

# GU4206-GR5206

*Name and UNI*

*3/02/2018*

The STAT GU4206/GR5206 Spring 2018 Midterm is open notes, open book(s), open computer and online resources are allowed. Students are **not** allowed to communicate with any other people regarding the exam with the exception of the instructor (Gabriel Young) and TA (Fan Gao). This includes emailing fellow students, using WeChat and other similar forms of communication. If there is any suspicion of one or more students cheating, further investigation will take place. If students do not follow the guidelines, they will receive a zero on the exam and potentially face more severe consequences. The exam will be posted on Canvas at 10:05AM. Students are required to submit both the .pdf and .Rmd files on Canvas (or .html if you must) by 12:40PM. Late exams will not be accepted. If for some reason you are unable to upload the completed exam on Canvas by 12:40PM, then immediately email markdown file to the course TA (fg2425).

## Part 1 (CDC Cancer Data - Subsetting and Plotting)

Consider the following dataset **BYSITE.TXT** taken directly from the Center of Disease Control's website. This dataset describes incidence and mortality crude rates of several types of cancer over time and also includes demographic variables such as **RACE** and **SEX**. The variables of interest in this exercise are: **YEAR**, **RACE**, **SITE**, **EVENT\_TYPE**, and **CRUDE\_RATE**.

Load in the **BYSITE.TXT** dataset. Also look at the levels of the variable **RACE**.

```
cancer <- read.table("BYSITE.TXT", sep = "|", header=T,
                    na.strings=c("~", "."))
dim(cancer)

## [1] 44982    13

levels(cancer$RACE)

## [1] "All Races"                "American Indian/Alaska Native"
## [3] "Asian/Pacific Islander"    "Black"
## [5] "Hispanic"                 "White"
```

### Problem 1.1

Create a new dataframe named **Prostate** that includes only the rows for prostate cancer. Check that the **Prostate** dataframe has 408 rows.

```
# levels(cancer$SITE)
# code goes here
Prostate <- cancer[cancer$SITE == 'Prostate', ]
cat('Check: ', dim(Prostate))

## Check:  408 13      nrow    ncol
```

### Problem 1.2

Using the **Prostate** dataframe from Problem 1.1, compute the average incidence crude rate for each level of **RACE**. To accomplish this task, use the appropriate function from the **apply** family. **Note:** first extract the rows that correspond to **EVENT\_TYPE** equals **Incidence**. Then use the appropriate function from the **apply** family with continuous variable **CRUDE\_RATE**.

```
tapply(ave_incidence$CRUDE_RATE, ave_incidence$RACE, mean)
```

```
# code goes here
ave_incidence <- Prostate[Prostate$EVENT_TYPE == 'Incidence', ]
by(ave_incidence$CRUDE_RATE, ave_incidence$RACE, mean)

## ave_incidence$RACE: All Races
## [1] 140.9294
## -----
## ave_incidence$RACE: American Indian/Alaska Native
## [1] 41.68824
## -----
## ave_incidence$RACE: Asian/Pacific Islander
## [1] 51.67647
## -----
## ave_incidence$RACE: Black
## [1] 152.4
## -----
## ave_incidence$RACE: Hispanic
## [1] 53.65882
## -----
## ave_incidence$RACE: White
## [1] 141.7588
```

