

Week 4 - Homework

STAT 420, Summer 2023, D. Unger

Directions

Students are encouraged to work together on homework. However, sharing, copying or providing any part of a homework solution or code is an infraction of the University's rules on Academic Integrity. Any violation will be punished as severely as possible.

- Be sure to remove this section if you use this `.Rmd` file as a template.
 - You may leave the questions in your final document.
-

Exercise 1 (Using `1m`)

For this exercise we will use the data stored in [nutrition-2018.csv](#). It contains the nutritional values per serving size for a large variety of foods as calculated by the USDA in 2018. It is a cleaned version totaling 5956 observations and is current as of April 2018.

The variables in the dataset are:

- `ID`
- `Desc` - short description of food
- `Water` - in grams
- `Calories`
- `Protein` - in grams
- `Fat` - in grams
- `Carbs` - carbohydrates, in grams
- `Fiber` - in grams
- `Sugar` - in grams
- `Calcium` - in milligrams
- `Potassium` - in milligrams
- `Sodium` - in milligrams
- `VitaminC` - vitamin C, in milligrams
- `Chol` - cholesterol, in milligrams
- `Portion` - description of standard serving size used in analysis

(a) Fit the following multiple linear regression model in R. Use `Calories` as the response and `Fat`, `Sugar`, and `Sodium` as predictors.

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

Here,

- Y_i is **Calories**.
- x_{i1} is **Fat**.
- x_{i2} is **Sugar**.
- x_{i3} is **Sodium**.

Use an F -test to test the significance of the regression. Report the following:

- The null and alternative hypotheses
- The value of the test statistic
- The p-value of the test
- A statistical decision at $\alpha = 0.01$
- A conclusion in the context of the problem

When reporting these, you should explicitly state them in your document, not assume that a reader will find and interpret them from a large block of **R** output.

(b) Output only the estimated regression coefficients. Interpret all $\hat{\beta}_j$ coefficients in the context of the problem.

(c) Use your model to predict the number of **Calories** in a Filet-O-Fish. According to [McDonald's publicized nutrition facts](#), the Filet-O-Fish contains 18g of fat, 5g of sugar, and 580mg of sodium.

(d) Calculate the standard deviation, s_y , for the observed values in the Calories variable. Report the value of s_e from your multiple regression model. Interpret both estimates in the context of this problem.

(e) Report the value of R^2 for the model. Interpret its meaning in the context of the problem.

(f) Calculate a 90% confidence interval for β_2 . Give an interpretation of the interval in the context of the problem.

(g) Calculate a 95% confidence interval for β_0 . Give an interpretation of the interval in the context of the problem.

(h) Use a 99% confidence interval to estimate the mean Calorie content of a food with 15g of fat, 0g of sugar, and 260mg of sodium, which is true of a medium order of McDonald's french fries. Interpret the interval in context.

(i) Use a 99% prediction interval to predict the Calorie content of a Crunchy Taco Supreme, which has 11g of fat, 2g of sugar, and 340mg of sodium according to [Taco Bell's publicized nutrition information](#). Interpret the interval in context.

Exercise 2 (More 1m for Multiple Regression)

For this exercise we will use the data stored in [goalies17.csv](#). It contains career data for goaltenders in the National Hockey League during the first 100 years of the league from the 1917-1918 season to the 2016-2017 season. It holds the 750 individuals who played at least one game as goalie over this timeframe. The variables in the dataset are:

- **Player** - Player's Name (those followed by * are in the Hall of Fame as of 2017)
- **First** - First year with game recorded as goalie
- **Last** - Last year with game recorded as goalie
- **Active** - Number of seasons active in the NHL
- **GP** - Games Played
- **GS** - Games Started

- W - Wins
- L - Losses (in regulation)
- TOL - Ties and Overtime Losses
- GA - Goals Against
- SA - Shots Against
- SV - Saves
- SV_PCT - Save Percentage
- GAA - Goals Against Average
- SO - Shutouts
- PIM - Penalties in Minutes
- MIN - Minutes

For this exercise we will consider three models, each with Wins as the response. The predictors for these models are:

- Model 1: Goals Against, Saves
- Model 2: Goals Against, Saves, Shots Against, Minutes, Shutouts
- Model 3: All Available

After reading in the data but prior to any modeling, you should clean the data set for this exercise by removing the following variables: `Player`, `GS`, `L`, `TOL`, `SV_PCT`, and `GAA`.

