

1. Use "Su_raw_matrix.txt" for the following questions (30 points).

```
setwd("C:\\\\Users\\\\DELL\\\\Desktop\\\\CSC587\\\\Week2\\\\datamining-main\\\\Rscripts")
getwd()
(a) Use read.delim function to read Su_raw_matrix.txt into a variable called su.
data.file <- file.path('data', 'Su_raw_matrix.txt')
su <- read.delim(data.file, header = TRUE)
(b) Use mean and sd functions to find mean and standard deviation of Liver_2.CEL
column
mean(su[["Liver_2.CEL"]])
sd(su[["Liver_2.CEL"]])
(c) Use colMeans and colSums functions to get the average and total values of each
column.
colMeans(su)
colSums(su)
```

```
R 4.3.2 · C:/Users/DELL/Desktop/CSC587/Week2/datamining-main/Rscripts/ 
> #1. Use "Su_raw_matrix.txt" for the following questions (30 points).
> setwd("C:\\\\Users\\\\DELL\\\\Desktop\\\\CSC587\\\\Week2\\\\datamining-main\\\\Rscripts")
> getwd()
[1] "C:/Users/DELL/Desktop/CSC587/week2/datamining-main/Rscripts"
> # Load in the data set from disk.
> # (a) Use read.delim function to read Su_raw_matrix.txt into a variable called su.
> data.file <- file.path('data', 'Su_raw_matrix.txt')
> su <- read.delim(data.file, header = TRUE)
> # (b) Use mean and sd functions to find mean and standard deviation of Liver_2.CEL column
> mean(su[["Liver_2.CEL"]])
[1] 241.8246
> sd(su[["Liver_2.CEL"]])
[1] 1133.352
> # (c) Use colMeans and colSums functions to get the average and total values of each column.
> colMeans(su)
  Brain_1.CEL   Brain_2.CEL Fetal_brain_1.CEL Fetal_brain_2.CEL Fetal_liver_1.CEL Fetal_liver_2.CEL   Liver_1.CEL
[1] 204.9763      315.0924     198.3439      267.6551      209.8722      399.1482      160.8558
  Liver_2.CEL
[1] 241.8246
> colSums(su)
  Brain_1.CEL   Brain_2.CEL Fetal_brain_1.CEL Fetal_brain_2.CEL Fetal_liver_1.CEL Fetal_liver_2.CEL   Liver_1.CEL
[1] 2588031      3978357     2504290       3379413      2649846       5039645      2030966
  Liver_2.CEL
[1] 3053278
```

2. Use rnorm(n, mean = 0, sd = 1) function in R to generate 10000 numbers for the following (mean, sigma) pairs and plot histogram for each, meaning you need to change the function parameter accordingly. Then comment on how these histograms are different from each other and state the reason. (20 points)

(a) mean=0, sigma=0.2

```
sigma1 <- data.frame(X = rnorm(10000, mean = 0, sd = 0.2))
```

(b) mean=0, sigma=0.5

```
sigma2 <- data.frame(X = rnorm(10000, mean = 0, sd = 0.5))
```

The sigma 0.5 is lower/shorter and wider.

Please save your figures as image from RStudio. (Hint: to see the difference in plots you may need to set the xlim parameter in plot function to c(-5,5))

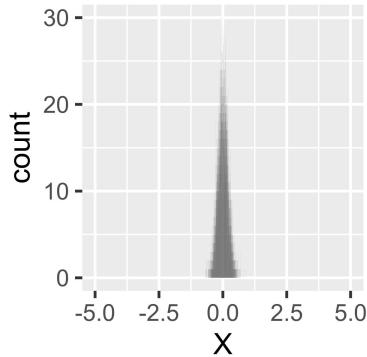
#Start visualizing data using the ggplot2 package.

```
library('ggplot2')
sigma1ggplot = ggplot(sigma1, aes(x = X)) + geom_histogram(binwidth = 0.001) +
  xlim(c(-5, 5))
```

