

# STATS 419 Survey of Multivariate Analysis

## Week 03 Assignment 02\_datasets

Justin Pickel  
[justin.pickel@wsu.edu](mailto:justin.pickel@wsu.edu)

Instructor: Monte J. Shaffer

18 September 2020

```
library(devtools)
local.path = "C:/_git_/WSU_STATS419_FALL2020/";
source(paste0(local.path, "functions/libraries.R"), local=T);
```

## 1 Martrix

Create the “rotate matrix” functions as described in lectures. Apply to the example “myMatrix”.

```
source( paste0(local.path, "functions/functions-matrix.R"), local=T );

myMatrix = matrix ( c (
  1, 0, 2,
  0, 3, 0,
  4, 0, 5
), nrow=3, byrow=T);
```

### 1.1 Matrix

```
rotateMatrix90(myMatrix)
```

```
##      [,1] [,2] [,3]
## [1,]    4    0    1
## [2,]    0    3    0
## [3,]    5    0    2
```

## 1.2 Matrix

```
rotateMatrix180(myMatrix)
```

```
##      [,1] [,2] [,3]
## [1,]    5    0    4
## [2,]    0    3    0
## [3,]    2    0    1
```

## 1.3 Matrix

```
rotateMatrix270(myMatrix)
```

```
##      [,1] [,2] [,3]
## [1,]    2    0    5
## [2,]    0    3    0
## [3,]    1    0    4
```

## 2 Iris graphic

Recreate the graphic for the IRIS Data Set using R. Same titles, same scales, same colors. `## Iris graphic`

```
library(datasets)
data("iris")

plot(iris[1:4],main="Iris Data (red=setosa,green=versicolor,blue=virginica)",col=c("black"),
pch=21,bg=c("red","#258425","blue")[iris$Species])
```