(a) Use an F -test to compare Models 1 and 2. Report the following:

- The null hypothesis
- The value of the test statistic
- The p-value of the test
- A statistical decision at $\alpha = 0.05$
- The model you prefer

(b) Use an F -test to compare Model 3 to your preferred model from part (a). Report the following:

- The null hypothesis
- The value of the test statistic
- The p-value of the test
- A statistical decision at $\alpha = 0.05$
- The model you prefer

(c) Use a t -test to test $H_0 : \beta_{SV} = 0$ vs $H_1 : \beta_{SV} \neq 0$ for the model you preferred in part (b). Report the following:

- The value of the test statistic
- The p-value of the test
- A statistical decision at $\alpha = 0.05$

Exercise 3 (Regression without lm)

For this exercise we will once again use the `Ozone` data from the `mlbench` package. The goal of this exercise is to fit a model with `ozone` as the response and the remaining variables as predictors.

```
data(Ozone, package = "mlbench")
Ozone = Ozone[, c(4, 6, 7, 8)]
colnames(Ozone) = c("ozone", "wind", "humidity", "temp")
Ozone = Ozone[complete.cases(Ozone), ]
```

(a) Obtain the estimated regression coefficients **without** the use of `lm()` or any other built-in functions for regression. That is, you should use only matrix operations. Store the results in a vector `beta_hat_no_lm`. To ensure this is a vector, you may need to use `as.vector()`. Return this vector as well as the results of `sum(beta_hat_no_lm ^ 2)`.

(b) Obtain the estimated regression coefficients **with** the use of `lm()`. Store the results in a vector `beta_hat_lm`. To ensure this is a vector, you may need to use `as.vector()`. Return this vector as well as the results of `sum(beta_hat_lm ^ 2)`.

(c) Use the `all.equal()` function to verify that the results are the same. You may need to remove the names of one of the vectors. The `as.vector()` function will do this as a side effect, or you can directly use `unname()`.

(d) Calculate s_e without the use of `lm()`. That is, continue with your results from (a) and perform additional matrix operations to obtain the result. Output this result. Also, verify that this result is the same as the result obtained from `lm()`.

(e) Calculate R^2 without the use of `lm()`. That is, continue with your results from (a) and (d), and perform additional operations to obtain the result. Output this result. Also, verify that this result is the same as the result obtained from `lm()`.

Exercise 4 (Regression for Prediction)

For this exercise use the `Auto` dataset from the `ISLR` package. Use `?Auto` to learn about the dataset. The goal of this exercise is to find a model that is useful for **predicting** the response `mpg`. We remove the `name` variable as it is not useful for this analysis. (Also, this is an easier to load version of data from the textbook.)

```
# load required package, remove "name" variable
library(ISLR)
```

```
## Warning: package 'ISLR' was built under R version 4.2.1
```

```
Auto = subset(Auto, select = -c(name))
dim(Auto)
```

```
## [1] 392 8
```

When evaluating a model for prediction, we often look at RMSE. However, if we both fit the model with all the data as well as evaluate RMSE using all the data, we're essentially cheating. We'd like to use RMSE as a measure of how well the model will predict on *unseen* data. If you haven't already noticed, the way we had been using RMSE resulted in RMSE decreasing as models became larger.

To correct for this, we will only use a portion of the data to fit the model, and then we will use leftover data to evaluate the model. We will call these datasets **train** (for fitting) and **test** (for evaluating). The definition of RMSE will stay the same

$$\text{RMSE}(\text{model}, \text{data}) = \sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2}$$

where

- y_i are the actual values of the response for the given data.
- \hat{y}_i are the predicted values using the fitted model and the predictors from the data.