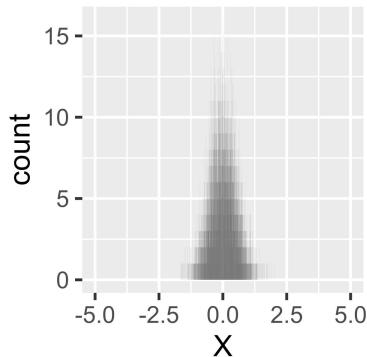
```

sigma2ggpot = ggplot(sigma2, aes(x = X)) + geom_histogram(binwidth = 0.001) +
  xlim(c(-5, 5))
ggsave("histogram_sigma1.png", plot = sigma1ggpot, width = 2, height = 2, dpi =
  5000)
ggsave("histogram_sigma2.png", plot = sigma2ggpot, width = 2, height = 2, dpi =
  5000)

```



(a) mean=0, sigma=0.2



(b) mean=0, sigma=0.5

3. Perform the steps below with "dat" dataframe which is just a sample data for you to observe how each plot function (3b through 3e) works. Notice that you need to have ggplot2 library installed on your system. Please refer slides how to install and import a library. Installation is done only once, but you need to import the library every time you need it by saying library(ggplot2). Then run the following commands for questions from 3a through 3e and observe how the plots are generated first. (20 points)

```

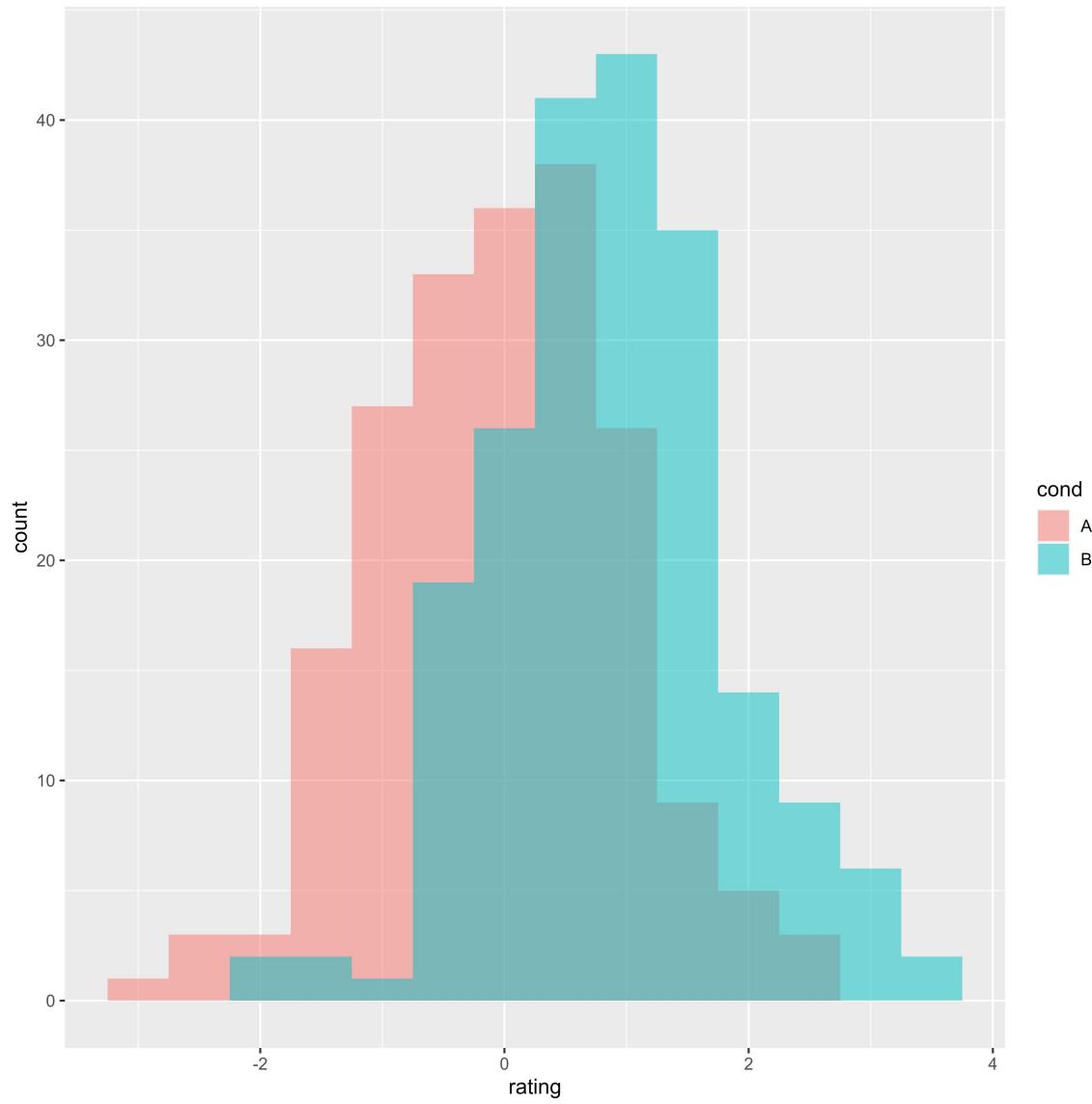
(a) dat <- data.frame(cond = factor(rep(c("A","B"), each=200)), rating =
c(rnorm(200),rnorm(200, mean=.8)))
dat <- data.frame(cond = factor(rep(c("A","B"), each=200)), rating =
c(rnorm(200),rnorm(200, mean=.8)))
(b) # Overlaid histograms
ggplot(dat, aes(x=rating, fill=cond)) + geom_histogram(binwidth=.5, alpha=.5,
position="identity")
t1 = ggplot(dat, aes(x=rating, fill=cond)) + geom_histogram(binwidth=.5, alpha=.5,
position="identity")