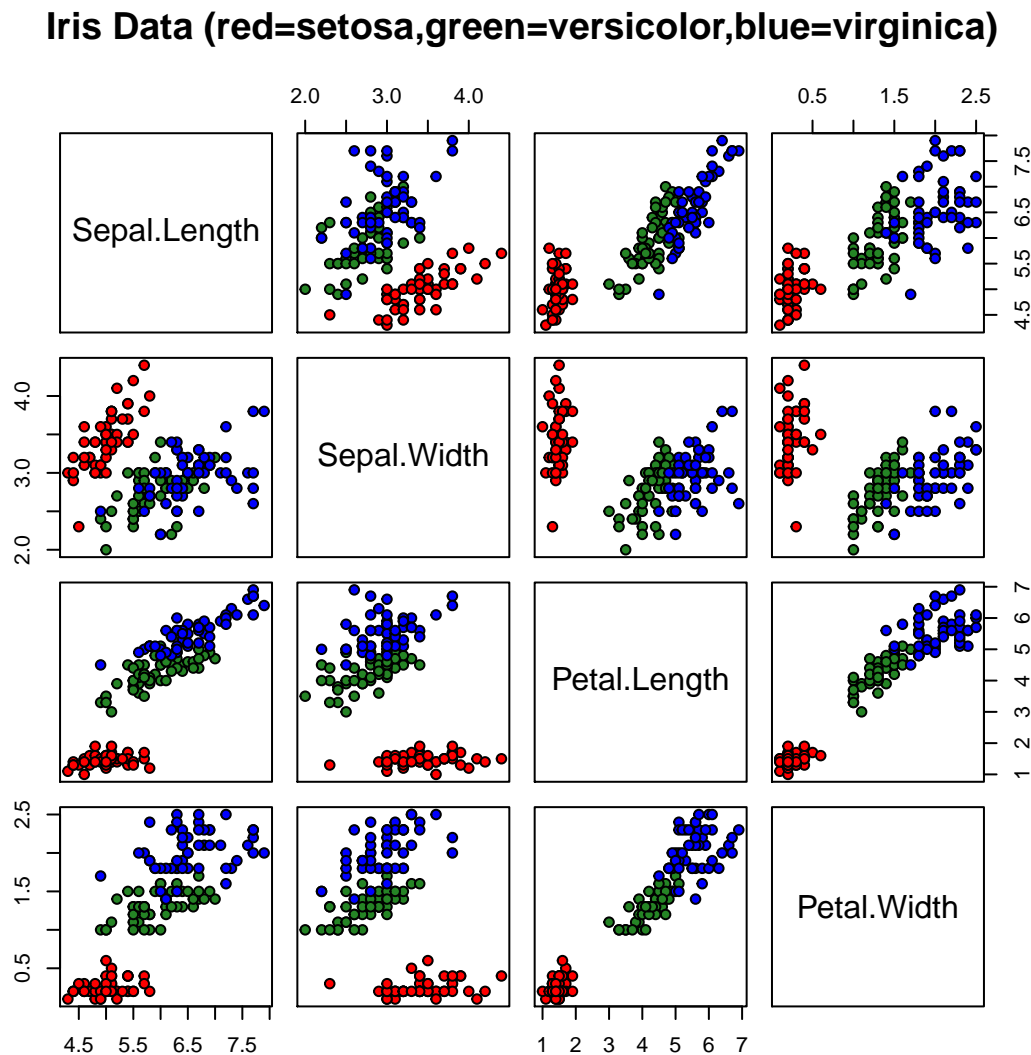


Figure 1: Iris data generated using ‘plot()’

### 3 Iris Summary

Write 2-3 sentences concisely defining the IRIS Data Set.

#### 3.1 Iris Summary

The Iris flower dataset contains 150 total data points, 50 data points for each of the three Iris flowers (setosa, virginica, versicolor). The data set records the measurement of four variables (Petal Length, Pedal Width, Sepal Width, Sepal length) for each of the 50 samples. The data set also contains a variable “class” indicating what species the measurements belongs to.

### 4 Clean Personality dataset

In the homework, for these tasks, report how many records your raw dataset had and how many records your clean dataset has.

## 4.1 Clean Personality dataset

```
source( paste0(local.path,"WEEK-03/functions/functions-CleanData.R"), local=T )

personality.raw = readfile(paste0(local.path,"datasets/personality-raw.txt"))

observations.raw = nrow(personality.raw)
observations.clean = nrow(cleanData(personality.raw))
```

The raw dataset has 838 records  
The clean dataset has 678 records

## 5 Functions and Z-score

Write functions for doSummary and sampleVariance and doMode ... test these functions in your homework on the "monte.shaffer@gmail.com" record from the clean dataset. Report your findings. For this "monte.shaffer@gmail.com" record, also create z-scores. Plot(x,y) where x is the raw scores for "monte.shaffer@gmail.com" and y is the z-scores from those raw scores. Include the plot in your assignment, and write 2 sentences describing what pattern you are seeing and why this pattern is present.

### 5.1 Functions doSummary(), doMode()

```
source( paste0(local.path,"WEEK-03/functions/functions-CleanData.R"), local=T )
source( paste0(local.path,"WEEK-03/functions/functions-summary.R"), local=T )

personality.raw = readfile(paste0(local.path,"datasets/personality-raw.txt"))
personality.clean = cleanData(personality.raw)
record = retrieveTopRecord(personality.clean)

doSummary(record) # returns summary data
```

```
##   length countNA mean median NaiveVar.sum_ NaiveVar.sumSq NaiveVar.variance
## 1      60      0 3.48   3.4      208.8      771.04      0.7528136
##   twoPassVar.sum twoPassVar.sum2 twoPassVar.variance    stDev mode
## 1      208.8      44.416      0.7528136 0.8676483 4.2
```

```
doMode(record) # calculates the mode of a vector.
```

```
## [1] 4.2
```

### 5.2 Z-scores

```
z.score = (record-mean(record))/sd(record)
plot(record,z.score,main="Z-score and raw score relationship",xlab="Raw score",ylab="Z-score")
```

