

Ideas:

- Calculate the MLE for a function by hand.
- For MLR, interpreting coefficients and R^2 values.
 - Potentially, why are estimates for the same feature different when more predictors are introduced?
- Identifying Non-Identifiability

In []:

Module 2: Peer Reviewed Assignment

Outline:

The objectives for this assignment:

1. Mathematically derive the values of $\hat{\beta}_0$ and $\hat{\beta}_1$
2. Enhance our skills with linear regression modeling.
3. Learn the uses and limitations of RSS, ESS, TSS and R^2 .
4. Analyze and interpret nonidentifiability.

General tips:

1. Read the questions carefully to understand what is being asked.
2. This work will be reviewed by another human, so make sure that you are clear and concise in what your explanations and answers.

Problem 1: Maximum Likelihood Estimates (MLEs)

Consider the simple linear regression model $Y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$ for $i = 1, \dots, n$, $\varepsilon_i \sim N(0, \sigma^2)$. In the videos, we showed that the least squares estimator in matrix-vector form is

$\hat{\beta} = (\beta_0, \beta_1)^T = (X^T X)^{-1} X^T \mathbf{Y}$. In this problem, you will derive the least squares estimators for simple linear regression without (explicitly) using linear algebra.

Least squares requires that we minimize

$$f(\mathbf{x}; \beta_0, \beta_1) = \sum_{i=1}^n \left(Y_i - [\beta_0 + \beta_1 x_i] \right)^2$$

over β_0 and β_1 .

1. (a) Taking Derivatives

Find the partial derivative of $f(\mathbf{x}; \beta_0, \beta_1)$ with respect to β_0 , and the partial derivative of $f(\mathbf{x}; \beta_0, \beta_1)$ with respect to β_1 . Recall that the partial derivative with respect to x of a multivariate function $h(x, y)$ is calculated by taking the derivative of h with respect to x while treating y constant.

$$\begin{aligned}\frac{\partial f}{\partial \beta_0} &= -2 \sum \left(Y_i - [\beta_0 + \beta_1 x_i] \right) \\ \frac{\partial f}{\partial \beta_1} &= -2 \sum \left(Y_i - [\beta_0 + \beta_1 x_i] \right) x_i\end{aligned}$$

1. (b) Solving for $\hat{\beta}_0$ and $\hat{\beta}_1$

Use 1. (a) to find the minimizers, $\hat{\beta}_0$ and $\hat{\beta}_1$, of f . That is, set each partial derivative to zero and solve for β_0 and β_1 . In particular, show

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(Y_i - \bar{Y})}{\sum_{i=1}^n (x_i - \bar{x})^2} \quad \text{and} \quad \hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{x}$$

First, we find $\hat{\beta}_0$:

$$\begin{aligned}-2 \sum \left(Y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i \right) &\stackrel{set}{=} 0 \implies \sum \left(Y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i \right) = 0 \\ \implies \sum Y_i - \sum \hat{\beta}_0 - \hat{\beta}_1 \sum x_i &= 0 \implies n\bar{Y} - n\hat{\beta}_0 - \hat{\beta}_1 n\bar{x} = 0 \\ &\implies \hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{x},\end{aligned}$$

where $\hat{\beta}_1$ is found below. Note that we added hats once we set the derivative to zero. Now, let's find $\hat{\beta}_1$.