```

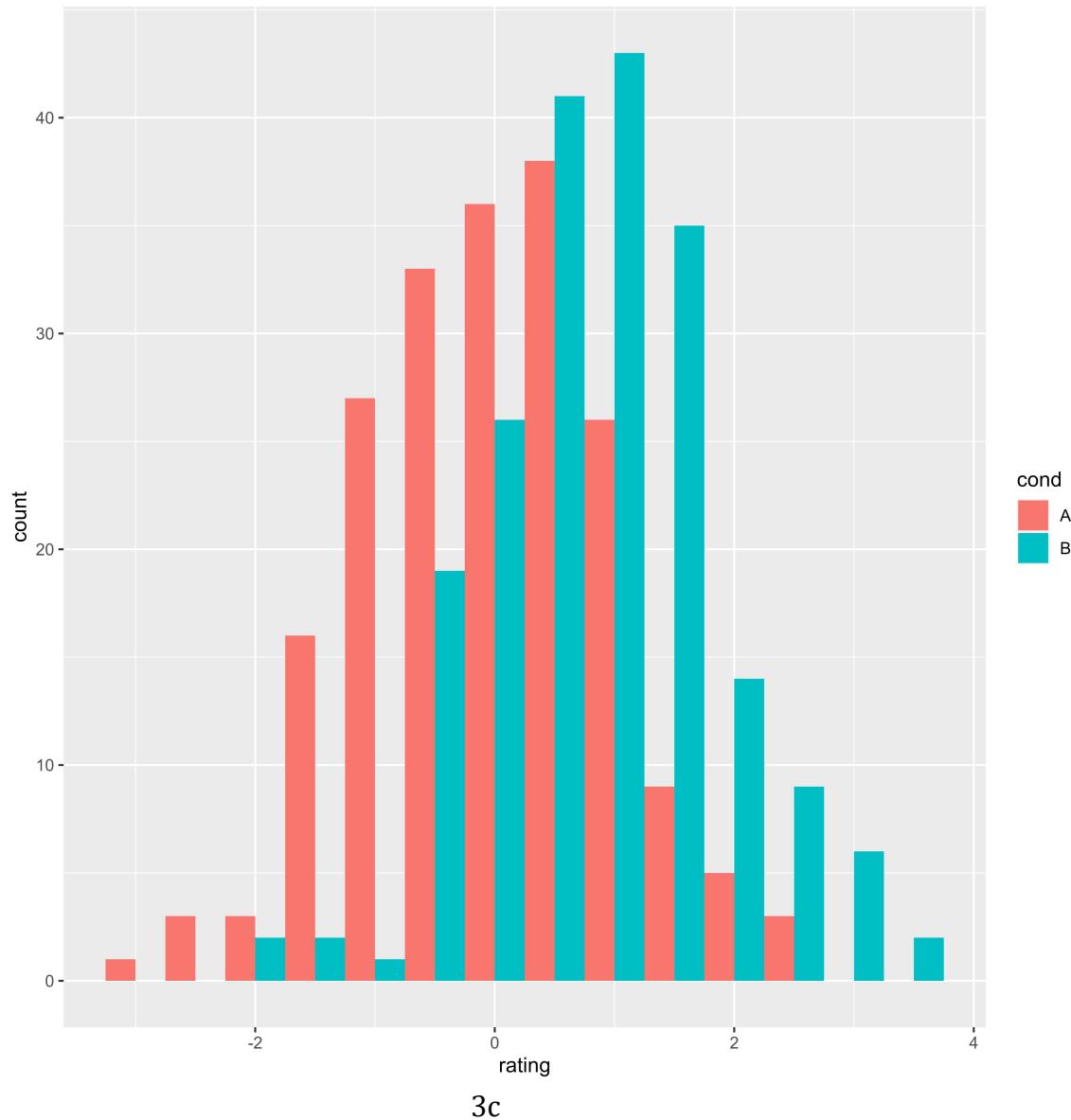
```

(c) # Interleaved histograms
ggplot(dat, aes(x=rating, fill=cond)) + geom_histogram(binwidth=.5,
position="dodge")
t2 = ggplot(dat, aes(x=rating, fill=cond)) + geom_histogram(binwidth=.5,
position="dodge")
(d) # Density plots
ggplot(dat, aes(x=rating, colour=cond)) + geom_density()
t3 = ggplot(dat, aes(x=rating, colour=cond)) + geom_density()
(e) # Density plots with semitransparent fill
ggplot(dat, aes(x=rating, fill=cond)) + geom_density(alpha=.3)
t4 = ggplot(dat, aes(x=rating, fill=cond)) + geom_density(alpha=.3)
(f) Read "diabetes_train.csv" into a variable called diabetes and apply the same
