

# Census Project

1. Load required libraries. Install package `install.packages("caret")` Install package `install.packages("corrplot")`  
Install package `install.packages('Boruta')`

```
library(ggplot2)
library(corrplot)
```

```
## corrplot 0.84 loaded
```

```
library(Boruta)
```

```
## Loading required package: ranger
```

```
library(caret)
```

```
## Loading required package: lattice
```

2. Load census data.

```
setwd("c:/Ryerson University/Semester 4/ProjectCode")
loc<-getwd()
censusdata <- read.csv(file="census.csv",header=TRUE,sep="," , na.string = "?")
```

- 2.1. Divide the data into train and test data.

```
inTrain <- createDataPartition(y=censusdata$income, p= 0.75, list=FALSE)
training <- censusdata[inTrain,]
testing <- censusdata[-inTrain,]
```

3. Display dimensions, summary of data, names and structure of data.

```
data <- training
dim(data)
```

```
## [1] 24421    15
```

```
nrow(data)
```

```
## [1] 24421
```

```
ncol(data)
```

```
## [1] 15
```

```
dim(testing)
```

```
## [1] 8140    15
```

```
summary(data)
```

```
##      age      workclass      fnlwt
##  Min.   :17.00   Private      :16985   Min.    : 12285
##  1st Qu.:28.00   Self-emp-not-inc: 1947   1st Qu.: 118088
##  Median :37.00   Local-gov       : 1570   Median : 178530
##  Mean   :38.63   State-gov       :  972   Mean   : 190094
##  3rd Qu.:48.00   Self-emp-inc    :  836   3rd Qu.: 237272
##  Max.   :90.00   (Other)         :  734   Max.   :1455435
##                NA's           : 1377
##      education  education.num      marital.status
##  HS-grad      :7866   Min.    : 1.00   Divorced           : 3309
```

```
## Some-college:5427 1st Qu.: 9.00 Married-AF-spouse : 14
## Bachelors :4041 Median :10.00 Married-civ-spouse :11220
## Masters :1301 Mean :10.08 Married-spouse-absent: 316
## Assoc-voc :1023 3rd Qu.:12.00 Never-married : 8043
## 11th : 872 Max. :16.00 Separated : 790
## (Other) :3891 Widowed : 729
## occupation relationship race
## Prof-specialty :3117 Husband :9877 Amer-Indian-Eskimo: 239
## Craft-repair :3066 Not-in-family :6250 Asian-Pac-Islander: 774
## Exec-managerial:3054 Other-relative: 755 Black : 2336
## Adm-clerical :2873 Own-child :3799 Other : 204
## Sales :2746 Unmarried :2567 White :20868
## (Other) :8184 Wife :1173
## NA's :1381
## sex capital.gain capital.loss hours.per.week
## Female: 8093 Min. : 0 Min. : 0.00 Min. : 1.00
## Male :16328 1st Qu.: 0 1st Qu.: 0.00 1st Qu.:40.00
## Median : 0 Median : 0.00 Median :40.00
## Mean : 1066 Mean : 87.29 Mean :40.52
## 3rd Qu.: 0 3rd Qu.: 0.00 3rd Qu.:45.00
## Max. :99999 Max. :4356.00 Max. :99.00
##
## native.country income
## United-States:21865 <=50K:18540
## Mexico : 477 >50K : 5881
## Philippines : 145
## Germany : 94
## Canada : 92
## (Other) : 1303
## NA's : 445
```

```
names(data)
```

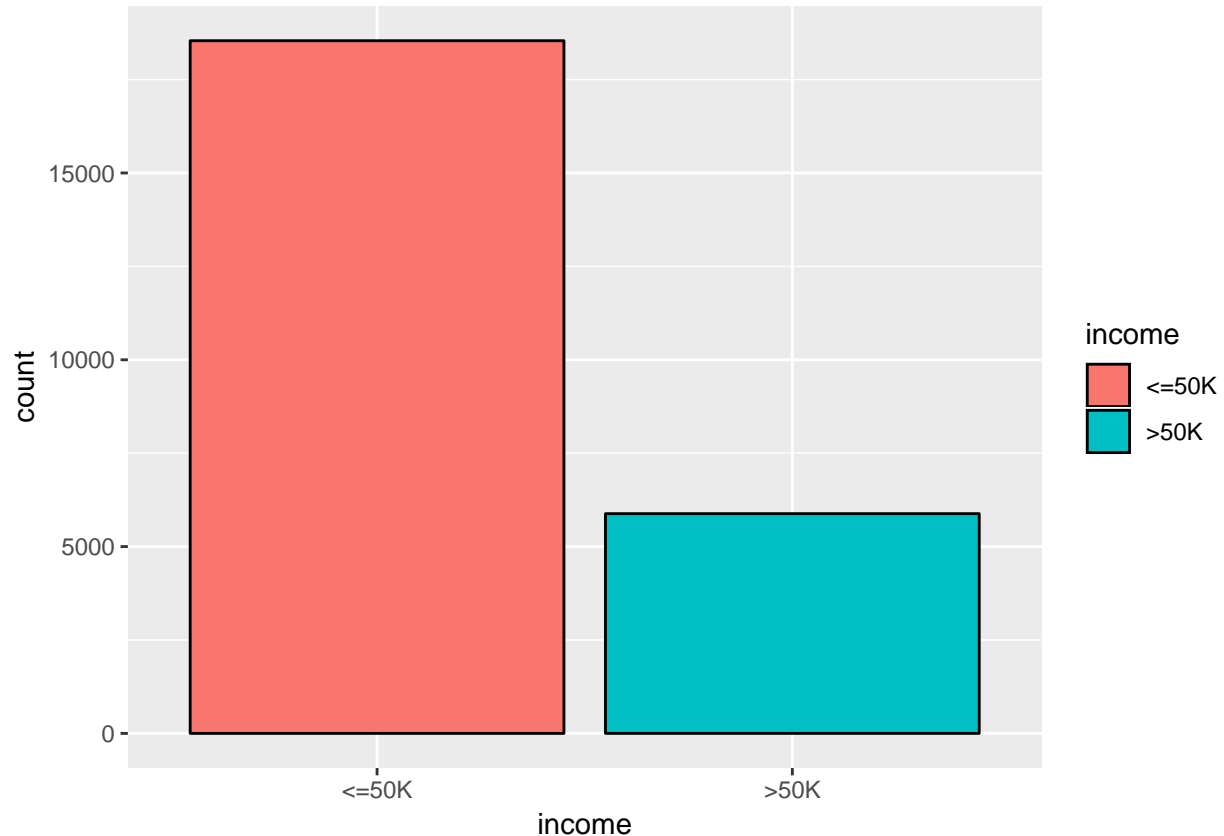
```
## [1] "age" "workclass" "fnlwgt" "education"
## [5] "education.num" "marital.status" "occupation" "relationship"
## [9] "race" "sex" "capital.gain" "capital.loss"
## [13] "hours.per.week" "native.country" "income"
```

```
str(data)
```

```
## 'data.frame': 24421 obs. of 15 variables:
## $ age : int 90 82 66 54 41 34 38 74 68 41 ...
## $ workclass : Factor w/ 8 levels "Federal-gov",...: NA 4 NA 4 4 4 4 7 1 4 ...
## $ fnlwgt : int 77053 132870 186061 140359 264663 216864 150601 88638 422013 70037 ...
## $ education : Factor w/ 16 levels "10th","11th",...: 12 12 16 6 16 12 1 11 12 16 ...
## $ education.num : int 9 9 10 4 10 9 6 16 9 10 ...
## $ marital.status: Factor w/ 7 levels "Divorced","Married-AF-spouse",...: 7 7 7 1 6 1 6 5 1 5 ...
## $ occupation : Factor w/ 14 levels "Adm-clerical",...: NA 4 NA 7 10 8 1 10 10 3 ...
## $ relationship : Factor w/ 6 levels "Husband","Not-in-family",...: 2 2 5 5 4 5 5 3 2 5 ...
## $ race : Factor w/ 5 levels "Amer-Indian-Eskimo",...: 5 5 3 5 5 5 5 5 5 5 ...
## $ sex : Factor w/ 2 levels "Female","Male": 1 1 1 1 1 1 2 1 1 2 ...
## $ capital.gain : int 0 0 0 0 0 0 0 0 0 0 ...
## $ capital.loss : int 4356 4356 4356 3900 3900 3770 3770 3683 3683 3004 ...
## $ hours.per.week: int 40 18 40 40 40 45 40 20 40 60 ...
## $ native.country: Factor w/ 41 levels "Cambodia","Canada",...: 39 39 39 39 39 39 39 39 39 NA ...
## $ income : Factor w/ 2 levels "<=50K",">50K": 1 1 1 1 1 1 1 2 1 2 ...
```