$$\begin{aligned}-2 \sum \left(Y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i \right) x_i &\stackrel{set}{=} 0 \implies \sum \left(Y_i x_i - \hat{\beta}_0 x_i - \hat{\beta}_1 x_i^2 \right) = 0 \\ \implies \sum Y_i x_i - \hat{\beta}_0 \sum x_i - \hat{\beta}_1 \sum x_i^2 &= 0 \implies \sum Y_i x_i - \hat{\beta}_0 n\bar{x} - \hat{\beta}_1 \sum x_i^2 = 0 \\ &\implies \sum Y_i x_i - (\bar{Y} - \hat{\beta}_1 \bar{x}) n\bar{x} - \hat{\beta}_1 \sum x_i^2 = 0 \\ &\implies \sum Y_i x_i - n\bar{Y}\bar{x} + \hat{\beta}_1 n\bar{x}^2 - \hat{\beta}_1 \sum x_i^2 = 0 \\ &\implies \sum Y_i x_i - n\bar{Y}\bar{x} - \hat{\beta}_1 \left(\sum x_i^2 - n\bar{x}^2 \right) = 0 \\ &\implies \hat{\beta}_1 \left(\sum x_i^2 - n\bar{x}^2 \right) = \sum Y_i x_i - n\bar{Y}\bar{x} \\ &\implies \hat{\beta}_1 = \frac{\sum Y_i x_i - n\bar{Y}\bar{x}}{\sum x_i^2 - n\bar{x}^2}\end{aligned}$$

The final form of $\hat{\beta}_1$ is found using the substitutions in problem 1 (b) and (c).

Problem 2: Oh My Goodness of Fit!

In the US, public schools have been slowly increasing class sizes over the last 15 years [https://stats.oecd.org/Index.aspx?DataSetCode=EDU_CLASS (https://stats.oecd.org/Index.aspx?DataSetCode=EDU_CLASS)]. The general cause for this is because it saves money to have more kids per teacher. But how much money does it save? Let's use some of our new regression skills to try and figure this out. Below is an explanation of the variables in the dataset.

Variables/Columns:

School

Per-Pupil Cost (Dollars)

Average daily Attendance

Average Monthly Teacher Salary (Dollars)

Percent Attendance

Pupil/Teacher ratio

Data Source: E.R. Enlow (1938). "Do Small Schools Mean Large Costs?," Peabody Journal of Education, Vol. 16, #1, pp. 1-11

```
In [41]: library(RCurl) #a package that includes the function getURL(), which allows for reading data from github.
url = getURL("https://raw.githubusercontent.com/bzaharatos/-Statistical-Modeling-for-Data-Science-Applications/master/Modern%20Regression%20Analysis%20/Datasets/school.csv")
school.data = read.csv(text = url, sep = ",", header = FALSE); school.data = school.data[, -7]
names(school.data) = c("school", "cost", "avg.attendance", "avg.salary", "pct.attendance", "pup.tch.ratio")
head(school.data)
summary(school.data)
```

school	cost	avg.attendance	avg.salary	pct.attendance	pup.tch.ratio
Calhoun	108.57	219.1	161.79	89.86	23.0
CapitolView	70.00	268.9	136.37	92.44	29.4
Connally	49.04	161.7	106.86	92.01	29.4
Couch	71.51	422.1	147.17	91.60	29.2
Crew	61.08	440.6	146.24	89.32	36.3
Davis	105.21	139.4	159.79	86.51	22.6

```

      school      cost      avg.attendance      avg.salary
Calhoun   : 1  Min.    : 49.04  Min.    : 139.4  Min.    :106.7
CapitolView: 1  1st Qu.: 59.36  1st Qu.: 282.9  1st Qu.:137.3
Connally   : 1  Median : 64.96  Median : 422.1  Median :146.8
Couch      : 1  Mean    : 67.91  Mean    : 433.4  Mean    :143.9
Crew       : 1  3rd Qu.: 73.27  3rd Qu.: 531.2  3rd Qu.:152.8
Davis      : 1  Max.    :108.57  Max.    :1065.4  Max.    :166.1
(Other)    :37
pct.attendance  pup.tch.ratio
Min.    :86.47  Min.    :22.60
1st Qu.:89.79  1st Qu.:29.20
Median :91.01  Median :30.70
Mean    :90.81  Mean    :30.77
3rd Qu.:92.03  3rd Qu.:32.80
Max.    :94.72  Max.    :36.70