functions 3b through 3e for the mass attribute of diabetes and save the images.
(Hint: instead of cond above, use the class attribute to color your groups. When you
have fill option, your plots should show same type of chart for both groups in
different colors on the same figure. Keep in mind that diabetes and dat are both
DataFrames)
data.file <- file.path('data', 'diabetes_train.csv')
diabetes <- read.csv(data.file, header = TRUE, sep = ',')
p1 = ggplot(diabetes , aes(x=mass, fill=class)) + geom_histogram(binwidth=.5,
alpha=.5, position="identity")
p2 = ggplot(diabetes , aes(x=mass, fill=class)) + geom_histogram(binwidth=.5,
position="dodge")
p3 = ggplot(diabetes , aes(x=mass, colour=class)) + geom_density()
p4 = ggplot(diabetes , aes(x=mass, fill=class)) + geom_density(alpha=.3)
ggsave("histogram_3fb.png", plot = p1, width = 8, height = 8, dpi = 1000)
ggsave("histogram_3fc.png", plot = p2, width = 8, height = 8, dpi = 1000)
ggsave("histogram_3fd.png", plot = p3, width = 8, height = 8, dpi = 1000)
ggsave("histogram_3fe.png", plot = p4, width = 8, height = 8, dpi = 1000)
ggsave("histogram_3b.png", plot = t1, width = 8, height = 8, dpi = 1000)
ggsave("histogram_3c.png", plot = t2, width = 8, height = 8, dpi = 1000)
ggsave("histogram_3d.png", plot = t3, width = 8, height = 8, dpi = 1000)
ggsave("histogram_3e.png", plot = t4, width = 8, height = 8, dpi = 1000)

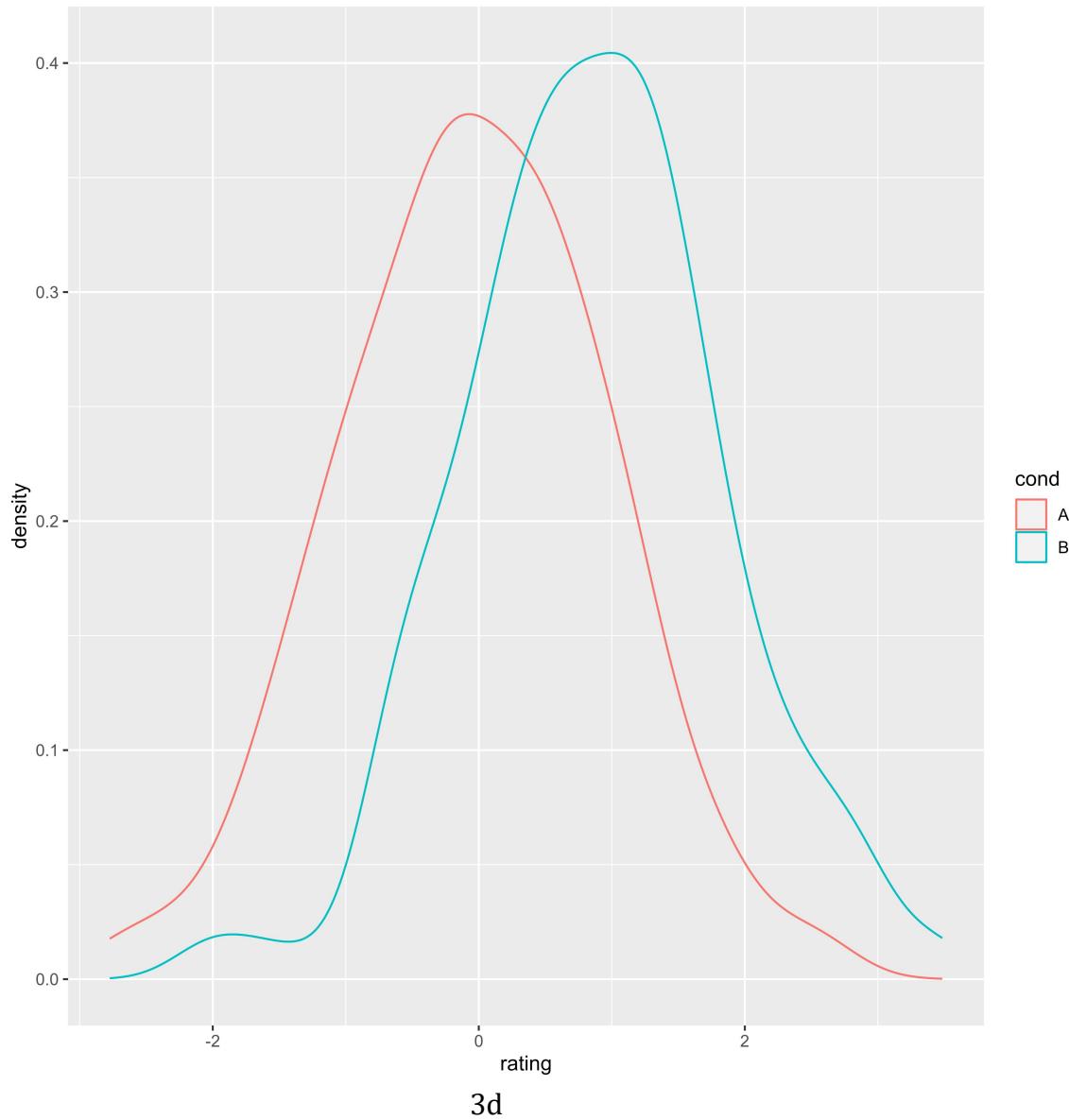
```