#### 4. Display Class Distributions.

```
# Imbalance data
result = summary(data$income)/nrow(data) * 100
ggplot(data=data,aes(income)) + geom_bar(aes(fill = income), color = "black")
```



result

```
##      <=50K      >50K
## 75.91827 24.08173
```

#### 5. Check and Cleaning missing values.

```
cat("Number of missing values in training set is:", sum(is.na(data)), "\n")
```

```
## Number of missing values in training set is: 3203
```

```
na_count <- sapply(data, function(y) sum(length(which(is.na(y)))))
na_count <- data.frame(na_count)
na_count
```

```
##           na_count
## age              0
## workclass      1377
## fnlwgt          0
## education       0
## education.num   0
## marital.status  0
## occupation     1381
## relationship    0
```

```
## race          0
## sex           0
## capital.gain  0
## capital.loss  0
## hours.per.week 0
## native.country 445
## income        0
```

```
nrow(data)
```

```
## [1] 24421
```

```
data <- na.omit(data)
nrow(data)
```

```
## [1] 22616
```

```
nrow(testing)
```

```
## [1] 8140
```

```
cat("Number of missing values in test set is:", sum(is.na(testing)), "\n")
```

```
## Number of missing values in test set is: 1059
```

```
na_count1 <- sapply(testing, function(y) sum(length(which(is.na(y)))))
na_count1
```

```
##          age      workclass      fnlwgt      education      education.num
##          0         459          0          0              0
## marital.status      occupation      relationship          race          sex
##          0         462          0          0              0
## capital.gain      capital.loss      hours.per.week      native.country      income
##          0          0          0          138              0
```

```
testingdata <- na.omit(testing)
nrow(testingdata)
```

```
## [1] 7546
```

5.1 Re-factoring the work class, occupation and native country after removing the NA values (exclude levels not required).

```
data$workclass <- factor(data$workclass)
data$occupation <- factor(data$occupation)
data$native.country <- factor(data$native.country)
```

5.1 Re-factoring the work class, occupation and native country after removing the NA values (exclude levels not required) for testing data also.

```
testingdata$workclass <- factor(testingdata$workclass)
testingdata$occupation <- factor(testingdata$occupation)
testingdata$native.country <- factor(testingdata$native.country)
```

## 6. Statistics of Numerical attributes

```
# statistics of numerical attributes
summary(data$age)
```

```
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##      17.00  28.00   37.00   38.49  47.00   90.00
```

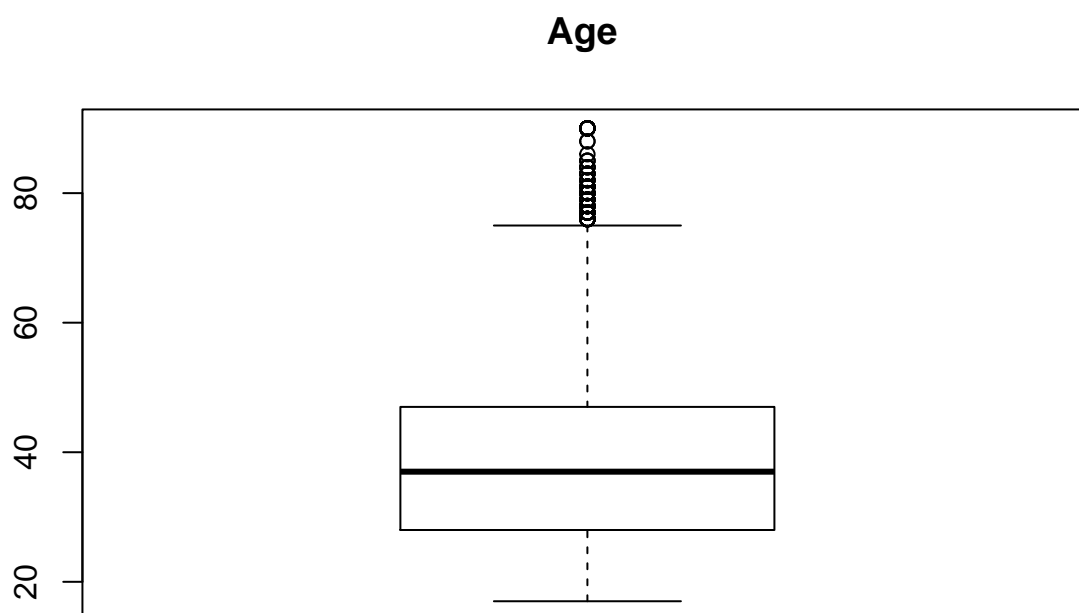
```
sd(data$age)
```

```
## [1] 13.15366
```

```
hist(data$age, main = "Distribution of Age",xlab = "Age of an Individual" ,col ="blue")
```



```
boxplot(data$age,main="Age ")
```



```
summary(data$education.num)
```

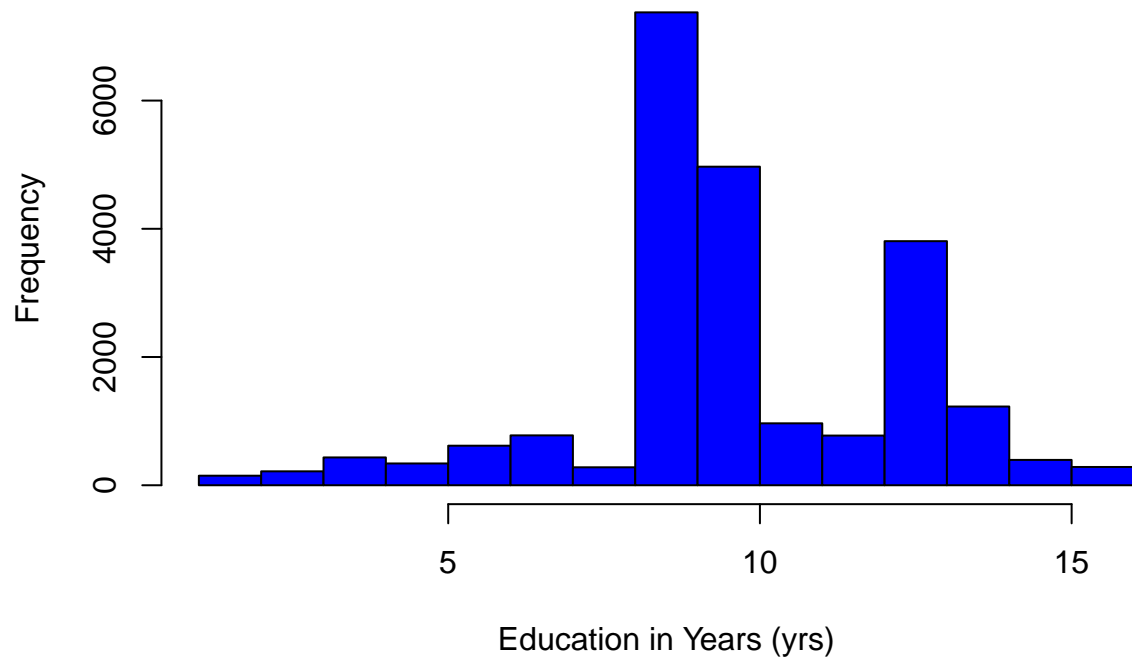
```
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##      1.00   9.00   10.00   10.12   13.00   16.00
```

```
sd(data$education.num)
```

```
## [1] 2.558803
```

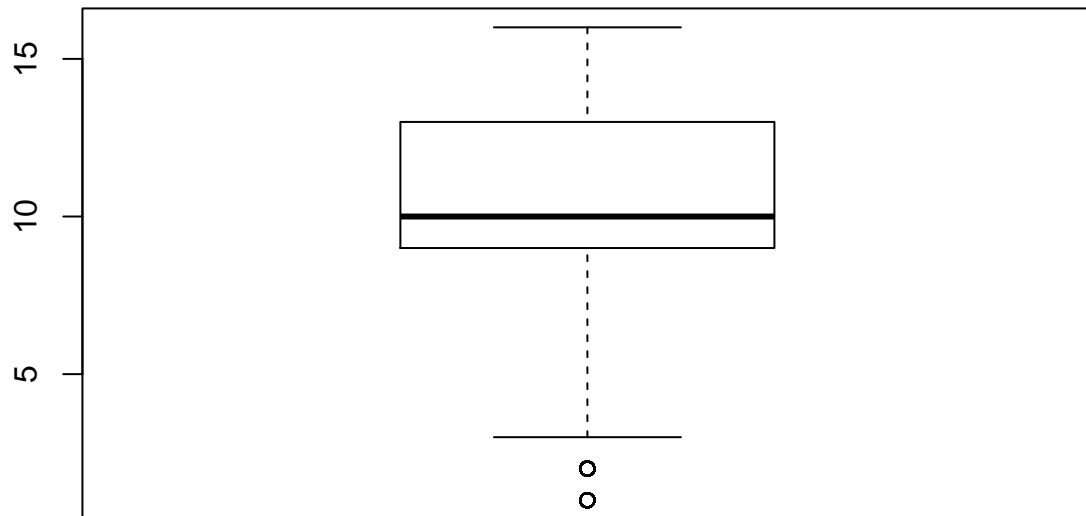
```
hist(data$education.num,main = "Distribution of Education in years",xlab="Education in Years (yrs)",col
```