```

2. (a) Create a model

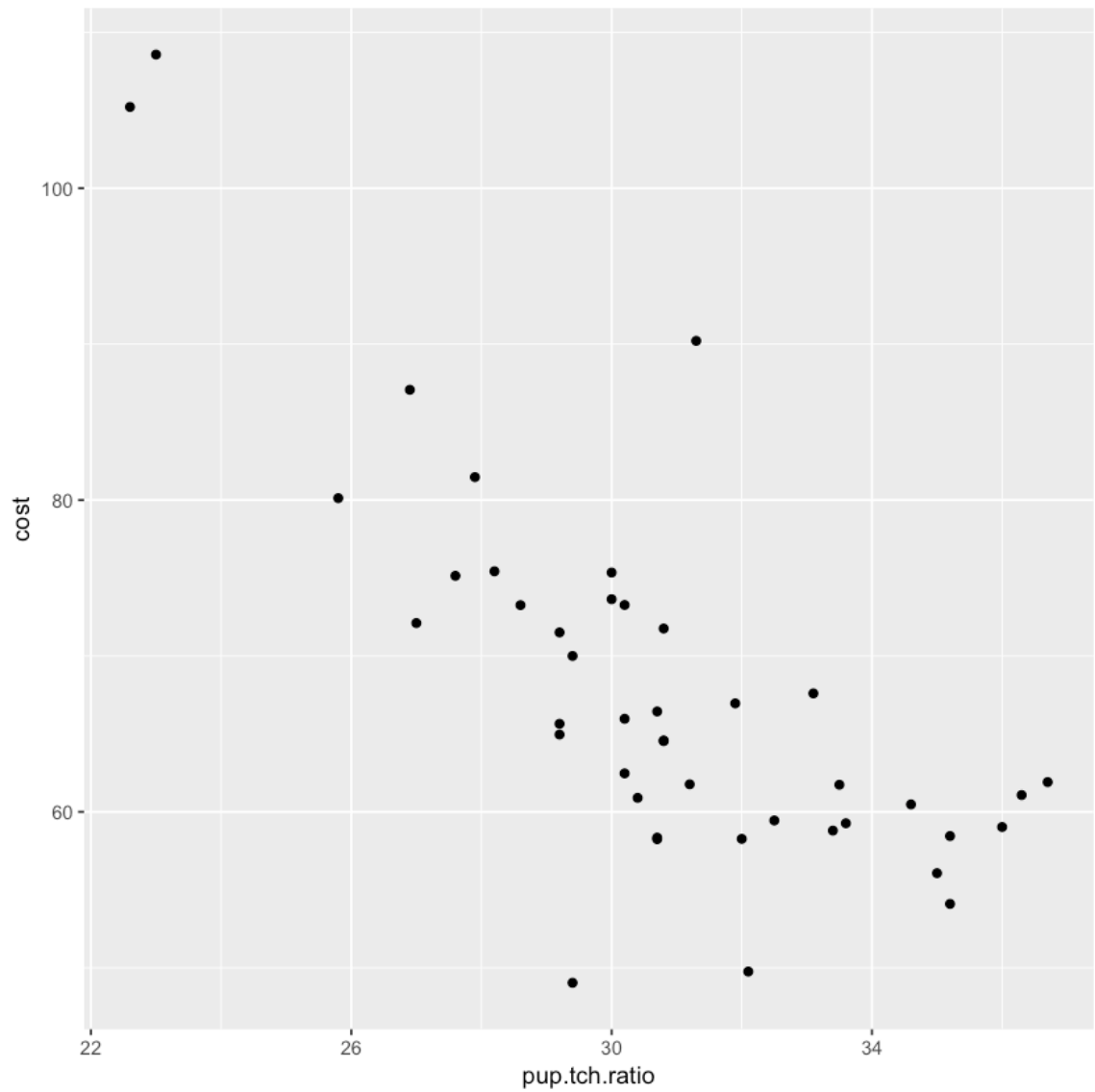
Begin by creating two figures for your model. The first with `pup.tch.ratio` on the x-axis and `cost` on the y-axis. The second with `avg.salary` on the x-axis and `cost` on the y-axis. Does there appear to be a relation between these two predictors and the response?