## Z-score and raw score relationship

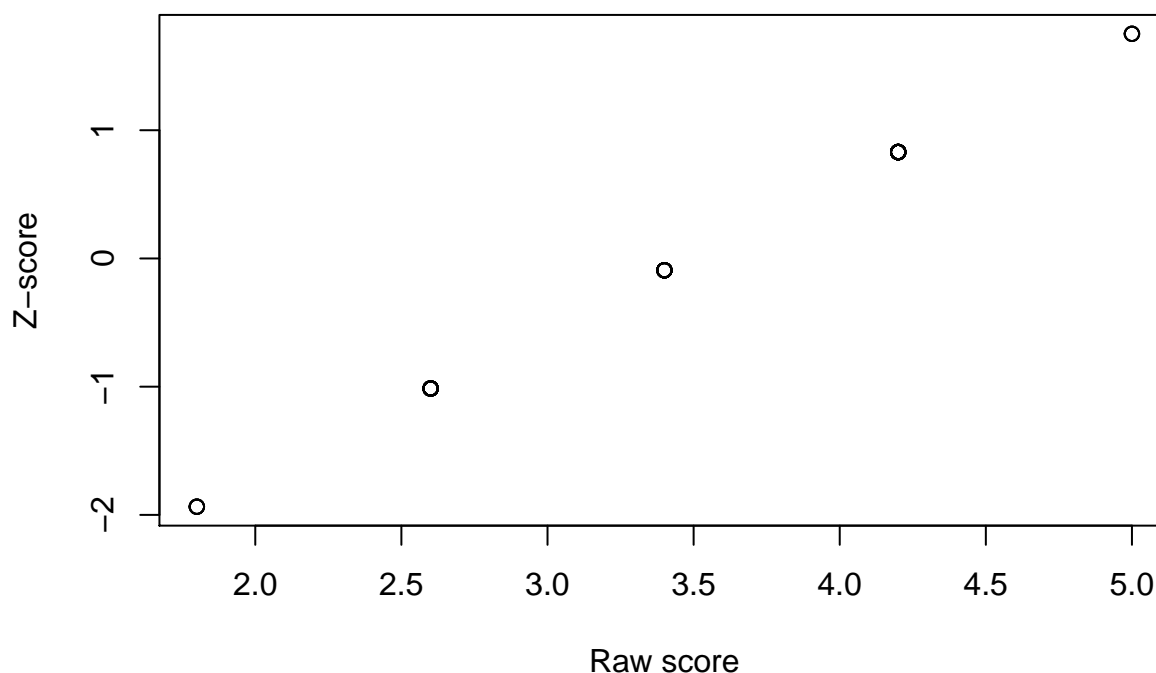


Figure 2: Relationship between z-scores and raw scores

## 6 Denzel Washington vs. Will Smith

Compare Will Smith and Denzel Washington. You will have to create a new variable `$millions.2000` that converts each movie's `$millions` based on the `$year` of the movie, so all dollars are in the same time frame. You will need inflation data from about 1980-2020 to make this work.

### 6.1 Side by Side comparison before conversion

```
source( paste0(local.path,"functions/functions-imdb.R"), local=T );

nmid = "nm0000226";
will = grabFilmsForPerson(nmid);

nmid = "nm0000243";
denzel = grabFilmsForPerson(nmid);

boxplot(denzel$movies.50$millions,will$movies.50$millions,
        names=c("Denzel Washington","Will Smith"),main="Movie Revenue",
        ylab= "Revenue (millions)",xlab = "Actors")
```

