum_{i=1}^p i = \frac{p(p + 1)}{2}$$

which is $\mathcal{O}(p^2)$ compared to $\mathcal{O}(2^p)$ for all possible regressions.

(iii) Backward direction: remove one variable at a time to model.

- Start with model that has p predictors.
- Fit p models that result from removing one variable from the regression; that is, each one has $(p - 1)$ variables.
- Pick the best of p models according to criterion.
 - Eliminate that variable x_j from model.
 - Fit $(p - 1)$ models that remove x_j and one other variable from model.
 - Pick best of $(p - 1)$ models (2 variables removed).

Continue removing 1 variable at a time in this way until we can no more variables improve the criterion. Same computational complexity as forward selection.

(iv) Forward-backwards (allows individual variables to be both added/removed)

- Start as in forward selection
- If we have k variables in model:
 - Backwards: fit k models with $(k - 1)$ variables. If any of these improve criterion, remove the variable.
 - Forwards: fit $(p - k)$ models with $(k + 1)$ variables. If any of these improve criterion, add that variable.
- These are the basic “stepwise” selection models to get a “good” (useful) model.
- Many other have sophisticated procedures available. For example, stochastic search, lasso.
- We’ve assumed that $n > p$ because otherwise $(X^\top X)$ is not invertible. More specialized methods needed if number of predictors is larger than sample size.