### Problem 1.3

Refine the **Prostate** dataframe by removing rows corresponding to **YEAR** level **2010-2014** and removing rows corresponding to **RACE** level **All Races**. After removing the rows, convert **YEAR** into a numeric variable. Check that the new **Prostate** dataframe has 320 rows.

```
#levels(cancer$YEAR)
# code goes here
Prostate <- cancer[cancer$SITE == 'Prostate', ]
Prostate <- Prostate[Prostate$YEAR != '2010-2014' & Prostate$RACE != 'All Races', ]
# ?????????????????????????????
Prostate$YEAR <- sapply(Prostate$YEAR, as.numeric) + 1998
head(Prostate, 5) as.numeric(as.character(P$YEAR))
```

```
##      YEAR      RACE SEX  SITE EVENT_TYPE
## 11731 1999 American Indian/Alaska Native Male Prostate Incidence
## 11732 1999 American Indian/Alaska Native Male Prostate Mortality
## 11733 2000 American Indian/Alaska Native Male Prostate Incidence
## 11734 2000 American Indian/Alaska Native Male Prostate Mortality
## 11735 2001 American Indian/Alaska Native Male Prostate Incidence
##      AGE_ADJUSTED_CI_LOWER AGE_ADJUSTED_CI_UPPER AGE_ADJUSTED_RATE COUNT
## 11731          96.9          116.0          106.2    590
## 11732          13.2          22.5          17.5     63
## 11733          83.4          99.9          91.4    577
## 11734          16.1          25.1          20.3     92
## 11735          88.2          104.4          96.1    671
##      POPULATION CRUDE_CI_LOWER CRUDE_CI_UPPER CRUDE_RATE
## 11731    1365653          39.8          46.8          43.2
## 11732    1410781           3.4           5.7           4.5
## 11733    1444839          36.7          43.3          39.9
## 11734    1493056           5.0           7.6           6.2
## 11735    1538799          40.4          47.0          43.6
```

```
dim(Prostate)
```

```
## [1] 320 13
```

## Problem 1.4

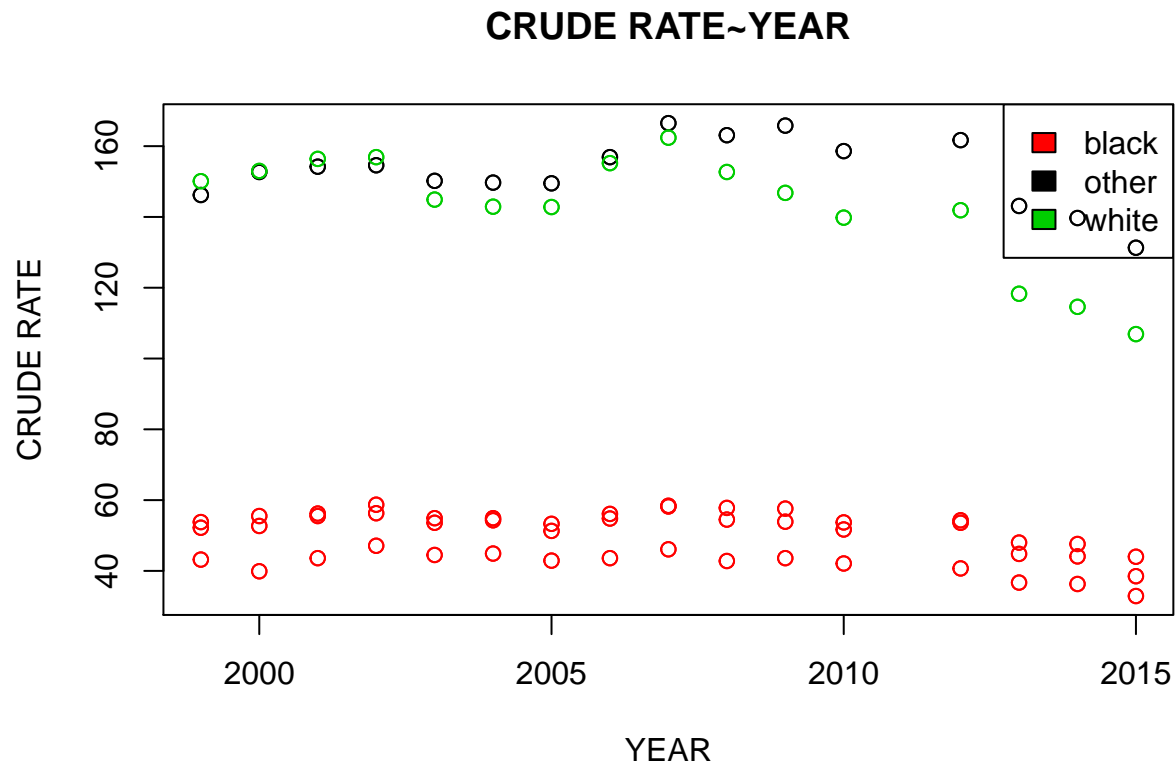
Create a new variable in the refined **Prostate** dataframe named **RaceNew** that defines three race levels: (1) white, (2) black, and (3) other. Construct a base-R plot that shows the incidence crude rate (not mortality) as a function of time (**YEAR**). Split the scatterplot by **RaceNew**. Make sure to include a legend and label the graphic appropriately.

```
# code goes here
```

```
# ?????? in not in
```

```
%in%
```

```
Prostate$RaceNew <- ifelse(Prostate$RACE == 'White', 'white', ifelse(Prostate$RACE == 'Black', 'black',
Prostate.Incidence <- Prostate[Prostate$EVENT_TYPE == 'Incidence',]
plot(x = Prostate.Incidence$YEAR, y = Prostate.Incidence$CRUDE_RATE, col = factor(Prostate.Incidence$RaceNew),
legend('topright', legend = levels(factor(Prostate.Incidence$RaceNew)), fill = unique(factor(Prostate.Incidence$RaceNew)))
```