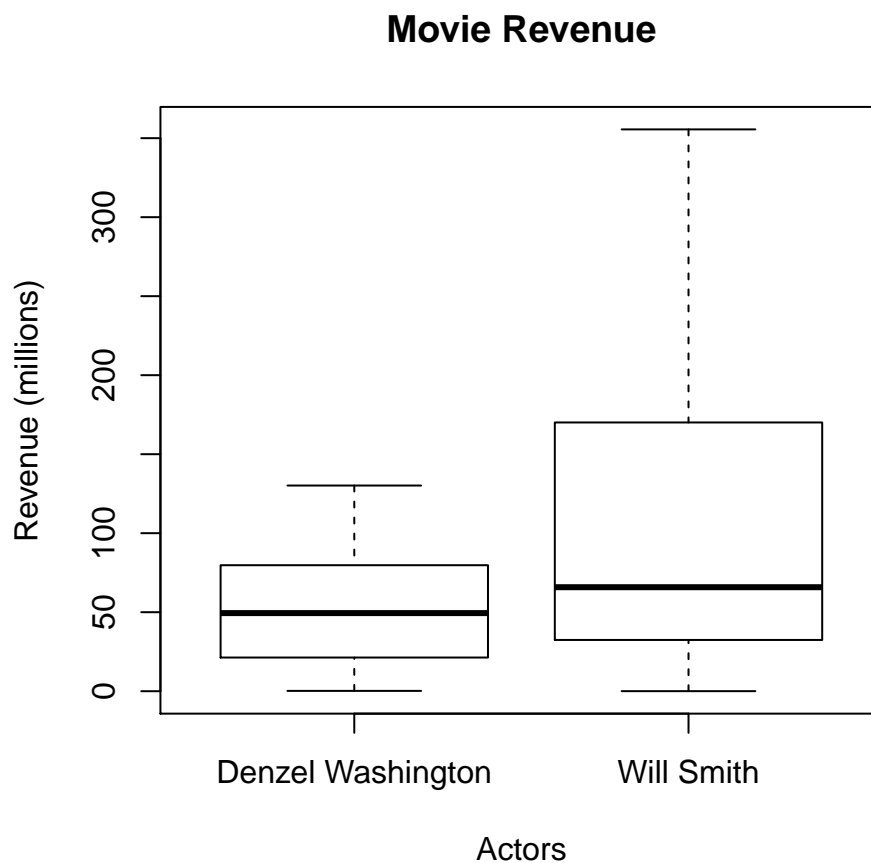


Figure 3: Side by Side comparison of Will Smith and Denzel Washington before conversion

## 6.2 Converting millions into year 2000 dollars

By converting all the revenue dollars from various movies across different years to a constant year allows us to compare apples to apples. If we do not convert the revenue to the same year our data will be skewed. I used inflation data from 1980 to 2020 by scraping the data from the web using code provided by Monte J. Shaffer to convert every year into the year 2000 dollars. After the conversion step we see when comparing Denzel Washington to Will Smith we can see most movies featuring Will Smith revenue higher than Denzel Washington. Will Smith also has a higher max revenue than Denzel Washington and the median revenue between the two actors is almost equal. We can also see that Will Smith has an outlier, meaning one movie generated way more revenue when compared to all his other top 50 movies (Figure 4).

```
source( paste0(local.path,"/WEEK-03/functions/functions-inflation.R"), local=T );

will.converted = convertMillions(will,2000)
denzel.converted = convertMillions(denzel,2000)
# show top 5 rows of Will Smith after converting to year 2000 dollars.
will.converted$movies.50[1:5,c(1,2,4,7,11,12)]
```

##	rank		title	year	genre	millions	millions.2000
## 1	1		I Am Legend	2007	Action, Adventure, Drama	256.39	212.9349
## 2	2		Suicide Squad	2016	Action, Adventure, Fantasy	325.10	233.2524
## 3	3		Independence Day	1996	Action, Adventure, Sci-Fi	306.17	336.0260
## 4	4		Men in Black	1997	Action, Adventure, Comedy	250.69	268.9646
## 5	5		I, Robot	2004	Action, Drama, Sci-Fi	144.80	131.9987

## 7 Denzel vs. Will Side by Side comparison of several Variables

Build side-by-side box plots on several of the variables (including #6) to compare the two movie stars. After each box plot, write 2+ sentences describing what you are seeing, and what conclusions you can logically make. You will need to review what the box plot is showing with the box portion, the divider in the box, and the whiskers.

### 7.1 Denzel vs. Will Side by Side comparison of revenue after conversion

Comparing the revenue generated from movies featuring Denzel Washington to Will Smith using a boxplot. We can see from Figure 4. the median for Will Smith is slightly higher than Denzel Washington. The actors have a median revenue of 56.5897 and 52.4536 million dollars, respectively. Fifty percent of the revenue for Denzel Washington is between 22.3942 and 70.0587 million dollars. Fifty percent of the revenue generated by Will Smith is between 30.4543 and 137.6564 million dollars. The max revenue for Denzel is 120.0872 while Will Smith has a max of 336.0260 million dollars. From these statistics we can conclude Will Smith has been featured in more films with higher revenue than Denzel Washington.

```
# Revenue in year 2000 dollars comparison
boxplot(denzel.converted$movies.50$millions.2000,will.converted$movies.50$millions.2000,
names=c("Denzel Washington","Will Smith"),main="Movie revenue in year 2000 dollars",
ylab= "Revenue (millions)",xlab = "Actors");
```