## Distribution of Education in years



```
boxplot(data$education.num,main="Distribution of Education")
```

## Distribution of Education



```
summary(data$capital.gain)
```

```
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##         0         0         0   1088         0   99999
```

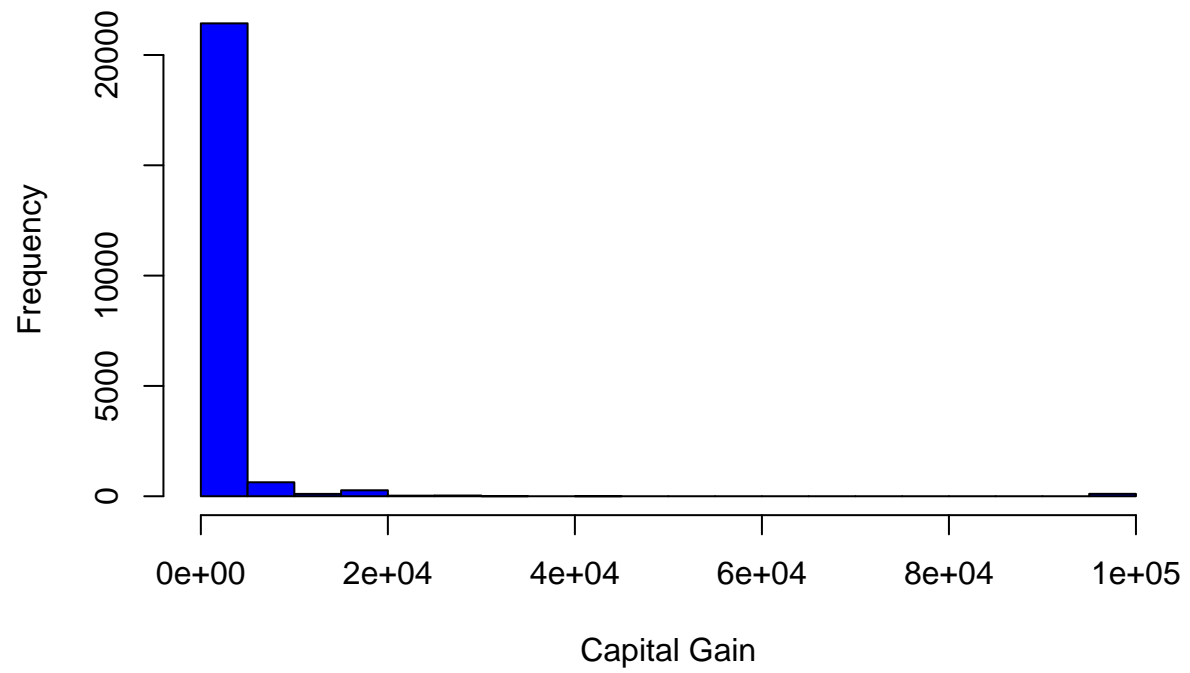
```
sd(data$capital.gain)
```

```
## [1] 7351.035
```

```
hist(data$capital.gain,main = "Distribution of Capital Gain",xlab="Capital Gain",col = "blue")
```

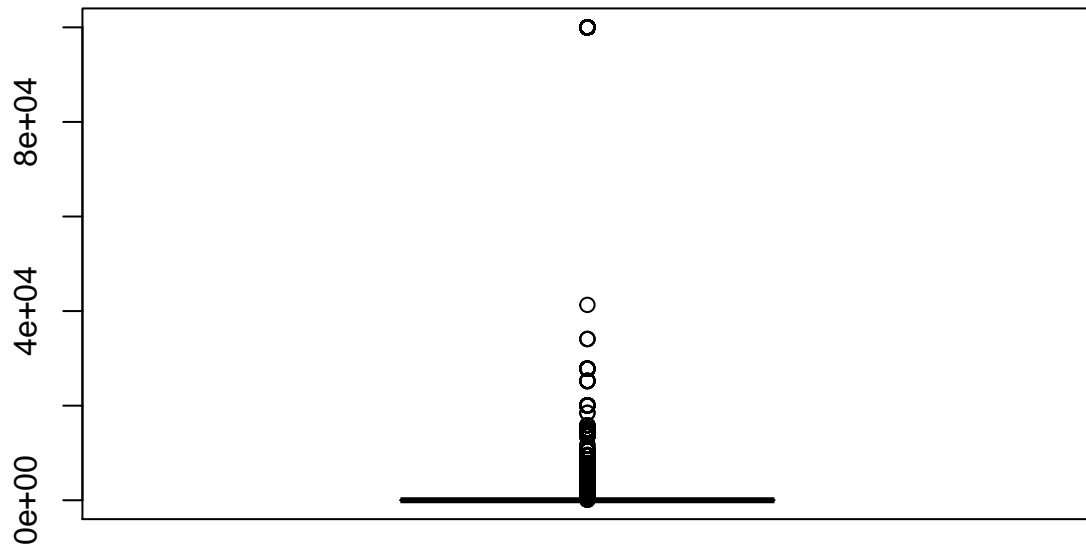


## Distribution of Capital Gain



```
boxplot(data$capital.gain,main="Capital Gain")
```

## Capital Gain



```
summary(data$capital.loss)
```

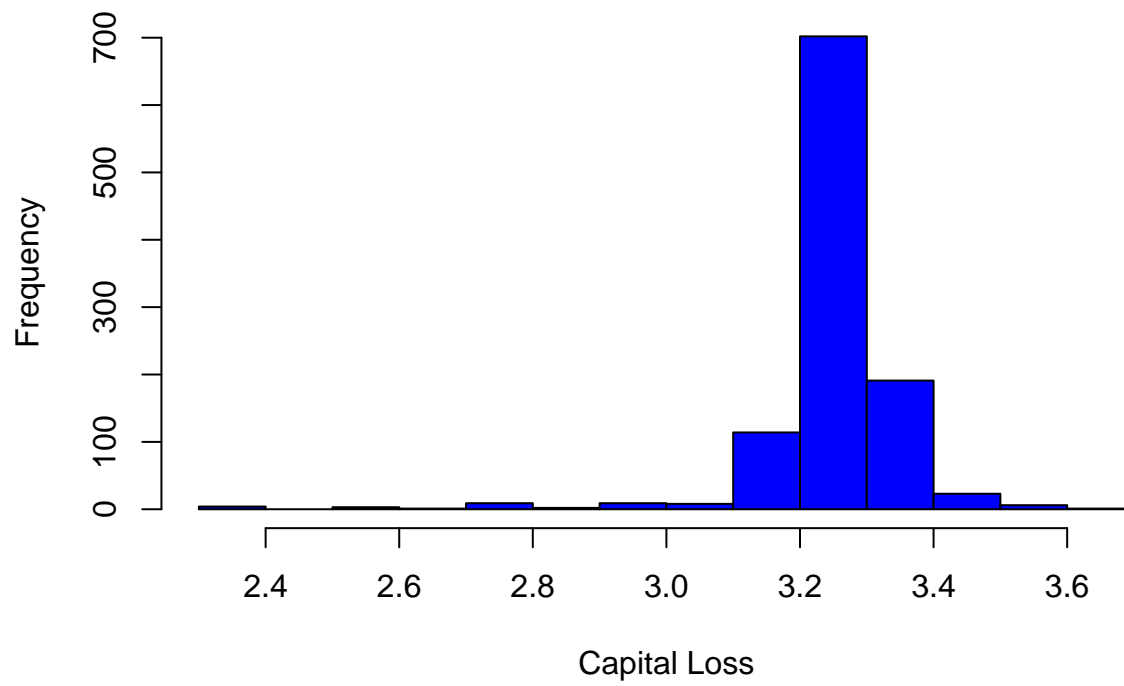
```
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##      0.00   0.00   0.00   88.53   0.00 4356.00
```

```
sd(data$capital.loss)
```

```
## [1] 404.9324
```