## Part 2 (Basic Web Scraping)

### Problem 2.1

Open up the **SP500.html** file to get an idea of what the data table looks like. This website shows the SP500 monthly average closing price for every year from 1871 to 2018. Use regular expressions and the appropriate character-data functions to scrape a “nice” dataset out of the html code. Your final dataframe should have two variables: (1) the variable **Time**, which ranges from 1871 to 2018; (2) the variable **Price** which are the

corresponding SP500 price values for each year. Name the final dataframe **SP500.df** and display both the head and the tail of this scrapped dataset.

```
SP500 <- readLines("SP500.html")
# head(SP500)

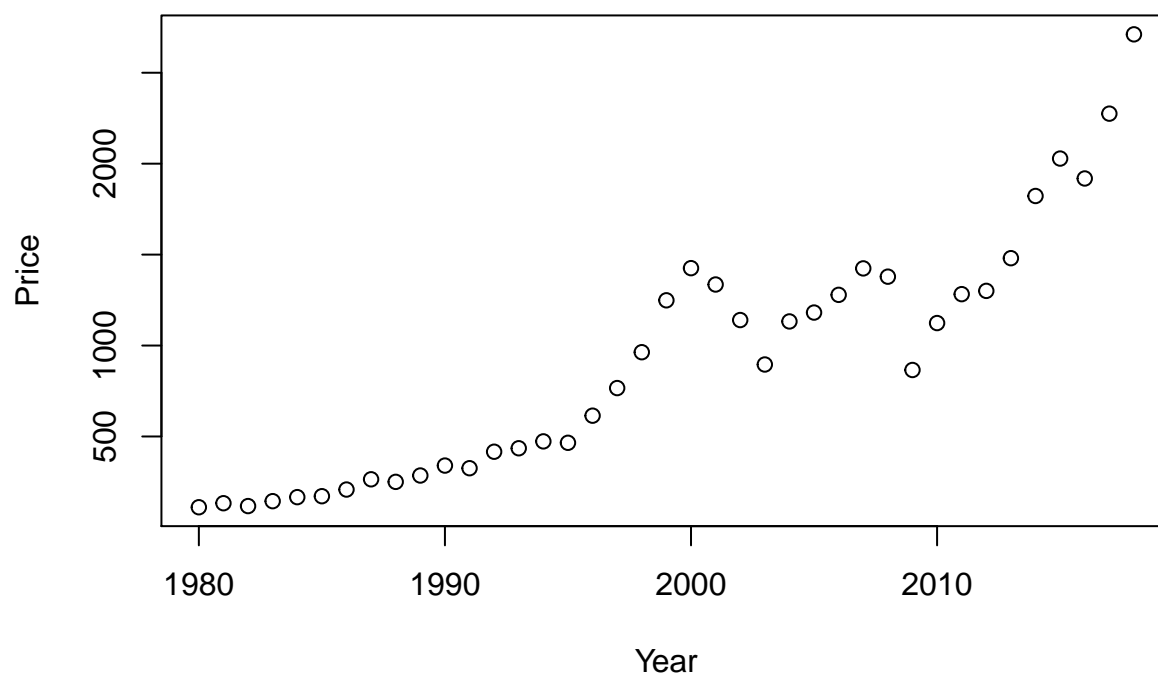
# code goes here
reg_data <- function (pattern, data) {
  sgrep <- grep(pattern, data)
  matches <- gregexpr(pattern = pattern, text = data[sgrep])
  reg.data <- unlist(regmatches(data[sgrep], matches), use.names = FALSE)
  return(reg.data)
}
time <- reg_data(pattern = '[A-Z][a-z]{2} [0-9]{1,2}[,] [0-9]{4}', SP500)
price <- reg_data(pattern = '<td class="right">[0-9]{1,3}[.][0-9]{2}', SP500)
price <- sapply(strsplit(price, '<td class="right">'), '[', 2)
SP500.df <- data.frame(time, price)
colnames(SP500.df) <- c('Time', 'Price')
```