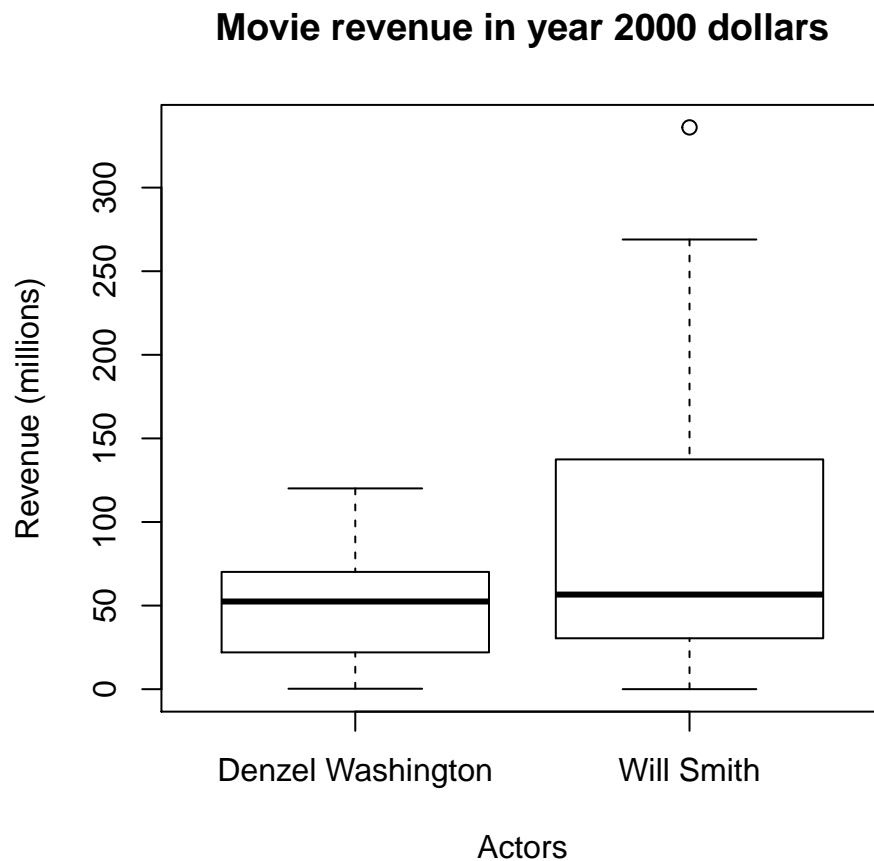


Figure 4: Side by Side comparison of Will and Denzel after conversion

## 7.2 Denzel vs. Will Side by Side comparison of ratings

Comparing the ratings from movies featuring Denzel Washington to Will Smith using a boxplot. We can see from Figure 5. the average ratings for films featuring Will Smith is slightly lower than Denzel Washington. The actors have an average rating of 6.326 and 6.852, respectively. Fifty percent of ratings data for Denzel Washington is between 6.525 and 7.3. Fifty percent of the ratings data for Will smith is between 6.0 and 6.875. The max rating for Denzel is a 7.8 while Will Smith has a max rating of 8.0. From these statistics we can conclude Denzel Washington has been in more movies with higher ratings than Will Smith has.

```
# Ratings comparison
boxplot(denzel.converted$movies.50$ratings,will.converted$movies.50$ratings,
names=c("Denzel Washington","Will Smith"),main="Movie ratings",
ylab= "Ratings",xlab = "Actors")
```



