

Homework Exercise 1

W203 K Iwasaki

January 9, 2016

W203 Statistics for Data Science

Unit 1 Homework

Exercise

Load the dataset found in the file, cars.csv.

```
setwd("C:/Users/K/Desktop/Berkeley/00_Academics/01_2017 Summer/W203 STATS/hw01")
data <- read.csv('cars.csv')

# check the data
head(data)
```

```
##      mpg cyl  disp  hp drat   wt  qsec vs am gear carb
## 1  21.0   6   160  110 3.90 2.620 16.46  0  1    4    4
## 2  21.0   6   160  110 3.90 2.875 17.02  0  1    4    4
## 3  22.8   4   108   93 3.85 2.320 18.61  1  1    4    1
## 4  21.4   6   258  110 3.08 3.215 19.44  1  0    3    1
## 5  18.7   8   360  175 3.15 3.440 17.02  0  0    3    2
## 6  18.1   6   225  105 2.76 3.460 20.22  1  0    3    1
```

```
nrow(data)
```

```
## [1] 25
```

```
?cars
```

```
## starting httpd help server ...
```

```
## done
```

1. What are the variables in the file?

```
str(data)
```

2. Find the mean, median, minimum, maximum, 1st quartile and 3rd quartile for the mpg variable.

3. Create a histogram of the mpg variable.

Distribution of Variable "mpg"



```
std = sd(data$mpg)
std
```

```
## [1] 6.047446
```

5. What is the variance of mpg variable?

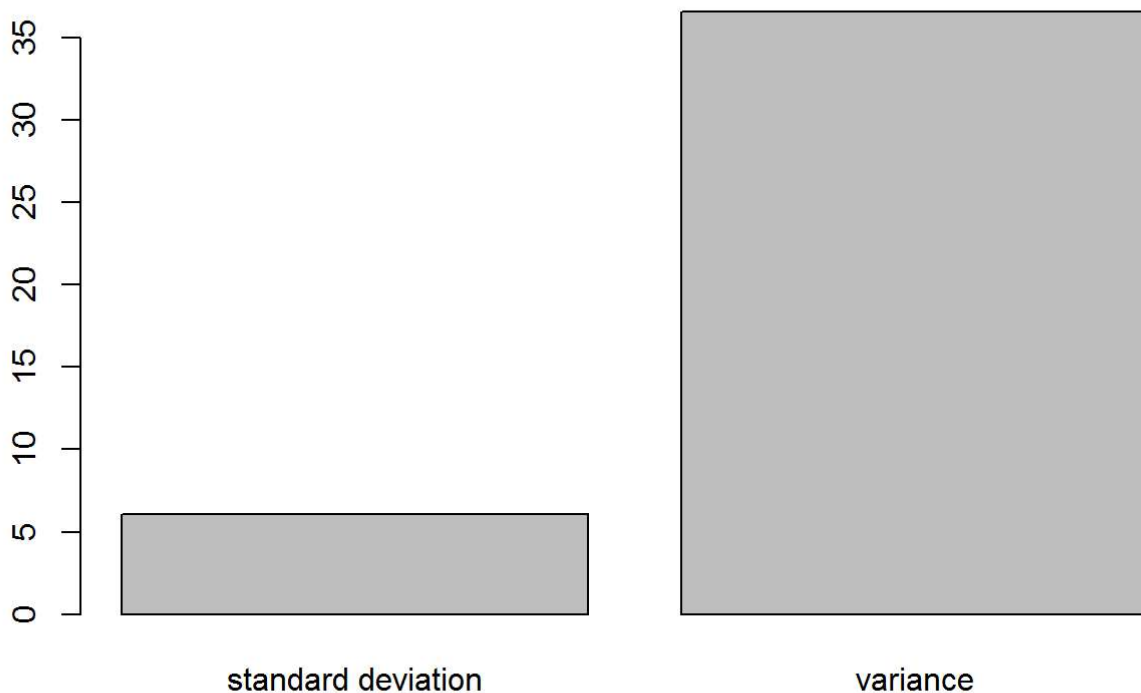
```
variance = var(data$mpg)
variance
```

```
## [1] 36.5716
```