## TA Problem 2.2

Create a time series plot of the monthly average SP500 closing price values over the years 1980 to 2018, i.e., use the first 40 lines of **SP500.df**.

```
# code goes here
year_pattern <- '[0-9]{4}'
SP500.year <- factor(reg_data(year_pattern, SP500.df$Time))
SP500.df$Year <- as.numeric(SP500.year) + 1870
SP500.df$Price <- as.character(SP500.df$Price)
# fix price
price_num <- c()
price_str <- strsplit(SP500.df$Price, ',')
for (i in 1:dim(SP500.df)[1]){
  if(is.na(price_str[[i]][2])){
    price_num[i] <- as.numeric(price_str[[i]][1])
  }else{
    price_num[i] <- as.numeric(paste(price_str[[i]][1], price_str[[i]][2], sep = ''))
  }
}
SP500.df$Price <- price_num
# groupby month
SP500.newdf <- aggregate(Price ~ Year, SP500.df[SP500.df$Year >= 1980, ], mean)
# time series plot
plot(SP500.newdf$Year, SP500.newdf$Price, xlab = 'Year', ylab = 'Price', main = 'S&P500 Time Series plot')
```

## S&P500 Time Series plot



## Part 3 (Knn Regression)

Recall the `kNN.decision` function from class. In the `kNN.decision` function, we classified the market direction using a non-parametric classification method known as “k-nearest neighbors.”

```
library(ISLR)
head(Smarket, 3)
```

```
##   Year Lag1  Lag2  Lag3  Lag4  Lag5 Volume Today Direction
## 1 2001 0.381 -0.192 -2.624 -1.055  5.010 1.1913 0.959      Up
## 2 2001 0.959  0.381 -0.192 -2.624 -1.055 1.2965 1.032      Up
## 3 2001 1.032  0.959  0.381 -0.192 -2.624 1.4112 -0.623     Down
```

```
KNN.decision <- function(Lag1.new, Lag2.new, K = 5,
                          Lag1 = Smarket$Lag1,
                          Lag2 = Smarket$Lag2,
                          Dir = Smarket$Direction) {
  n <- length(Lag1)
  stopifnot(length(Lag2) == n, length(Lag1.new) == 1,
            length(Lag2.new) == 1, K <= n)

  dists <- sqrt((Lag1-Lag1.new)^2 + (Lag2-Lag2.new)^2)

  neighbors <- order(dists)[1:K]
  neighb.dir <- Dir[neighbors]
  choice <- names(which.max(table(neighb.dir)))
}
```

```

    return(choice)
}
KNN.decision(Lag1.new=2,Lag2.new=4.25)

```

```
## [1] "Down"
```

### Problem 3.1

In our setting, we consider two datasets that describe yearly US congressional approval ratings over the years 1974 to 2012. The first file **Congress\_train.csv** is the training (or model building) dataset and the second file “**Congress\_test.csv**” is the test (or validation) dataset. The code below reads in the data and plots each set on separate graphs.