Then fit a multiple linear regression model with `cost` as the response and `pup.tch.ratio` and `avg.salary` as predictors.

```
In [44]: library(ggplot2)
ggplot(school.data, aes(x = pup.tch.ratio, y = cost)) +
  geom_point()

ggplot(school.data, aes(x = avg.salary, y = cost)) +
  geom_point()

lm_mlr = lm(cost ~ pup.tch.ratio + avg.salary, data = school.data)
summary(lm_mlr)
```



Call:
lm(formula = cost ~ pup.tch.ratio + avg.salary, data = school.data)

Residuals:

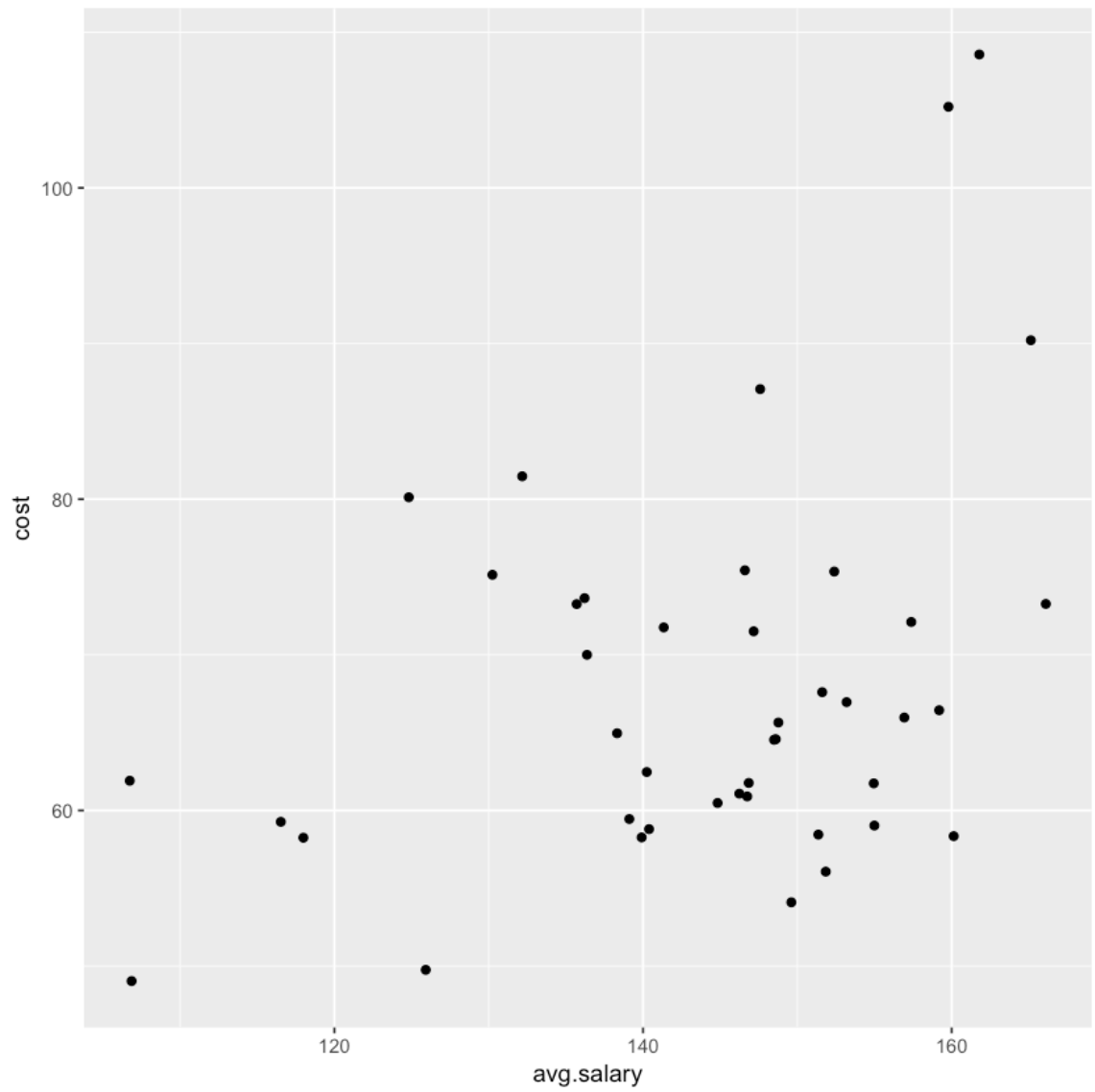
	Min	1Q	Median	3Q	Max
	-13.8290	-5.2752	-0.8332	3.8253	19.6986