6. What is the relationship of the standard deviation to the variance? Why does the standard deviation and variance of the mpg variable differ?

```
# create a table and compare the values visually
mpg_df <- data.frame(variable=c("standard deviation", "variance"), val = c(std, variance))

barplot(mpg_df$val, names.arg = mpg_df$variable)
```



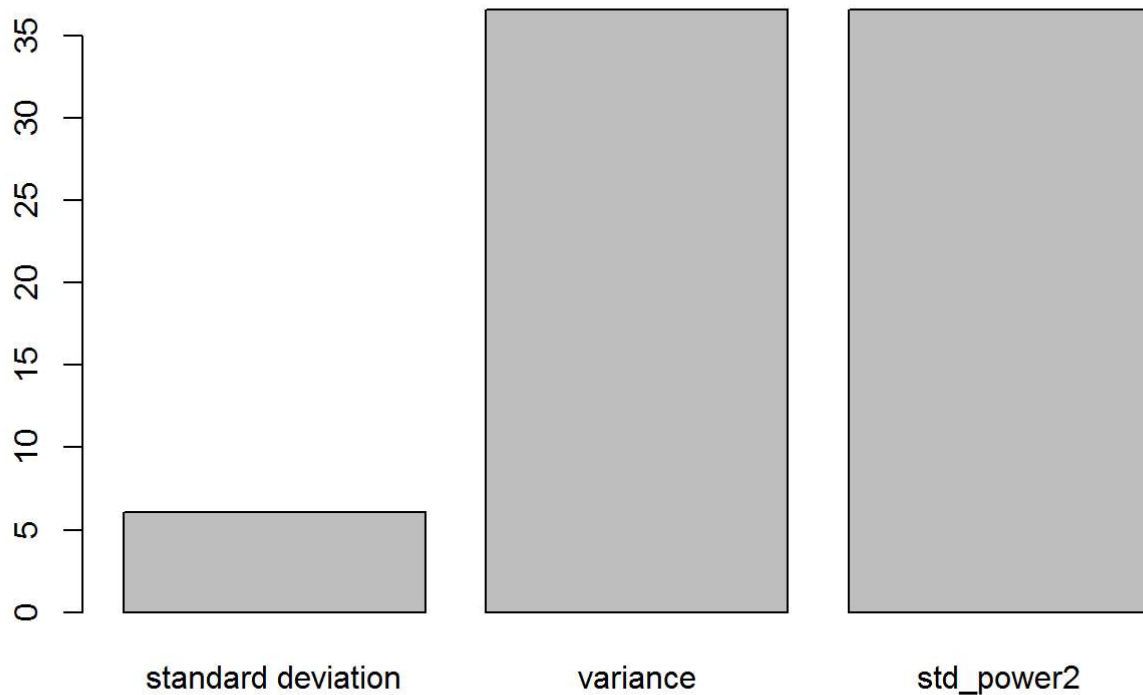
```
# The standard deviation is the square root of the variance by definition.
# Let's check if this is true here.

std_power2 <- std **2
check <- std_power2 == variance
check
```

```
## [1] TRUE
```

```
# viualize this
new.row1 <- data.frame(variable = "std_power2", val = std_power2)
mpg_df <- rbind(mpg_df, new.row1)

barplot(mpg_df$val, names.arg = mpg_df$variable)
```



7. How many data points are there for the cyl variable?

```
length(data$cyl)
```

```
## [1] 25
```

```
summary(data$cyl) # note there is two NAs
```

```
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.     NA's
##    4.000   4.000   6.000   6.261   8.000   8.000         2
```

8. What is the mean of the cyl variable?

```
mean(data$cyl, na.rm=T)
```

```
## [1] 6.26087
```