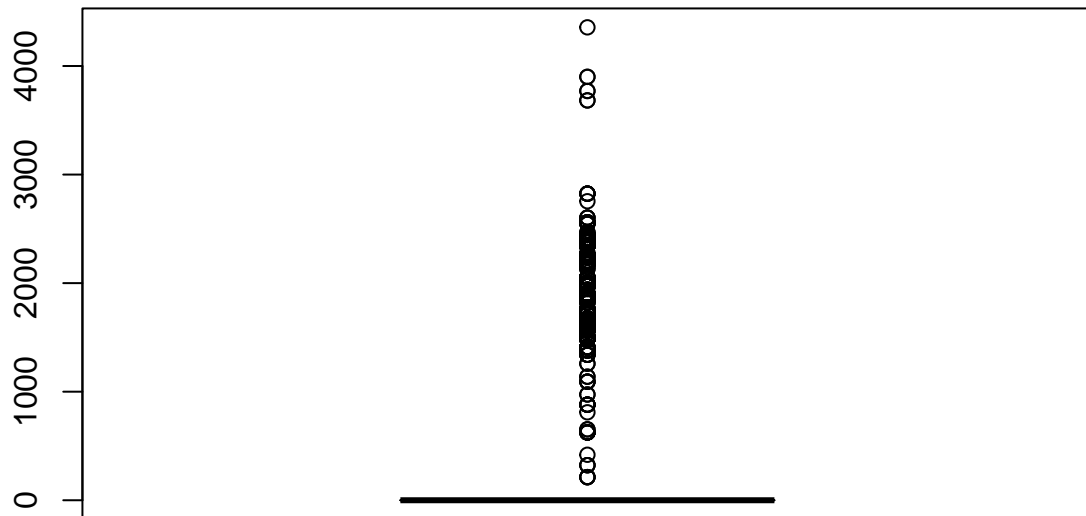
```
hist(log10(data$capital.loss),main = "Distribution of Capital Loss",xlab="Capital Loss",col = "blue")
```

## Distribution of Capital Loss



```
boxplot(data$capital.loss,main="Capital Loss")
```

## Capital Loss



```
summary(data$hours.per.week)
```

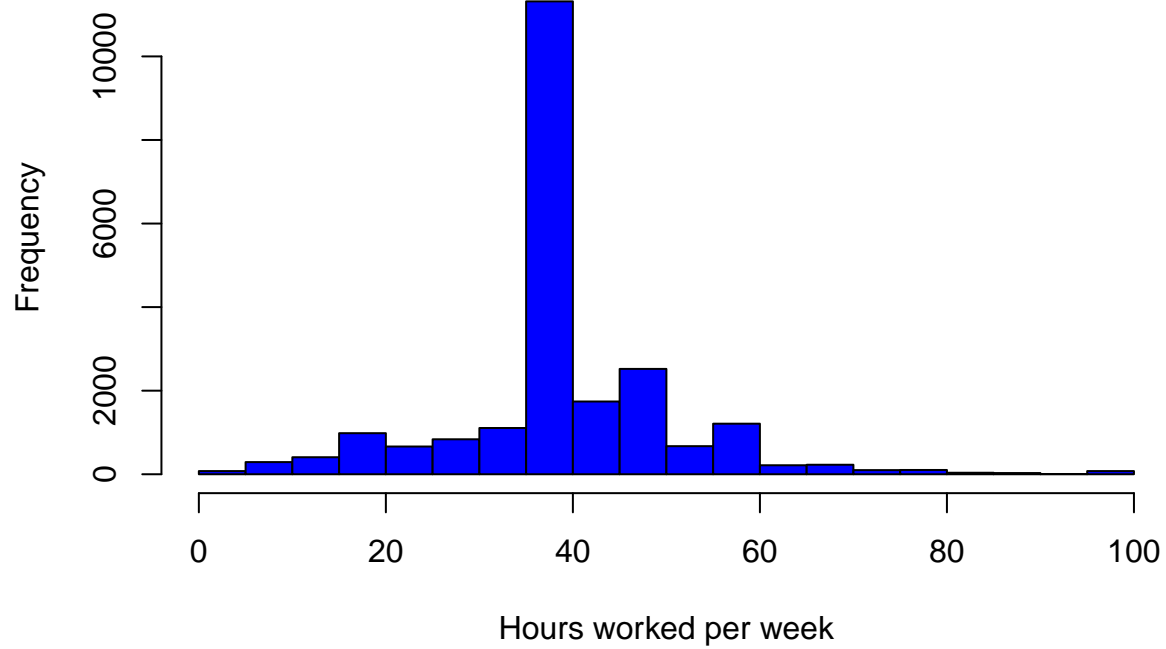
```
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##      1.00  40.00   40.00   41.02  45.00   99.00
```

```
sd(data$hours.per.week)
```

```
## [1] 11.98793
```

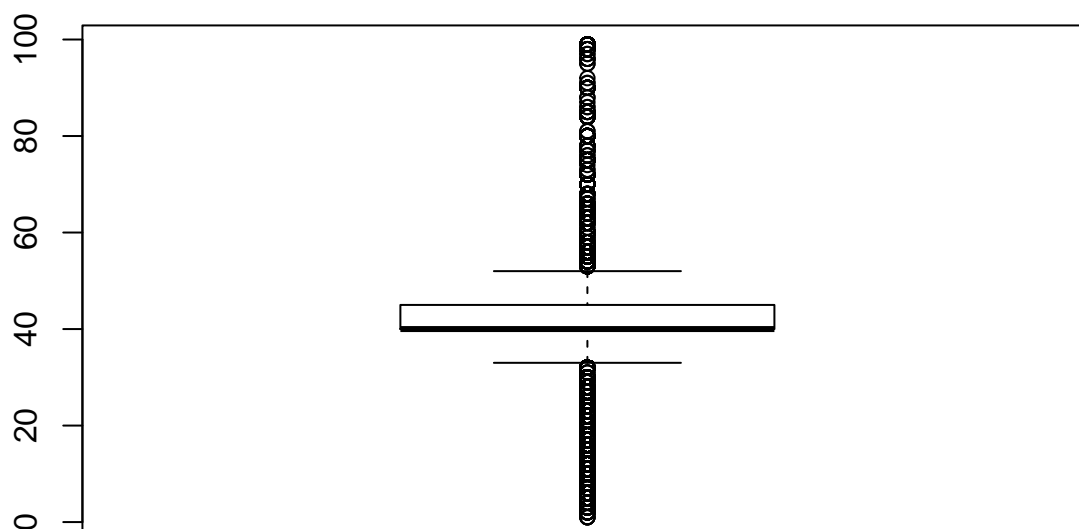
```
hist(data$hours.per.week,main = "Distribution of Hours Worked per Week",xlab="Hours worked per week",col="red",border="black",las=1)
```

## Distribution of Hours Worked per Week



```
boxplot(data$hours.per.week,main="Hours Worked per Week")
```

## Hours Worked per Week



8. Exploratory analysis of the attribute native country.

```
summary(data$native.country)
```

##	Cambodia	Canada
##	16	80
##	China	Columbia
##	54	37
##	Cuba	Dominican-Republic
##	73	49
##	Ecuador	El-Salvador
##	17	85
##	England	France
##	69	21
##	Germany	Greece
##	88	25
##	Guatemala	Haiti
##	42	29
##	Holand-Netherlands	Honduras
##	1	8
##	Hong	Hungary
##	14	11
##	India	Iran
##	81	30
##	Ireland	Italy
##	19	48
##	Jamaica	Japan

```
##           57           42
##           Laos           Mexico
##           15           453
##           Nicaragua Outlying-US(Guam-USVI-etc)
##           22           11
##           Peru           Philippines
##           24           138
##           Poland           Portugal
##           47           27
##           Puerto-Rico           Scotland
##           76           9
##           South           Taiwan
##           64           29
##           Thailand           Trinidad&Tobago
##           13           15
##           United-States           Vietnam
##           20617           46
##           Yugoslavia
##           14
```

9.1 Reducing/Combining levels of native country in training data.

```
data$native.country <- as.character(data$native.country)
asia <- c("Cambodia", "China", "Hong", "India", "Iran", "Japan", "Laos", "Philippines", "Taiwan", "Thailand")
northAmerica <- c("Canada", "Cuba", "Dominican-Republic", "El-Salvador", "Guatemala", "Haiti", "Honduras")
southAmerica <- c("Columbia", "Ecuador", "Peru")
europe <- c("England", "France", "Germany", "Greece", "Holand-Netherlands", "Hungary", "Ireland", "Italy", "Yugoslavia")
other <- c("South")
data$native.country[data$native.country %in% northAmerica] <- "North America"
data$native.country[data$native.country %in% asia] <- "Asia"
data$native.country[data$native.country %in% southAmerica] <- "South America"
data$native.country[data$native.country %in% europe] <- "Europe"
data$native.country[data$native.country %in% other] <- "Other"
table(data$native.country)
```