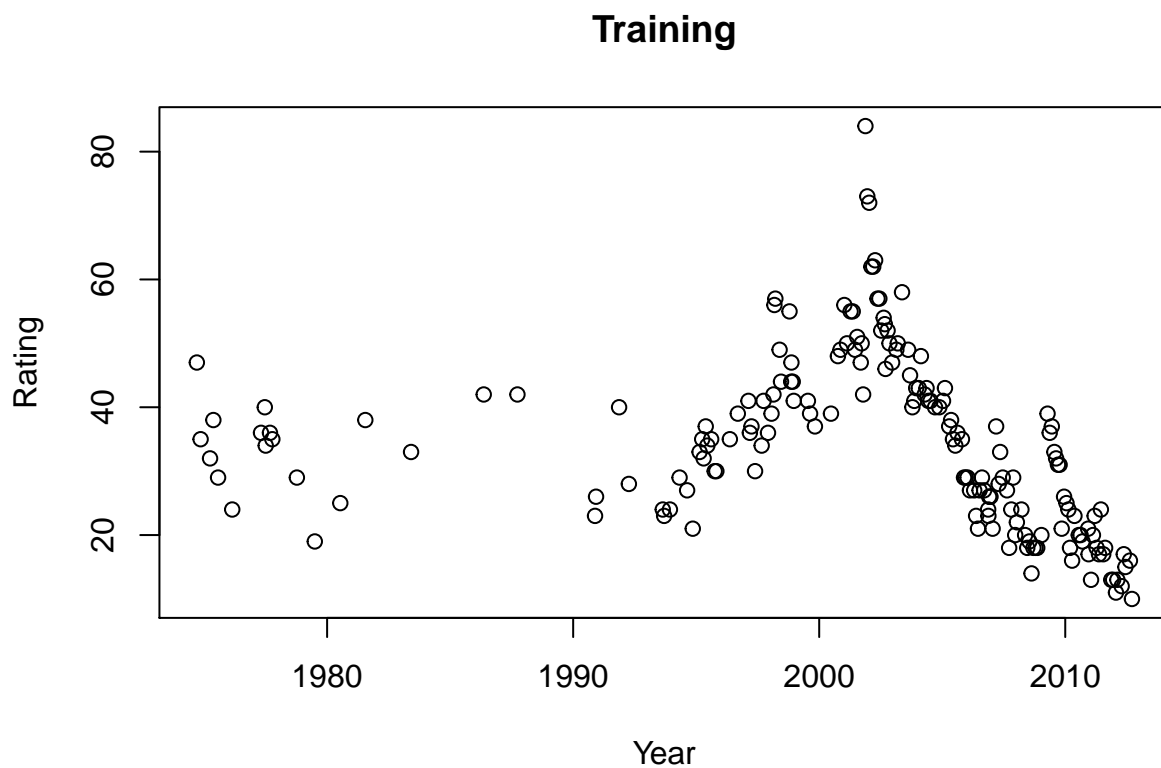
```

Congress_train <- read.csv("Congress_train.csv")
n_train <- nrow(Congress_train)
n_train

```

```
## [1] 181
```

```
plot(Congress_train$Year,Congress_train$Rating,xlab="Year",ylab="Rating",main="Training")
```