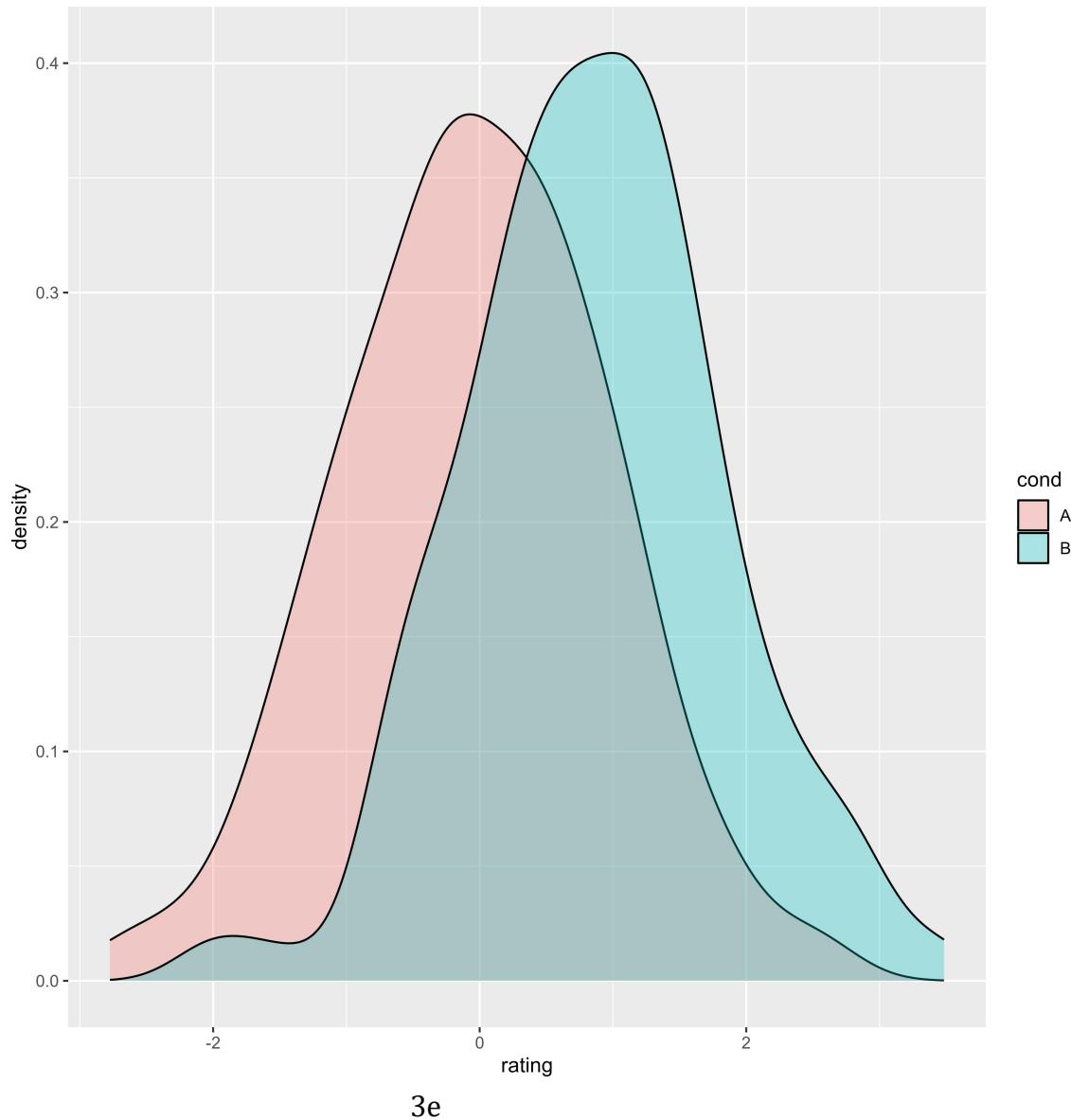