However, we will now evaluate it on both the **train** set and the **test** set separately. So each model you fit will have a **train** RMSE and a **test** RMSE. When calculating **test** RMSE, the predicted values will be found by predicting the response using the **test** data with the model fit using the **train** data. *Test data should never be used to fit a model.*

- Train RMSE: Model fit with *train* data. Evaluate on **train** data.
- Test RMSE: Model fit with *train* data. Evaluate on **test** data.

Set a seed of 22, and then split the **Auto** data into two datasets, one called **auto_trn** and one called **auto_tst**. The **auto_trn** data frame should contain 290 randomly chosen observations. The **auto_tst** data will contain the remaining observations. Hint: consider the following code:

```
set.seed(22)
auto_trn_idx = sample(1:nrow(Auto), 290)
```

Fit a total of five models using the training data.

- One must use all possible predictors.
- One must use only **displacement** as a predictor.
- The remaining three you can pick to be anything you like. One of these should be the *best* of the five for predicting the response.

For each model report the **train** and **test** RMSE. Arrange your results in a well-formatted markdown table. Argue that one of your models is the best for predicting the response.

Exercise 5 (Simulating Multiple Regression)

For this exercise we will simulate data from the following model:

$$Y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \beta_3 x_{i3} + \beta_4 x_{i4} + \beta_5 x_{i5} + \epsilon_i$$

Where $\epsilon_i \sim N(0, \sigma^2)$. Also, the parameters are known to be:

- $\beta_0 = 2$
- $\beta_1 = -0.75$
- $\beta_2 = 1.6$
- $\beta_3 = 0$
- $\beta_4 = 0$

- $\beta_5 = 2$
- $\sigma^2 = 25$

We will use samples of size $n = 40$.

We will verify the distribution of $\hat{\beta}_1$ as well as investigate some hypothesis tests.

(a) We will first generate the X matrix and data frame that will be used throughout the exercise. Create the following nine variables:

- **x0**: a vector of length n that contains all 1
- **x1**: a vector of length n that is randomly drawn from a normal distribution with a mean of 0 and a standard deviation of 2
- **x2**: a vector of length n that is randomly drawn from a uniform distribution between 0 and 4
- **x3**: a vector of length n that is randomly drawn from a normal distribution with a mean of 0 and a standard deviation of 1
- **x4**: a vector of length n that is randomly drawn from a uniform distribution between -2 and 2
- **x5**: a vector of length n that is randomly drawn from a normal distribution with a mean of 0 and a standard deviation of 2
- **X**: a matrix that contains **x0**, **x1**, **x2**, **x3**, **x4**, and **x5** as its columns
- **C**: the C matrix that is defined as $(X^T X)^{-1}$
- **y**: a vector of length n that contains all 0
- **sim_data**: a data frame that stores **y** and the **five predictor** variables. **y** is currently a placeholder that we will update during the simulation.

Report the sum of the diagonal of **C** as well as the 5th row of **sim_data**. For this exercise we will use the seed 420. Generate the above variables in the order listed after running the code below to set a seed.

```
set.seed(400)
sample_size = 40
```

(b) Create three vectors of length 2500 that will store results from the simulation in part (c). Call them **beta_hat_1**, **beta_3_pval**, and **beta_5_pval**.

(c) Simulate 2500 samples of size $n = 40$ from the model above. Each time update the **y** value of **sim_data**. Then use **lm()** to fit a multiple regression model. Each time store:

- The value of $\hat{\beta}_1$ in **beta_hat_1**
- The p-value for the two-sided test of $\beta_3 = 0$ in **beta_3_pval**
- The p-value for the two-sided test of $\beta_5 = 0$ in **beta_5_pval**

(d) Based on the known values of X , what is the true distribution of $\hat{\beta}_1$?

(e) Calculate the mean and variance of **beta_hat_1**. Are they close to what we would expect? Plot a histogram of **beta_hat_1**. Add a curve for the true distribution of $\hat{\beta}_1$. Does the curve seem to match the histogram?

(f) What proportion of the p-values stored in **beta_3_pval** is less than 0.10? Is this what you would expect?

(g) What proportion of the p-values stored in **beta_5_pval** is less than 0.01? Is this what you would expect?