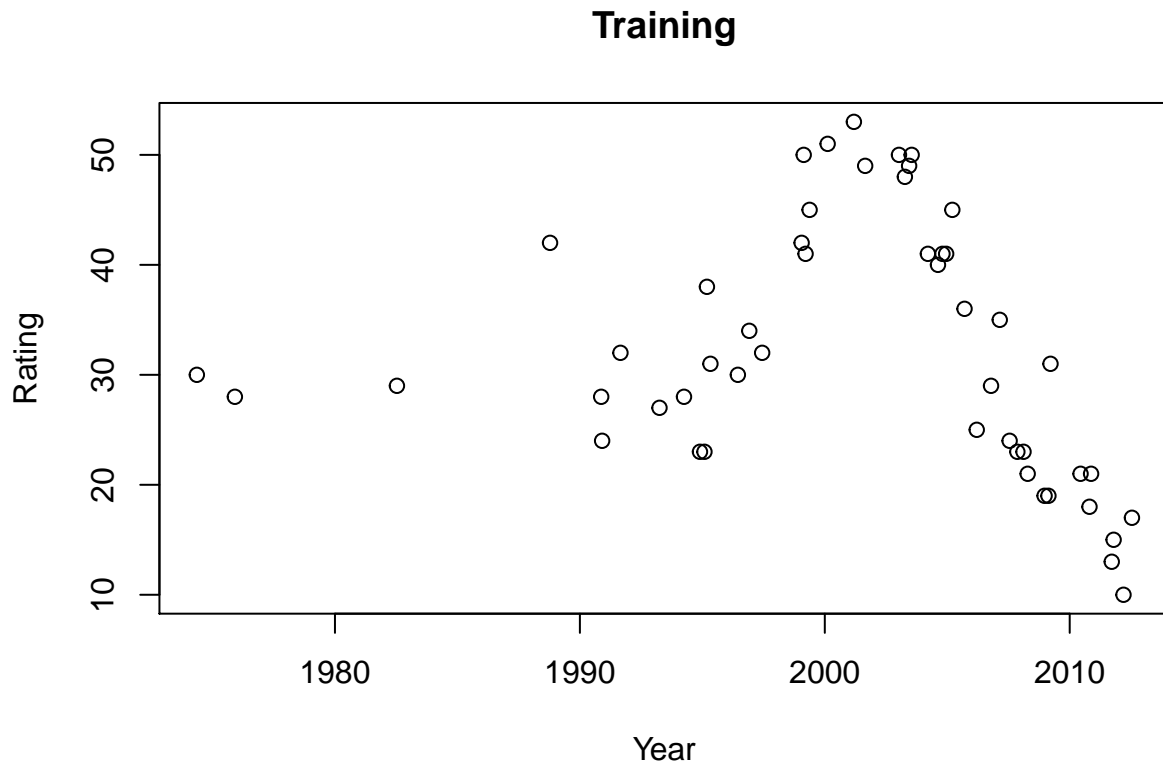
```

Congress_test <- read.csv("Congress_test.csv")
n_test <- nrow(Congress_test)
n_train

```

```
## [1] 181
```

```
plot(Congress_test$Year, Congress_test$Rating, xlab="Year", ylab="Rating", main="Training")
```



Write a function called **kNN.regression** which fits a non-parametric curve to a continuous response. Here you will fit a “moving average” to the yearly congressional approval ratings over the years 1974 to 2012. There is only one feature in this exercise, i.e., **Year** is the only independent variable. Thus for a test time say  $t = t_0$ , we compute the **arithmetic average rating** of the  $K$  closest neighbors of  $t_0$ . Using the **Congress\_train** dataset, train your model to predict the approval rating when  $t = 2000$ . Set the tuning parameter to  $K = 5$ .

**Note:** to receive full credit, you must extend off of the **kNN.decision** function. You cannot just look up a moving average function online. The new function should also include euclidean distance and the **order** function.

```
# code goes here
KNN.decision.reg <- function(newX, K = 5,
                             data = Congress_train,
                             nameX = 'Year',
                             nameY = 'Rating') {
  # calculate distance
  dists <- sqrt((data[nameX] - newX)^2)
  # find neighbors
  neighbors <- order(dists)[1:K]
  neighb.dir <- data[neighbors, 'Rating']
  choice <- names(which.max(table(neighb.dir)))
  return(choice)
}
KNN.decision.reg(newX = 2000)
```

```
## [1] "39"
```