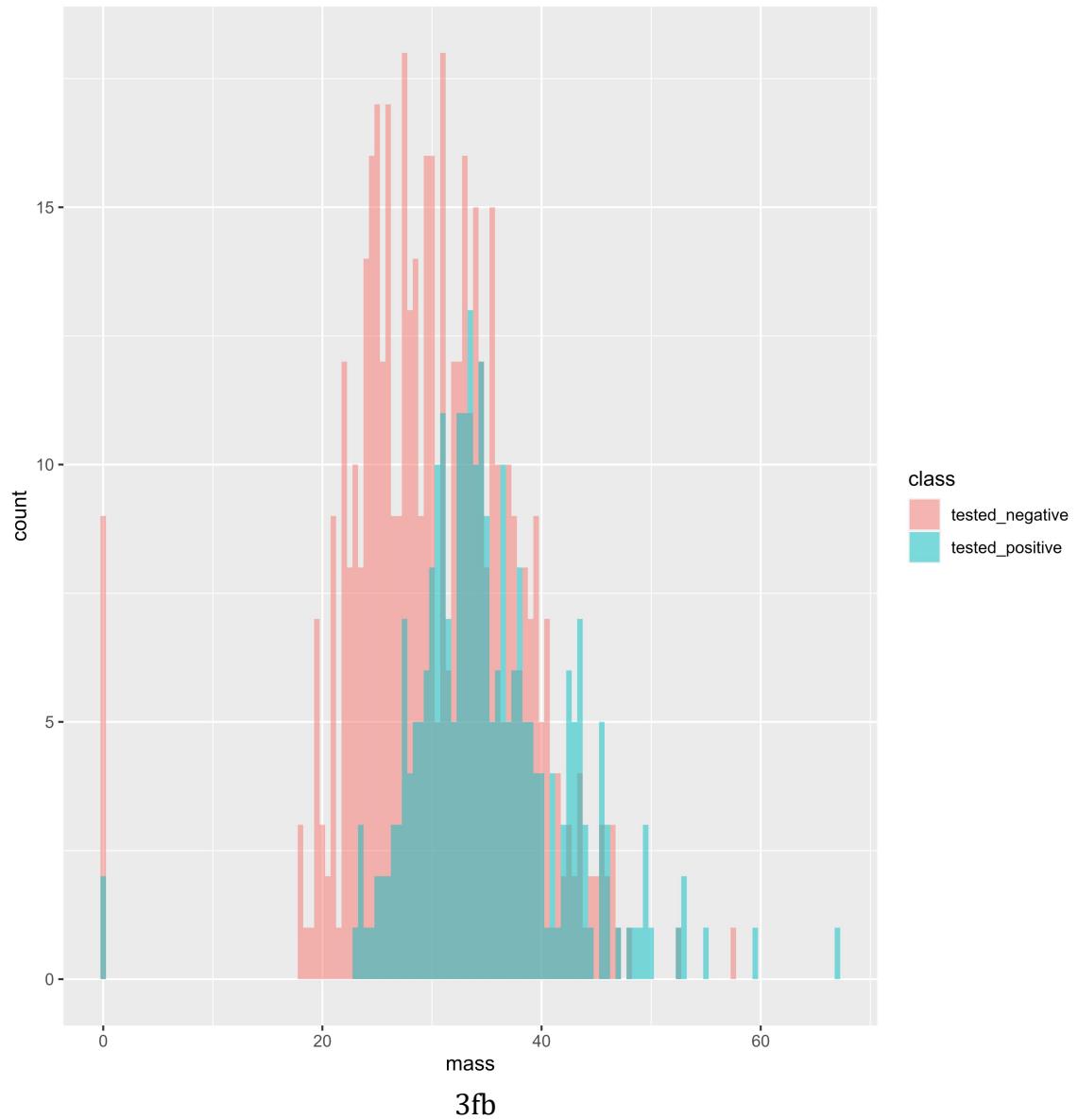
3b

