

Assigned:  
January 20, 2018

Homework 1

Due:  
February 03, 2018

---

Please complete the assigned problems to the best of your abilities. Ensure that the work you do is entirely your own, external resources are only used as permitted by the instructor, and all allowed sources are given proper credit for non-original content.

## 1 Recitation Problems

These problems are to be found in: **Introduction to Statistical Learning, 7<sup>th</sup> Printing (Online Edition)** by *Gareth James, Daniela Witten, Trevor Hastie, Robert Tibshirani*.

### 1.1 Chapter 2

Problems: 1,2,4,6

## 2 Practicum Problems

These problems will primarily reference the *lecture materials and the examples given in class* using **R** and **CRAN**. It is suggested that a *RStudio* session be used for the programmatic components.

### 2.1 Problem 1

Load the *iris* sample dataset into **R** using a dataframe (it is a built-in dataset). Create a boxplot of each of the 4 features, and highlight the feature with the largest empirical *IQR*. Calculate the parametric standard deviation for each feature - do your results agree with the empirical values? Use the *ggplot2* library from **CRAN** to create a colored boxplot for each feature, with a box-whisker per flower species. Which flower type exhibits a significantly different *Petal Length/Width* once it is separated from the other classes?

### 2.2 Problem 2

Load the *trees* sample dataset into **R** using a dataframe (it is a built-in dataset), and produce a 5-number summary of each feature. Create a histogram of each variable - which variables appear to be normally distributed based on visual inspection? Do any variables exhibit positive or negative skewness? Install the *moments* library from **CRAN** use the *skewness* function to calculate the skewness of each variable. Do the values agree with the visual inspection?

### 2.3 Problem 3

Load the *auto-mpg* sample dataset from the UCI Machine Learning Repository (**auto-mpg.data**) into **R** using a dataframe (**Hint**: You will need to use *read.csv* with *url*, and set the appropriate values for **header**, **as.is**, and **sep**).

Assigned:  
January 20, 2018

Homework 1

Due:  
February 03, 2018

---

The *horsepower* feature has a few missing values with a ? - and will be treated as a string. Use the *as.numeric* casting function to obtain the column as a numeric vector, and replace all NA values with the median. How does this affect the value obtained for the mean vs the original mean when the records were ignored?