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)	
(Intercept)	120.23756	17.73230	6.781	3.79e-08	***
pup.tch.ratio	-2.82585	0.37714	-7.493	3.90e-09	***
avg.salary	0.24061	0.08396	2.866	0.0066	**

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

Residual standard error: 7.721 on 40 degrees of freedom
Multiple R-squared: 0.6372, Adjusted R-squared: 0.6191
F-statistic: 35.13 on 2 and 40 DF, p-value: 1.559e-09



2. (b) RSS, ESS and TSS

In the code block below, manually calculate the RSS, ESS and TSS for your MLR model. Print the results.

```
In [81]: n = dim(school.data)[1]
p = 2
ess = sum((fitted(lm_mlr) - mean(school.data$cost))^2);
rss = sum(residuals(lm_mlr)^2);
tss = with(school.data, sum((cost - mean(cost))^2));

sigma = sqrt(rss/(n - p - 1)); sigma

cat(paste("The explained sum of squares is", round(ess,2), "."),
    paste("The residual sum of squares is ", round(rss,2),
"."),
    paste("The total sum of squares is ", round(tss,2), "."),
    paste("The estimate of the error standard deviation is", ro
und(sigma,2), "."),
    sep = "\n"
)
```

7.72107009767512

The explained sum of squares is 4188.57 .
The residual sum of squares is 2384.6 .
The total sum of squares is 6573.17 .
The estimate of the error standard deviation is 7.72 .

2. (c) Are you Squared?

Using the values from **2.b**, calculate the R^2 value for your model. Check your results with those produced from the `summary()` statement of your model.

In words, describe what this value means for your model.

```
In [49]: r2 = 1-rss/tss;
cat(paste("The coefficient of determination is", round(r2,2), "."))
```

The coefficient of determination is 0.64 .

64% of the variability in cost can be explained by the Pupil/Teacher ratio and the Average Monthly Teacher Salary .

2. (d) Conclusions

Describe at least two advantages and two disadvantages of the R^2 value.

- R^2 can be high (close to one) when the model does not fit properly.
- R^2 can be low (close to zero) when the model does fit properly.

We should consider R^2 only when we have a model that is appropriate for the data (e.g., linearity has been satisfied, we have roughly the correct predictors, etc.).

Problem 3: Identifiability

This problem might require some outside-of-class research if you haven't taken a linear algebra/matrix methods course.

Matrices and vectors play an important role in linear regression. Let's review some matrix theory as it might relate to linear regression.

Consider the system of linear equations

$$Y_i = \beta_0 + \sum_{j=1}^p \beta_j x_{i,j} + \varepsilon_i,$$

for $i = 1, \dots, n$, where n is the number of data points (measurements in the sample), and $j = 1, \dots, p$, where

1. $p + 1$ is the number of parameters in the model.
2. Y_i is the i^{th} measurement of the *response variable*.
3. $x_{i,j}$ is the i^{th} measurement of the j^{th} *predictor variable*.
4. ε_i is the i^{th} *error term* and is a random variable, often assumed to be $N(0, \sigma^2)$.
5. $\beta_j, j = 0, \dots, p$ are *unknown parameters* of the model. We hope to estimate these, which would help us characterize the relationship between the predictors and response.