## Problem 3.2

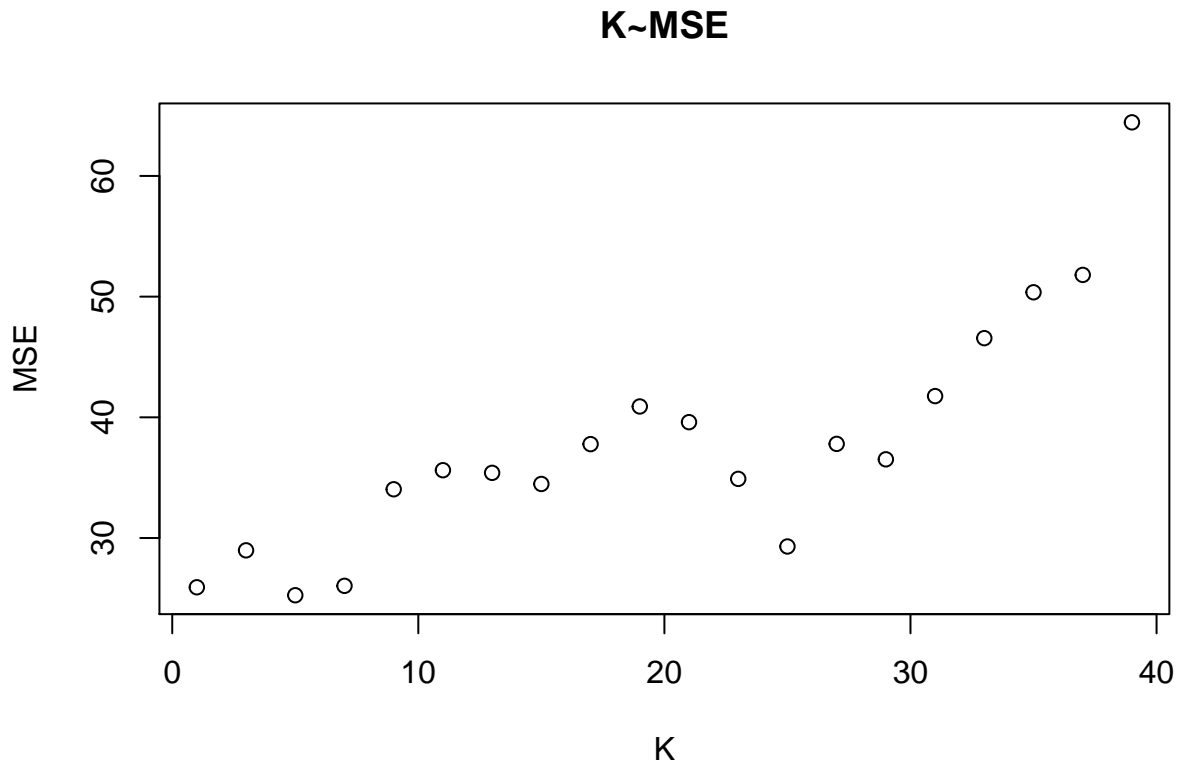
Compute the **test mean square error** using neighborhood sizes  $K = 1, 3, 5, \dots, 39$ . In this exercise you will train the model using **Congress\_train** and assess its performance on the test data **Congress\_test** with the different tuning parameters  $K$ . Plot the test mean square error as a function of  $K$  and choose the best value of the tuning parameter based on this output.

```
# code goes here
# parameters include global variables
MSE <- function(k) {
  y_hat <- c()
  for (i in 1:length(Congress_test$Year)) {
    y_hat[i] <- KNN.decision.reg(newX = Congress_test$Year[i], K = k)
  }
  y_hat <- as.numeric(y_hat)
  y <- Congress_test$Rating
  MSE <- sum((y - y_hat)**2)/length(y)
  return(MSE)
}
# initialization
MSEvec <- c()
count <- 1
for (i in seq(1, 39, by = 2)){
  MSEvec[count] <- MSE(i)
  count <- count + 1
}
# create df for result
MSE.df <- data.frame(seq(1, 39, by = 2), MSEvec)
colnames(MSE.df) <- c('K', 'MSE')
# best
MSE.df[order(MSE.df$MSE)[1],]

##      K      MSE
## 3 5 25.26

# plot
plot(MSE.df$K, MSE.df$MSE, xlab = 'K', ylab = 'MSE', main = 'K~MSE')
```