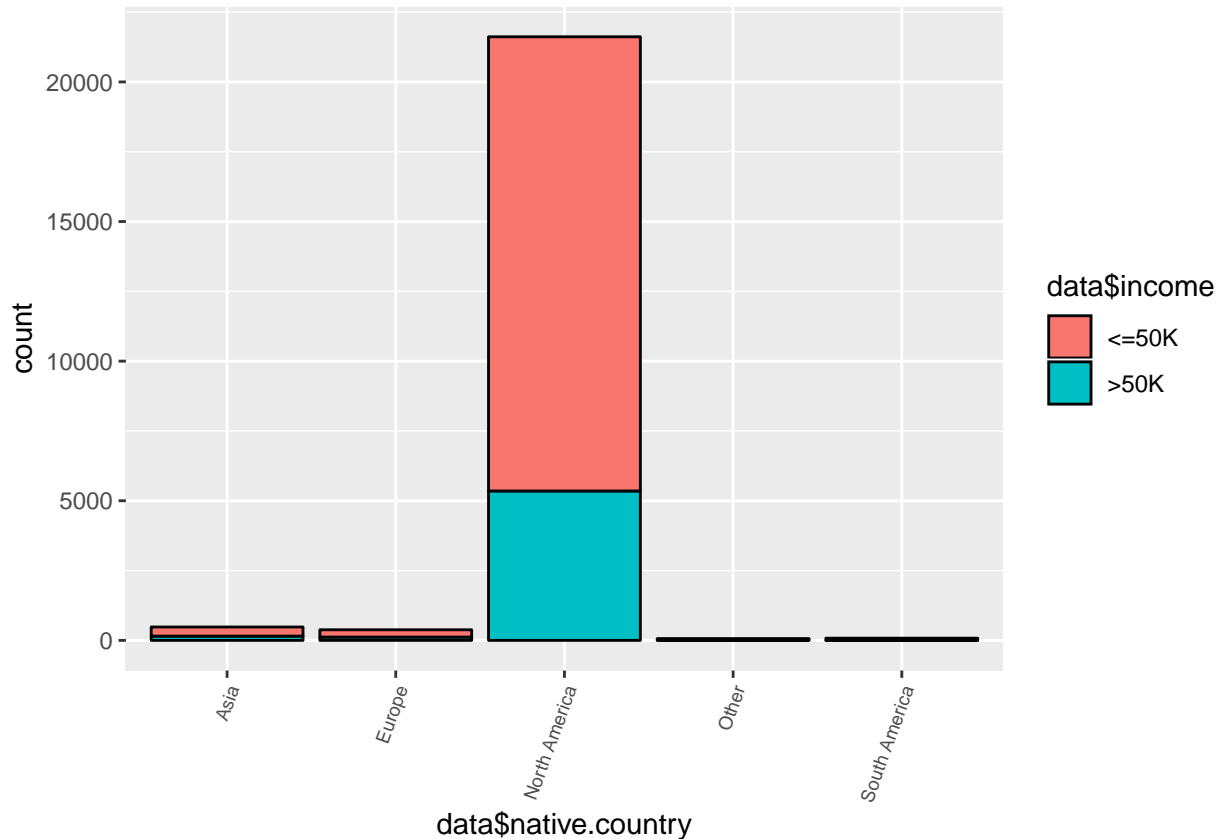
```
##
##           Asia           Europe North America           Other South America
##           478           379           21617           64           78
```

```
data$native.country <- as.factor(data$native.country)
levels(data$native.country)
```

```
## [1] "Asia"           "Europe"           "North America" "Other"
## [5] "South America"
```

```
## -----
## Reduce the level of native country into 5 levels "Asia", "Europe", "North America"
## "Other", "South America"
## -----
```

```
ggplot(data, aes(x=data$native.country, fill=data$income)) + geom_bar(position = "stack", color = "black")
```



```
## -----
## Native country of the Majority of the population is North America.
## -----
```

9.2. Reducing/Combining levels of native country in testing data.

```
testingdata$native.country <- as.character(testingdata$native.country)
testingdata$native.country[testingdata$native.country %in% northAmerica] <- "North America"
testingdata$native.country[testingdata$native.country %in% asia] <- "Asia"
testingdata$native.country[testingdata$native.country %in% southAmerica] <- "South America"
testingdata$native.country[testingdata$native.country %in% europe] <- "Europe"
testingdata$native.country[testingdata$native.country %in% other] <- "Other"
table(testingdata$native.country)
```

```
##
##      Asia      Europe North America      Other South America
##      156       114       7234          7         35
```

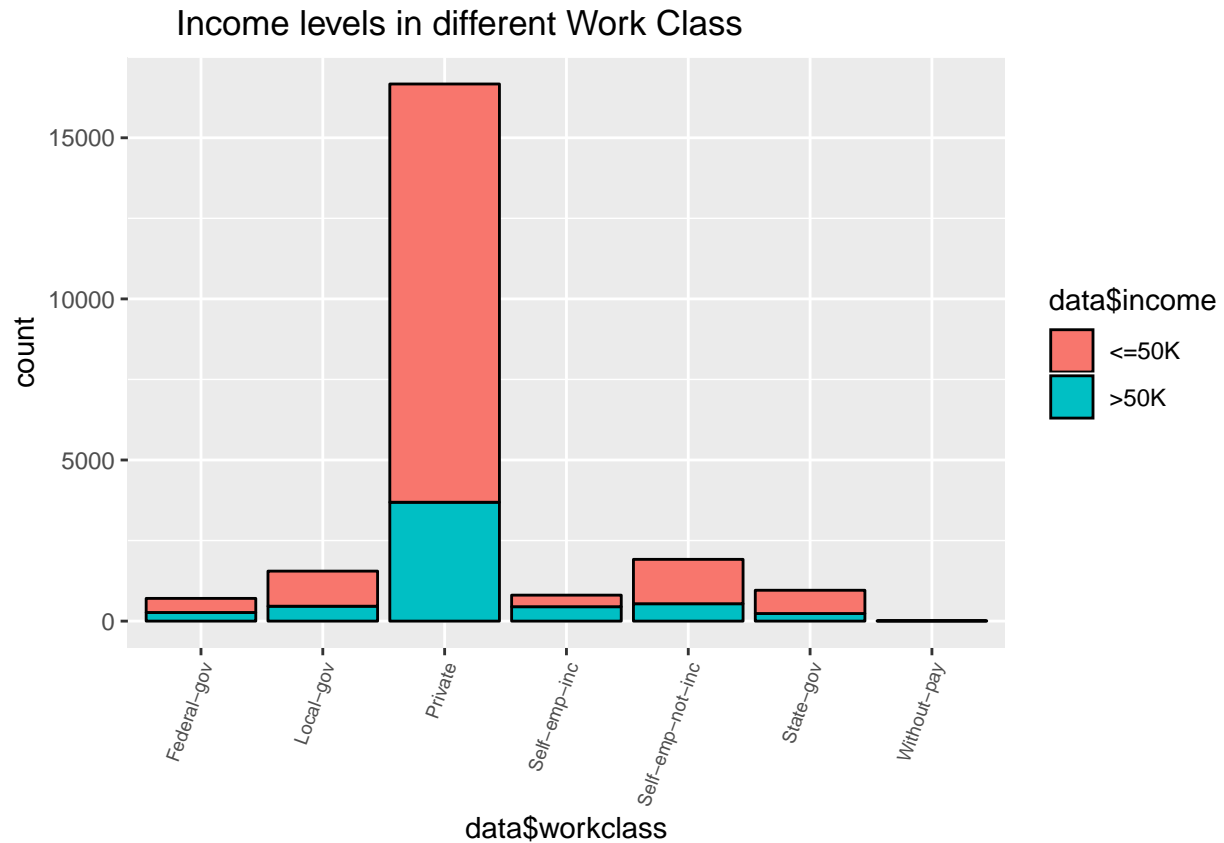
```
testingdata$native.country <- as.factor(testingdata$native.country)
levels(testingdata$native.country)
```

```
## [1] "Asia"      "Europe"    "North America" "Other"
## [5] "South America"
```

10.1. Combining categories of work class in training data.

```
ggplot(data, aes(x=data$workclass,fill=data$income)) + geom_bar(position = "stack", color = "black") +
```