3. (a) MLR Matrix Form

Write the equation above in matrix vector form. Call the matrix including the predictors X , the vector of Y_i s \mathbf{Y} , the vector of parameters β , and the vector of error terms ε . (This is more LaTeX practice than anything else...)

$$\mathbf{Y} = X\beta + \varepsilon,$$

$$\text{where } \mathbf{Y} = (Y_1, \dots, Y_n)^T, \quad X = \begin{bmatrix} 1 & X_{1,1} & X_{1,2} & \dots & X_{1,p} \\ 1 & X_{2,1} & X_{2,2} & \dots & X_{2,p} \\ \vdots & \vdots & \vdots & \dots & \vdots \\ 1 & X_{n,1} & X_{n,2} & \dots & X_{n,p} \end{bmatrix}, \quad \beta = (\beta_1, \dots, \beta_n)^T \text{ and}$$

$$\varepsilon = (\varepsilon_1, \dots, \varepsilon_n)^T$$

3. (b) Properties of this matrix

In lecture, we will find that the OLS estimator for β in MLR is $\hat{\beta} = (X^T X)^{-1} X^T Y$. Use this knowledge to answer the following questions:

1. What condition must be true about the columns of X for the "Gram" matrix $X^T X$ to be invertible?
2. What does this condition mean in practical terms, i.e., does X contain a deficiency or redundancy?
3. Suppose that the number of measurements (n) is less than the number of model parameters ($p + 1$). What does this say about the invertibility of $X^T X$? What does this mean on a practical level?
4. What is true about $\hat{\beta}$ if $X^T X$ is not invertible?

1. For $X^T X$ to be invertible, the columns of X must be linearly independent. That means that no column of X --i.e., no measured predictor--can be written as a linear combination of other columns. This implies that $n > (p + 1)$.
2. If we've measured a predictor that is simply a linear combination of others, that means that that predictor is not adding any new information that's not already contained in the other predictors. Imagine a simple case: X_1 is a predictor of measured weights in pounds, and X_2 is a predictor of measured weights in kilograms. Thus, $X_1 = 2.2046X_2$. Measuring X_1 doesn't give any new information.
3. This implies that we have more parameters than data/measurements, and thus $X^T X$ will not be invertible.
4. The formula for $\hat{\beta} = (X^T X)^{-1} X^T Y$ is derived from the "normal equations":

$$(X^T X)\beta = X^T Y.$$

The normal equations have a unique solutions if and only if $X^T X$ is invertible. If it's not invertible, then either the normal equations have no solutions or infinitely many solutions.

Problem 4: Downloading...

The following [data \(https://dasl.datadescription.com/datafile/downloading/\)](https://dasl.datadescription.com/datafile/downloading/) were collected to see if time of day made a difference on file download speed. A researcher placed a file on a remote server and then proceeded to download it at three different time periods of the day. They downloaded the file 48 times in all, 16 times at each Time of Day (`time`), and recorded the Time in seconds (`speed`) that the download took.

4. (a) Initial Observations

The `downloading` data is loaded in and cleaned for you. Using `ggplot`, create a boxplot of `speed` vs. `time`. Make some basic observations about the three categories.

```
In [50]: # Load in the Suntan data and format it
library(RCurl) #a package that includes the function getURL(), which
allows for reading data from github.
url = getURL("https://raw.githubusercontent.com/bzaharatos/-Statistical-Modeling-for-Data-Science-Applications/master/Modern%20Regression%20Analysis%20/Datasets/downloading.txt")
downloading = read.csv(text = url, sep = "\t")
names(downloading) = c("time", "speed")
# Change the types of brand and form to categories, instead of real
numbers
downloading$time = as.factor(downloading$time)
summary(downloading)
```

	time	speed
Early (7AM)	:16	Min. : 68.0
Evening (5 PM)	:16	1st Qu.:129.8
Late Night (12 AM):16		Median :198.0
		Mean :193.2
		3rd Qu.:253.0
		Max. :367.0

```
In [52]: summary(lm(speed ~ time, data = downloading))
boxplot(speed ~ time, data = downloading)
```

Call:

```
lm(formula = speed ~ time, data = downloading)
```

Residuals:

	Min	1Q	Median	3Q	Max
	-83.312	-34.328	-5.187	26.250	103.625

Coefficients:

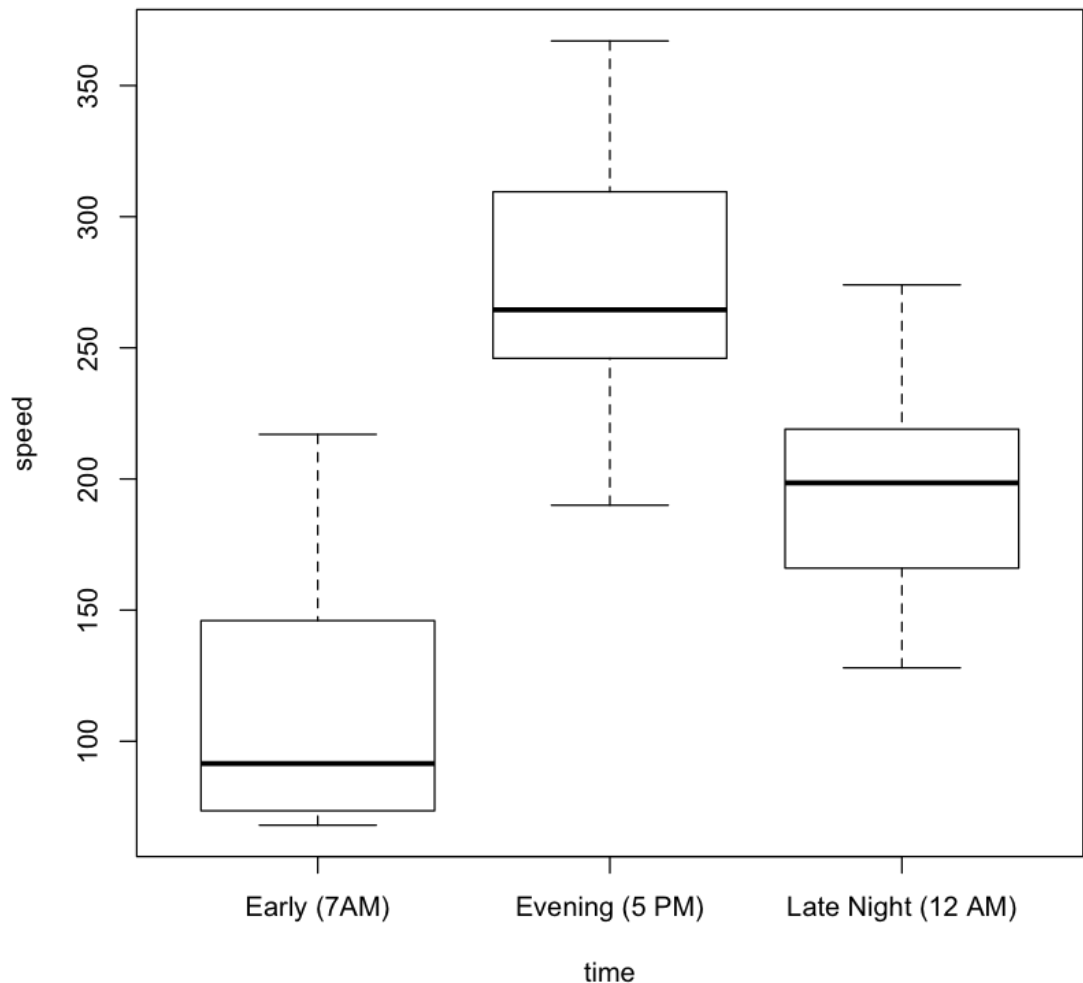
	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	113.38	11.79	9.619	1.73e-12 ***
timeEvening (5 PM)	159.94	16.67	9.595	1.87e-12 ***
timeLate Night (12 AM)	79.69	16.67	4.781	1.90e-05 ***

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

Residual standard error: 47.15 on 45 degrees of freedom

Multiple R-squared: 0.6717, Adjusted R-squared: 0.6571

F-statistic: 46.03 on 2 and 45 DF, p-value: 1.306e-11



It appears that, on average, download speeds are highest in the evening, followed by late night, with the morning having the lowest download speeds.

4. (b) How would we model this?