### Problem 3.2

Compute the **training mean square error** using neighborhood sizes  $K = 1, 3, 5, \dots, 39$ . In this exercise you will train the model using **Congress\_train** and assess its performance on the training data **Congress\_train** with the different tuning parameters  $K$ . Plot both the test mean square error and the training mean square error on the same graph. Comment on any interesting features/patterns displayed on the graph.

```
# code goes here
MSE.train <- function (k) {
  y_hat <- c()
  for (i in 1:length(Congress_train$Year)) {
    y_hat[i] <- KNN.decision.reg(newX = Congress_train$Year[i], K = k)
  }
  y_hat <- as.numeric(y_hat)
  y <- Congress_train$Rating
  MSE <- sum((y - y_hat)**2)/length(y)
  return(MSE)
}

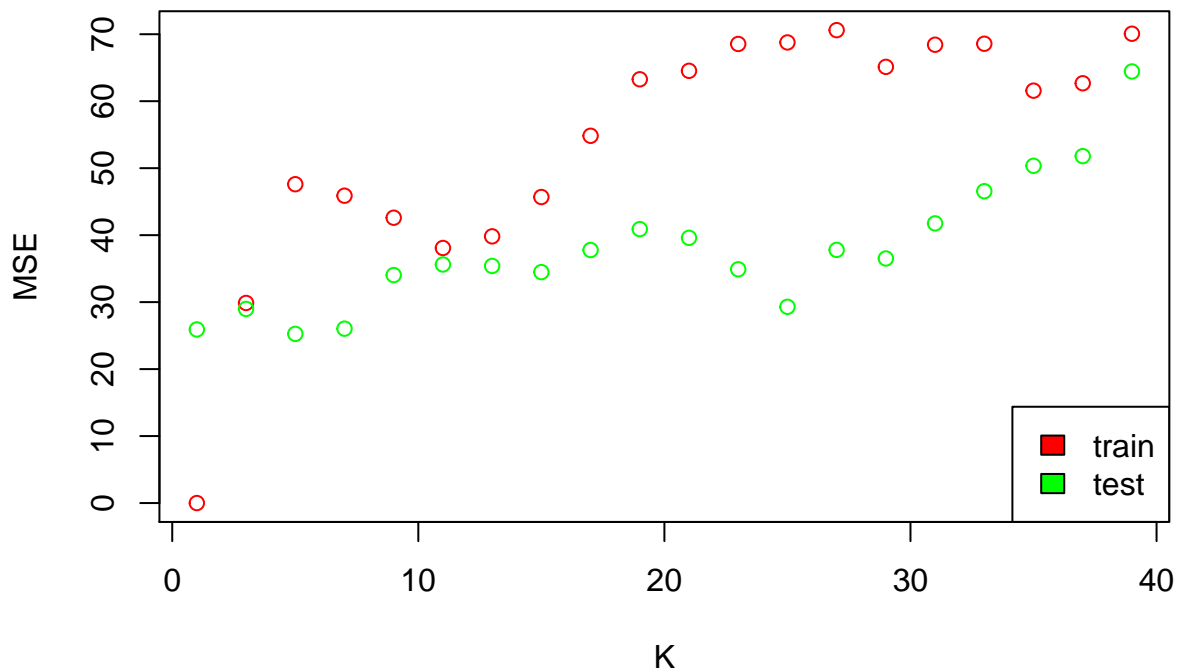
# initialization
MSEvec.train <- c()
count <- 1
for (i in seq(1, 39, by = 2)){
  MSEvec.train[count] <- MSE.train(i)
  count <- count + 1
}

# create df for result
```

```

MSE_train.df <- data.frame(seq(1, 39, by = 2), MSEvec.train)
colnames(MSE_train.df) <- c('K', 'MSE_train')
test.train <- cbind(MSE_train.df, MSE.df)
# plot
plot(test.train$K, test.train$MSE_train, col = 'red', xlab = 'K', ylab = 'MSE')
points(test.train$K, test.train$MSE, col = 'green')
legend('bottomright', legend = c('train', 'test'), fill = c('red', 'green'))

```



### Problem 3.3 (Extra Credit)

#### Code

Plot the kNN-regression over the training data set **Congress\_train** using optimal tuning parameter  $K$ . In this plot, the years must be refined so that the smoother shows predictions for all years from 1973 to 2015.

```
# code goes here
```

```

Rating.Predict <- c()
for(i in seq(1973,2015,0.5))
{
  Rating.Predict <- c(Rating.Predict,knn.decision(Congress_train,i,K_Best))
}
plot(seq(1973,2015,0.5),Rating.Predict,type = "l",xlab = "Year", ylab = "Rating", col="blue")
points(Congress_train$Year,Congress_train$Rating)

```