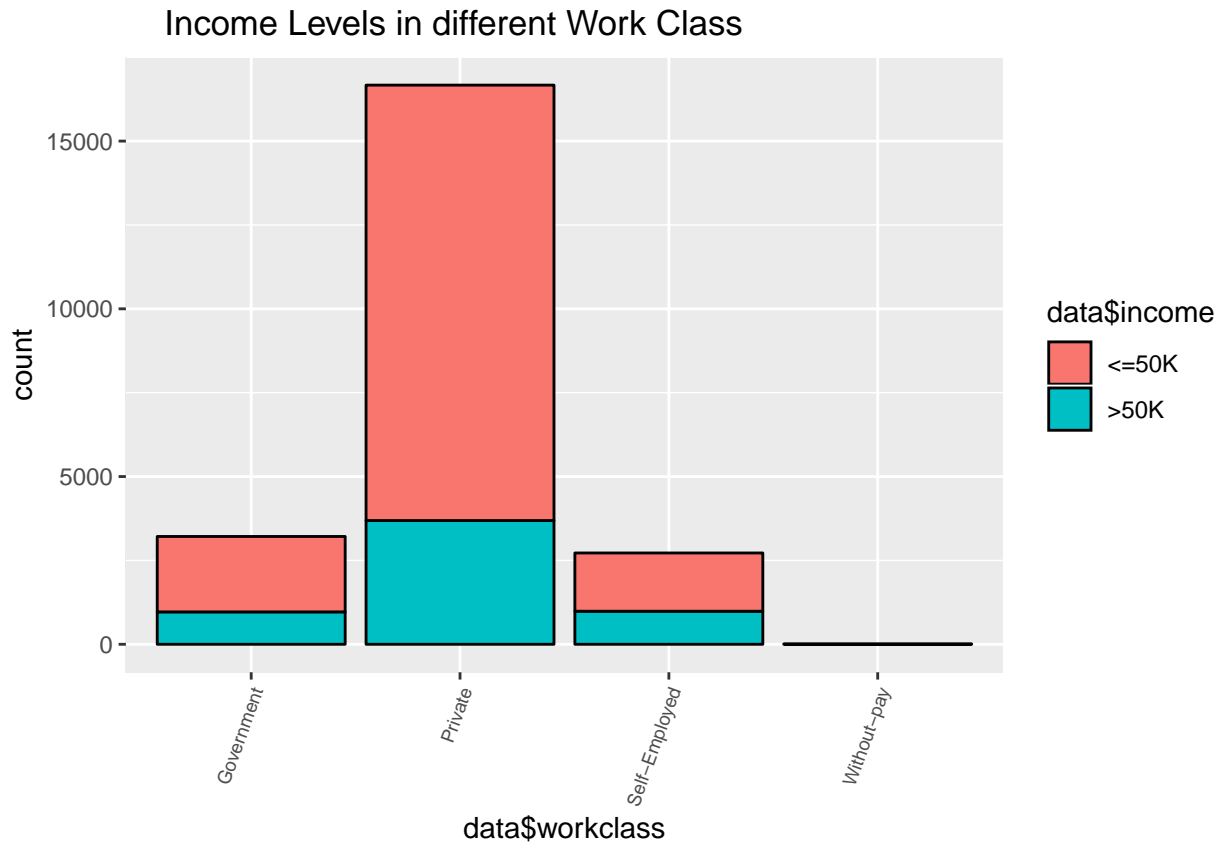


```
data$workclass <- gsub('^Federal-gov', 'Government', data$workclass)
data$workclass <- gsub('^Local-gov', 'Government', data$workclass)
data$workclass <- gsub('^State-gov', 'Government', data$workclass)

data$workclass <- gsub('^Self-emp-inc', 'Self-Employed', data$workclass)
data$workclass <- gsub('^Self-emp-not-inc', 'Self-Employed', data$workclass)

data$workclass <- gsub('^Other', 'Other', data$workclass)
data$workclass <- gsub('^Unknown', 'Other', data$workclass)

data$workclass <- as.factor(data$workclass)
ggplot(data, aes(x=data$workclass,fill=data$income)) + geom_bar(position = "stack", color = "black") +
```



```
## -----
## Replace Federal-gov, Local_gov and State_gov into government.
## Self-emp-inc and self-emp-not-inc into Self-Employed.
## other and unknown into other.
## -----

## Observations: Most of the people earning more than 50K are in private sector
## after that self employment and then Government.
```

10.2. Combining categories of work class in testing data.

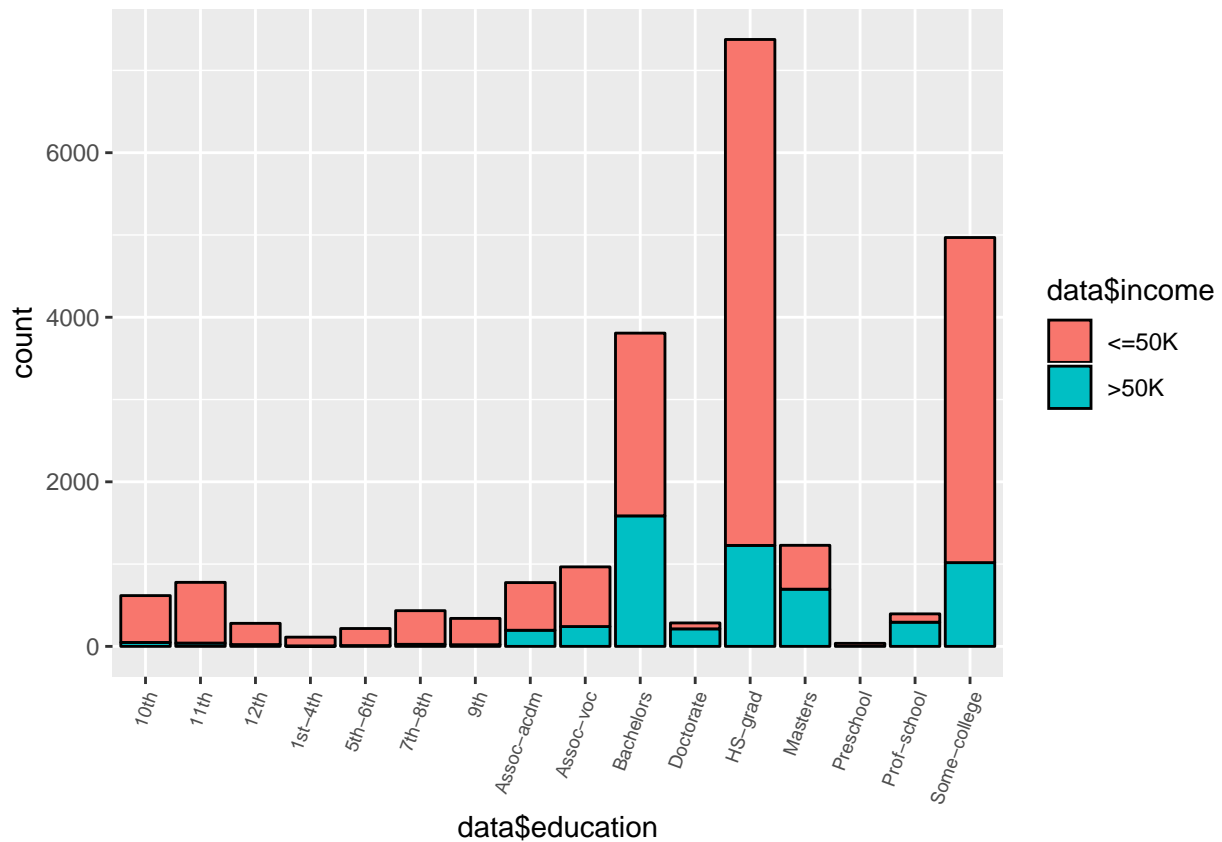
```
data$workclass <- gsub('^Federal-gov', 'Government', data$workclass)
data$workclass <- gsub('^Local-gov', 'Government', data$workclass)
data$workclass <- gsub('^State-gov', 'Government', data$workclass)

data$workclass <- gsub('^Self-emp-inc', 'Self-Employed', data$workclass)
data$workclass <- gsub('^Self-emp-not-inc', 'Self-Employed', data$workclass)

data$workclass <- gsub('^Other', 'Other', data$workclass)
data$workclass <- gsub('^Unknown', 'Other', data$workclass)
data$workclass <- as.factor(data$workclass)
```