Fit a regression to these data that uses `speed` as the response and `time` as the predictor. Print the summary. Notice that the result is actually *multiple* linear regression, not simple linear regression. The model being used here is:

$$Y_i = \beta_0 + \beta_1 X_{i,1} + \beta_2 X_{i,2} + \varepsilon_i$$

where

1. $X_{i,1} = 1$ if the i^{th} download is made in the evening (5 pm).
2. $X_{i,2} = 1$ if the i^{th} download is made at night (12 am).

Note: If $X_{i,1} = 0$ and $X_{i,2} = 0$, then the i^{th} download is made in the morning (7am).

To confirm this is the model being used, write out the explicit equation for your model - using the parameter estimates from part (a) - and print out it's design matrix.

```
In [69]: X = model.matrix(lm(speed ~ time, data = downloading))
         head(X)
```

(Intercept)	timeEvening (5 PM)	timeLate Night (12 AM)
1	0	0
1	0	0
1	0	0
1	0	0
1	0	0
1	0	0

$$Y_i = 113.38 + 159.94X_{i,1} + 79.69X_{i,2} + \varepsilon_i$$

4. (c) Only two predictors?

We have three categories, but only two predictors. Why is this the case? To address this question, let's consider the following model:

$$Y_i = \beta_0 + \beta_1 X_{i,1} + \beta_2 X_{i,2} + \beta_3 X_{i,3} + \varepsilon_i$$

where

1. $X_{i,1} = 1$ if the i^{th} download is made in the evening (5 pm).
2. $X_{i,2} = 1$ if the i^{th} download is made at night (12 am).
3. $X_{i,3} = 1$ if the i^{th} download is made in the morning (7 am).

Construct a design matrix to fit this model to the response, `speed`. Determine if something is wrong with it. Hint: Analyze the design matrix.

```
In [74]: X_wrong = matrix(NA, nrow = dim(downloading), ncol = 4)

X_wrong[,1] = 1;
X_wrong[,2] = X[,2]
X_wrong[,3] = X[,3]
X_wrong[,4] = ifelse(downloading$time == "Early (7AM)", 1, 0)

solve(t(X_wrong)%*%X_wrong)
```

```
Error in solve.default(t(X_wrong) %*% X_wrong): system is computatio
nally singular: reciprocal condition number = 9.25186e-18
Traceback:
```

1. solve(t(X_wrong) %*% X_wrong)
2. solve.default(t(X_wrong) %*% X_wrong)

Notice that `X_wrong` is singular, which means that one of the columns of `X` is dependent and we don't need it.

4. (d) Interpretation

Interpret the coefficients in the model from **4.b**. In particular:

1. What is the difference between the mean download speed at 7am and the mean download speed at 5pm?
2. What is the mean download speed (in seconds) in the morning?
3. What is the mean download speed (in seconds) in the evening?
4. What is the mean download speed (in seconds) at night?

```
In [75]: summary(lm(speed ~ time, data = downloading))
```

Call:

```
lm(formula = speed ~ time, data = downloading)
```

Residuals:

	Min	1Q	Median	3Q	Max
	-83.312	-34.328	-5.187	26.250	103.625

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)	
(Intercept)	113.38	11.79	9.619	1.73e-12	***
timeEvening (5 PM)	159.94	16.67	9.595	1.87e-12	***
timeLate Night (12 AM)	79.69	16.67	4.781	1.90e-05	***

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

Residual standard error: 47.15 on 45 degrees of freedom

Multiple R-squared: 0.6717, Adjusted R-squared: 0.6571

F-statistic: 46.03 on 2 and 45 DF, p-value: 1.306e-11

1. The difference between the mean download speed at 7am and the mean download speed at 5pm is 159.94.
2. The mean download speed (in seconds) in the morning is 113.38.
3. The mean download speed (in seconds) in the evening is $113.38 + 159.94 = 273.32$.
4. The mean download speed (in seconds) at night is $113.38 + 79.69 = 193.07$.

```
In [80]: #To Confirm
```

```
mean(downloading$speed[downloading$time == "Early (7AM)"])  
mean(downloading$speed[downloading$time == "Evening (5 PM)"])  
mean(downloading$speed[downloading$time == "Late Night (12 AM)"])
```

113.375

273.3125

193.0625