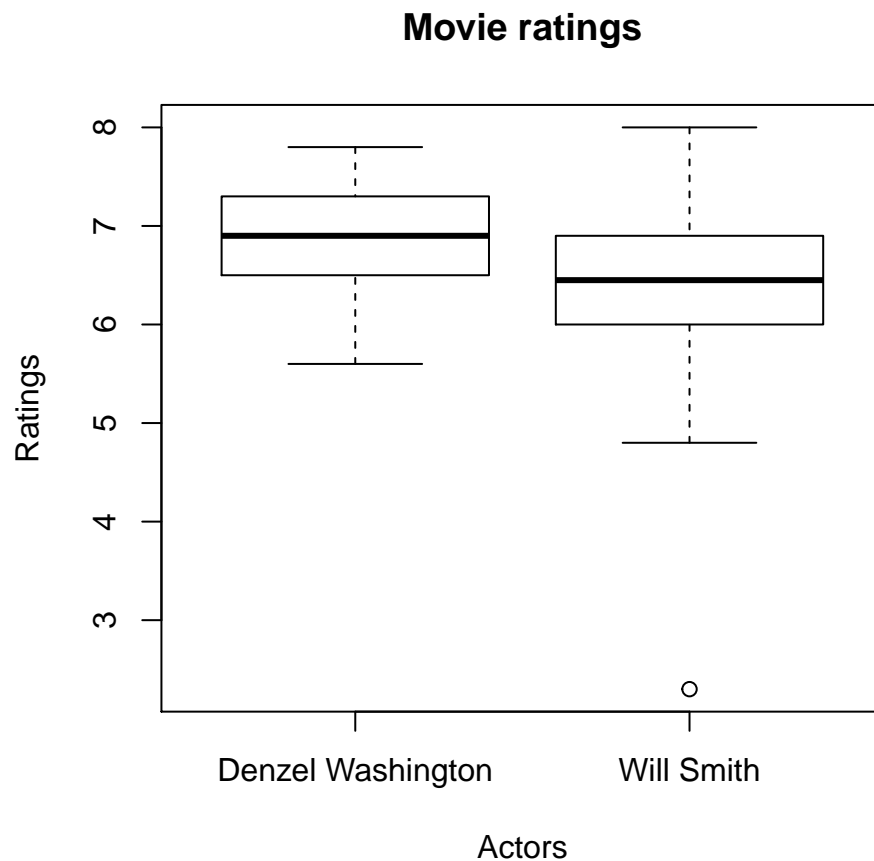


Figure 5: Side by Side comparison of ratings between Will and Denzel

### 7.3 Denzel vs. Will Side by Side comparison of runtime

Comparing the runtime from movies featuring Denzel Washington to Will Smith using a boxplot. We can see from Figure 6. the median for Will Smith is slightly lower than Denzel Washington. The actors have a median runtime of 113 and 118 minutes, respectively. Fifty percent of the runtime for Denzel Washington is between 106 and 129 minutes. Fifty percent of the runtime for Will smith is between 100.2 and 123 minutes. The max runtime for Denzel is 202 minutes while Will Smith has a max of 157 minutes. From these statistics we can conclude that Denzel Washington has been in movies with longer runtimes than movies featuring Will Smith.

```
# Runtime comparison
boxplot(denzel.converted$movies.50$minutes,will.converted$movies.50$minutes,
names=c("Denzel Washington","Will Smith"),main="Movie runtime minutes",
ylab= "Minutes",xlab = "Actors")
```

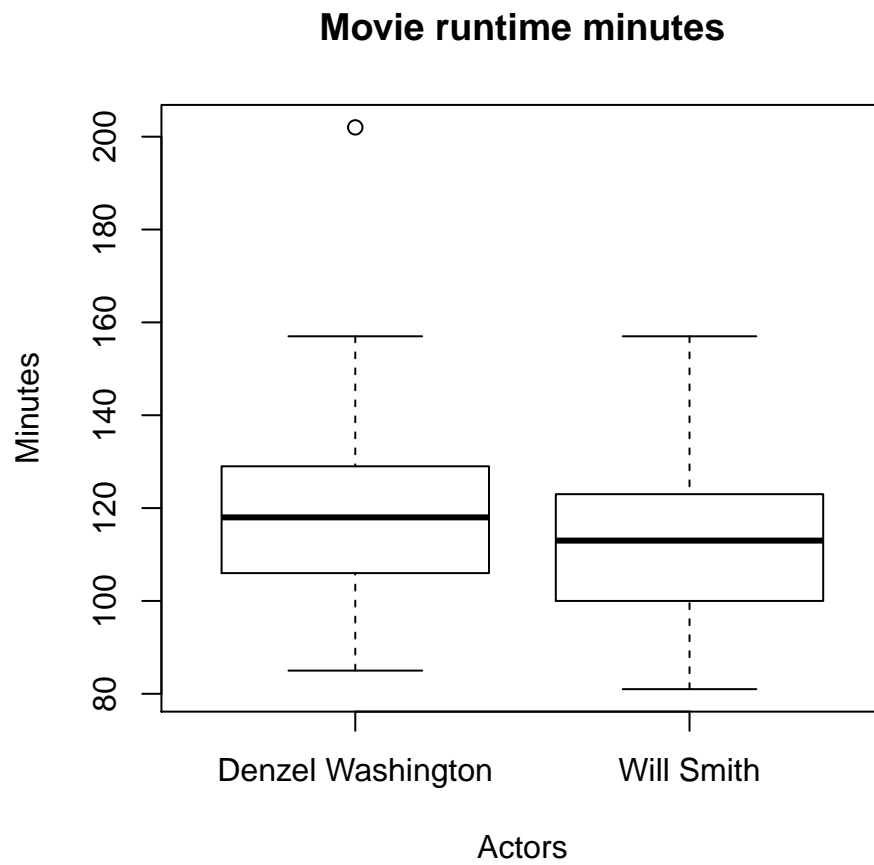


Figure 6: Side by Side comparison of runtime between Will and Denzel