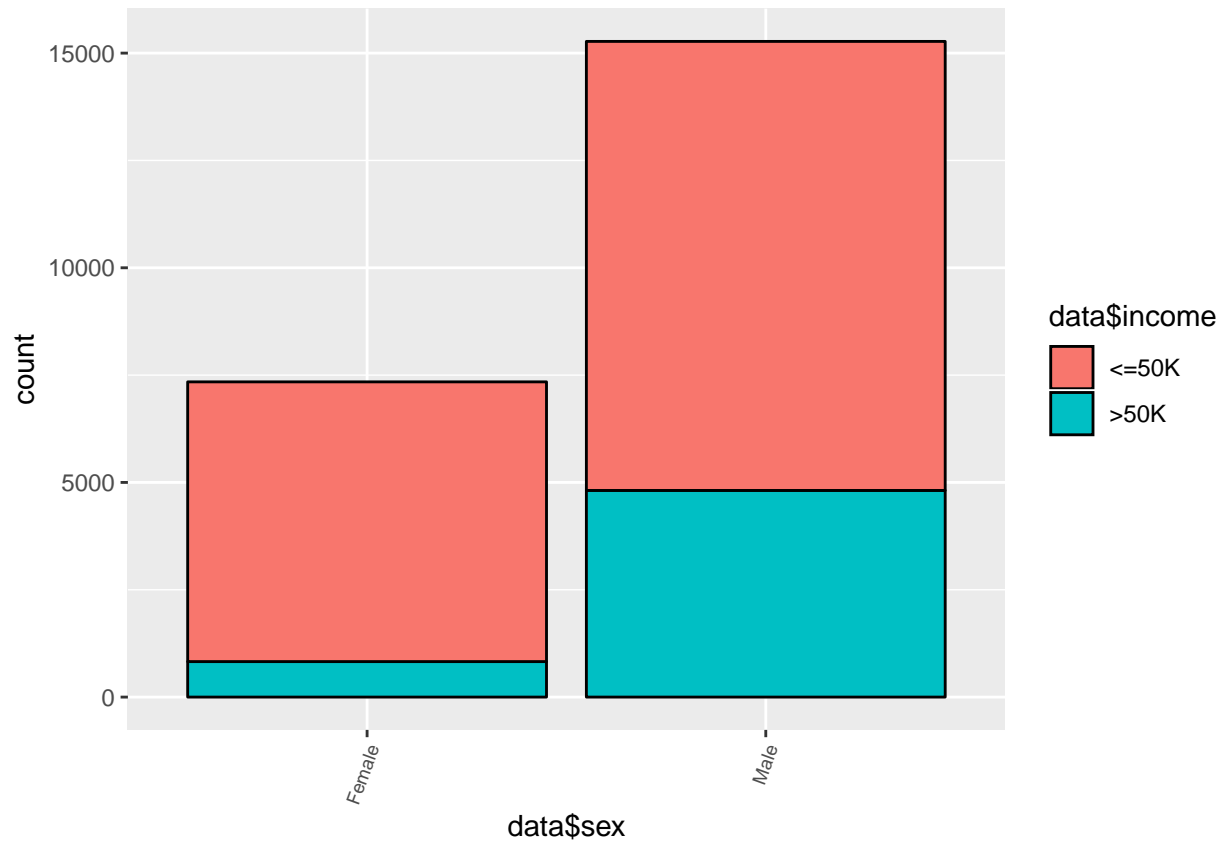
11. Relationship between categorical variables and income.

```
ggplot(data, aes(x=data$education,fill=data$income)) + geom_bar(position = "stack", color = "black") +
```



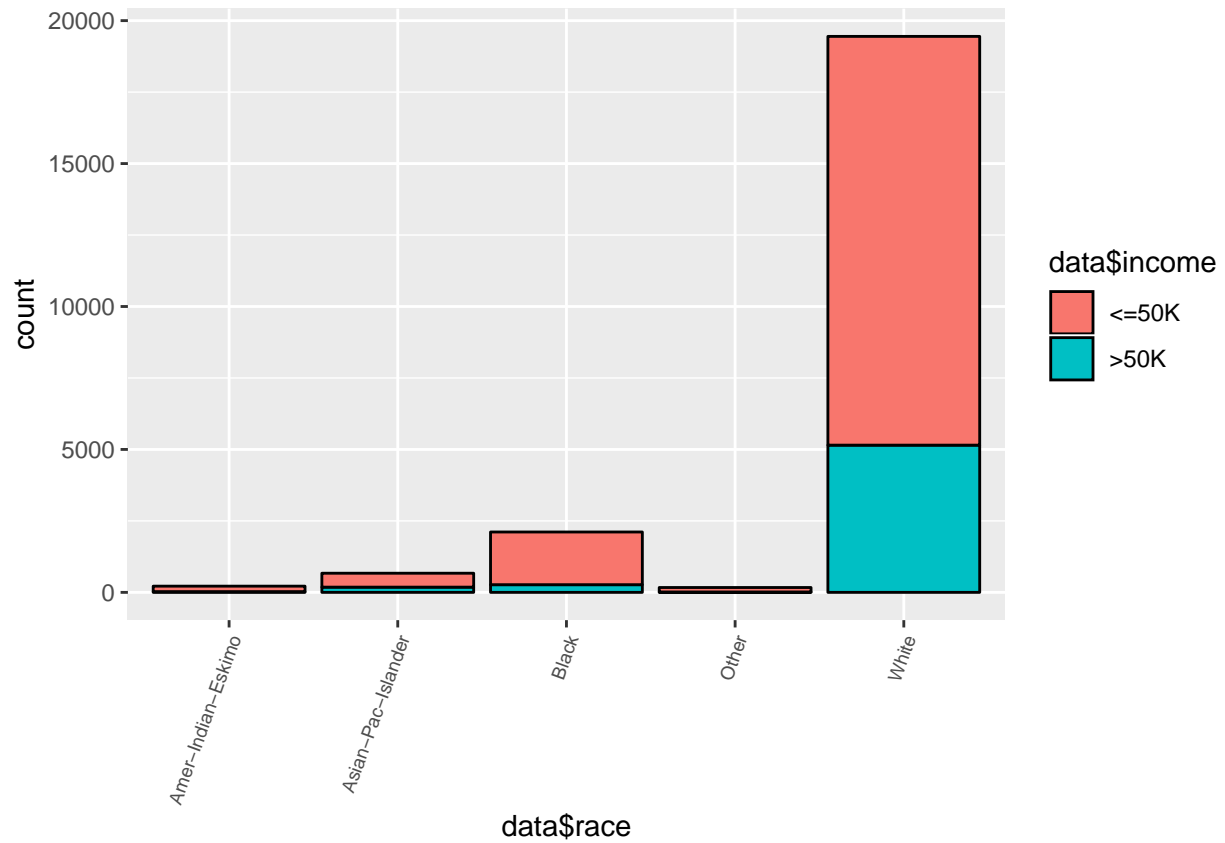
```
##-----
## The plot shows that the maximum number of adults earning income greater than 50K
## have bachelor's degree.
## In doctorate and masters also, the largest proportion is earning greater than 50 K.
## In lower education levels the largest proportion have income less than 50K.
## Higher education results in higher income.
##-----
```

```
ggplot(data, aes(x=data$sex,fill=data$income)) + geom_bar(position = "stack", color = "black") + theme(
```



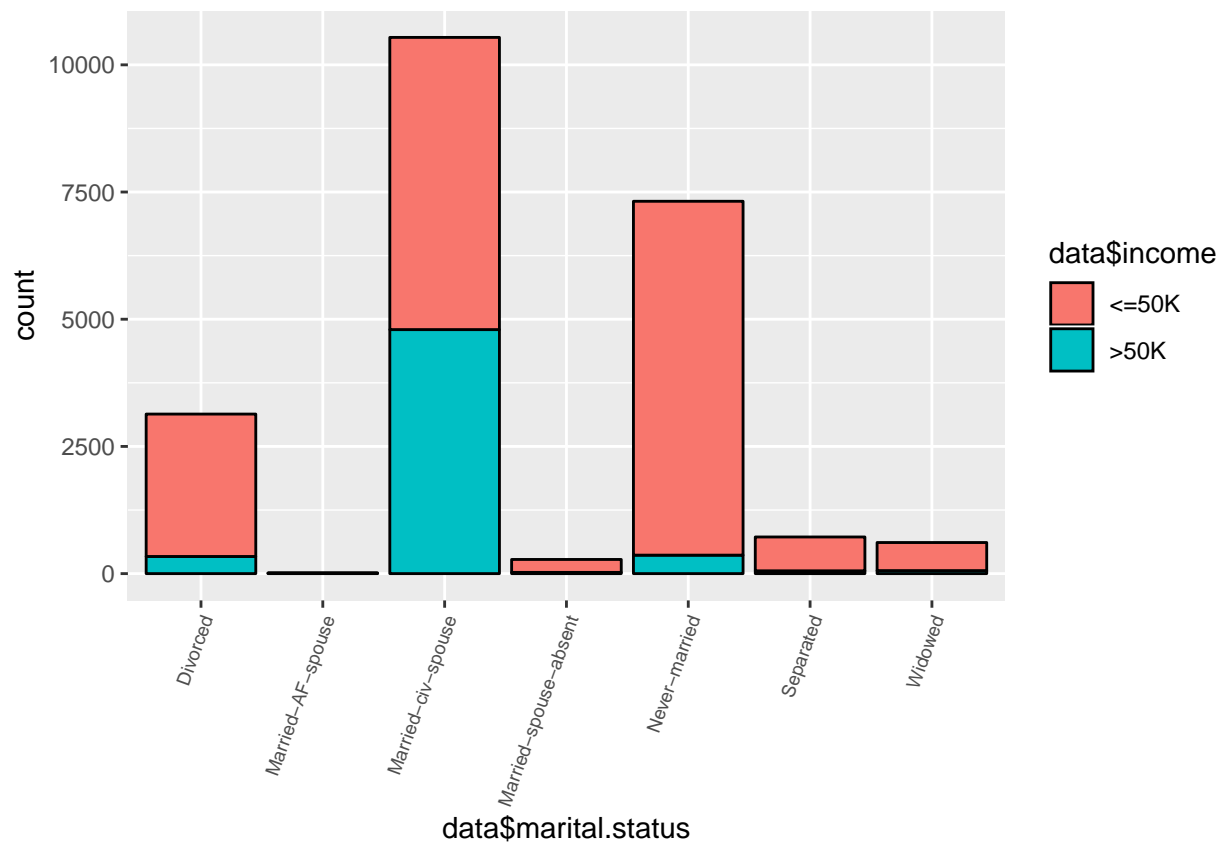
```
##-----
## Ratio of Males earning income greater than 50K are more as compare to female.
##-----
```

```
ggplot(data, aes(x=data$sex,fill=data$income)) + geom_bar(position = "stack", color = "black") + theme
```

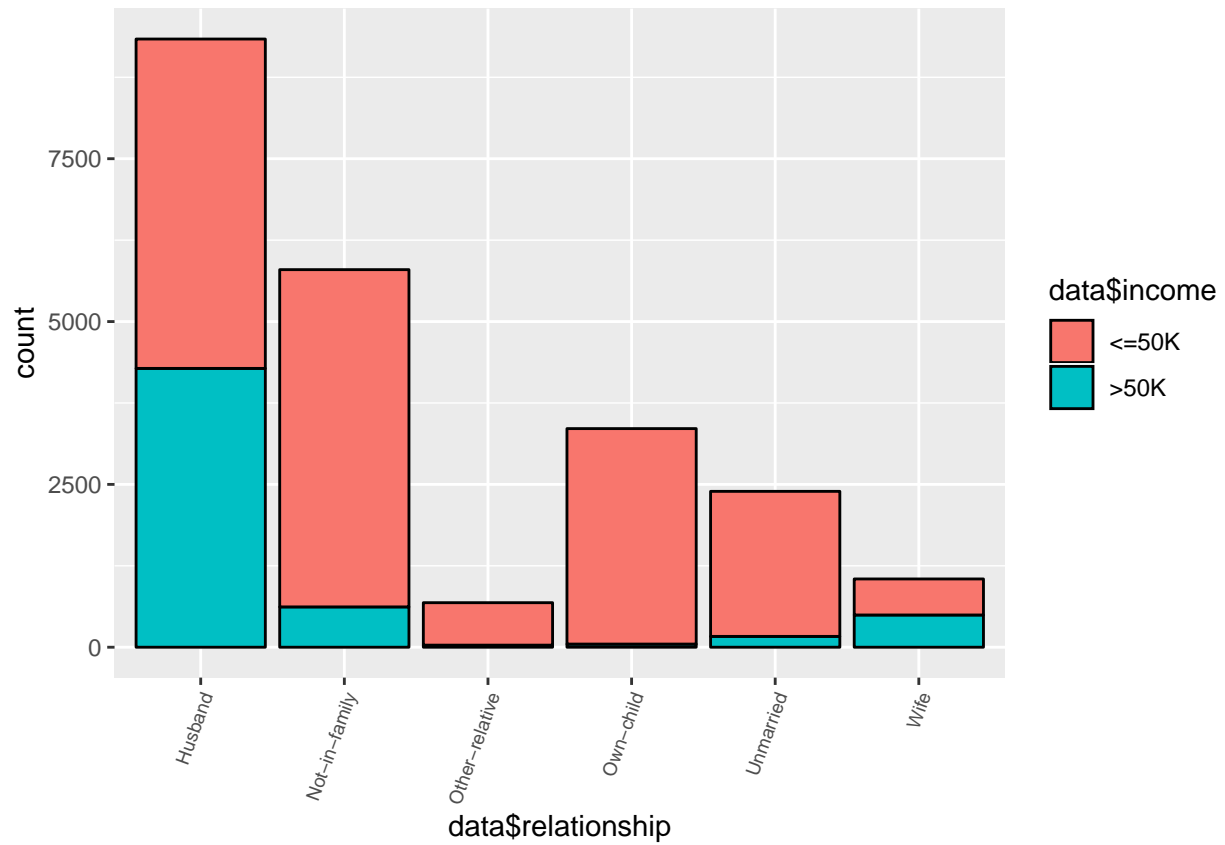


```
##-----
## Observations: Plot shows that in terms of race the highest earning people are
## from race White, then Asian-pacific and black.
##-----
```

```
ggplot(data, aes(x=data$marital.status, fill=data$income)) + geom_bar(position = "stack", color = "black")
```

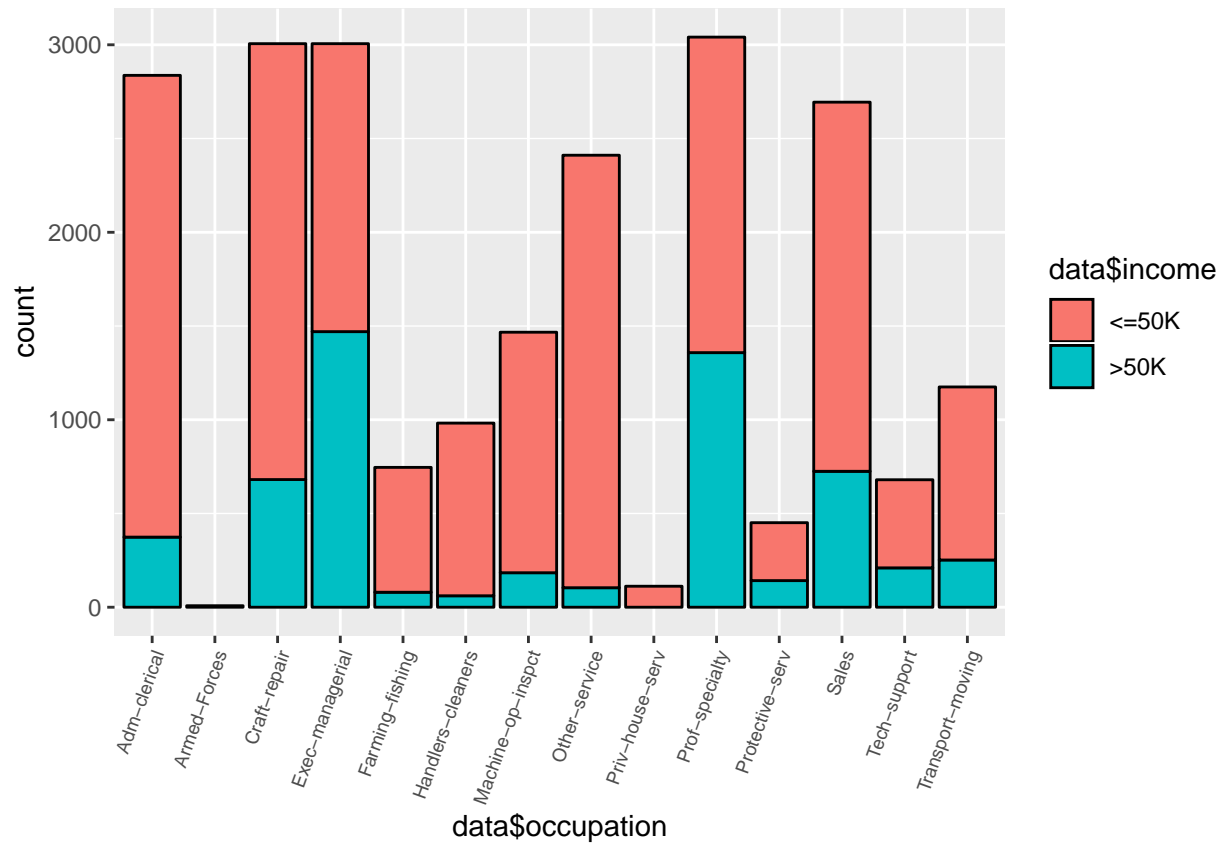


```
ggplot(data, aes(x=data$relationship, fill=data$income)) + geom_bar(position = "stack", color = "black")
```



```
##-----
## Observations: Plots shows that married people are earning more than 50k.
##-----
```

```
ggplot(data, aes(x=data$occupation,fill=data$income)) + geom_bar(position = "stack", color = "black") +
```



```
##-----
## Observations: Plots shows that in terms of occupation people with managerial job and
## professors are earning more than 50 K in the highest ratio.
## Showing that people at highest post are earning more.
##-----
```