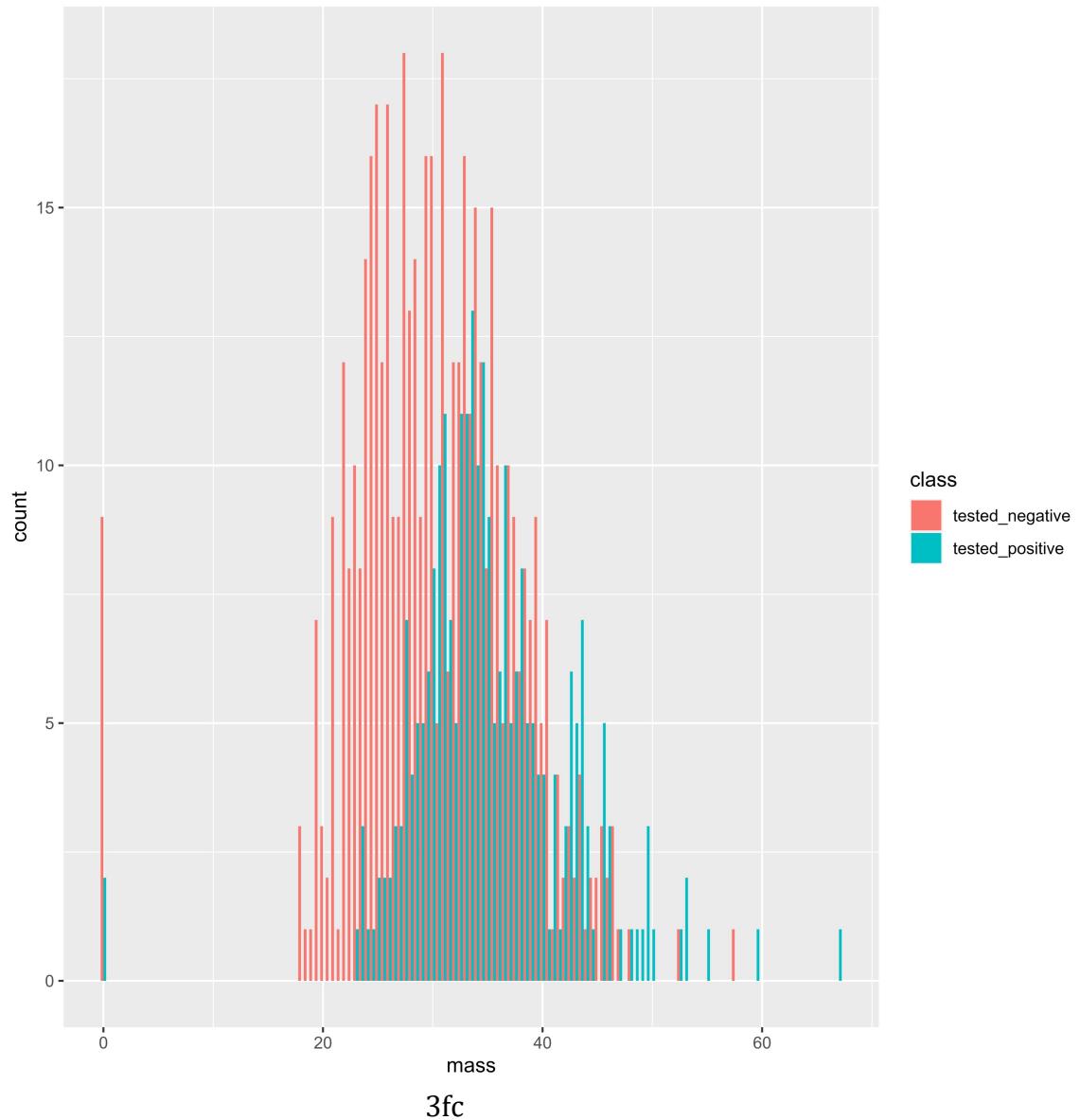


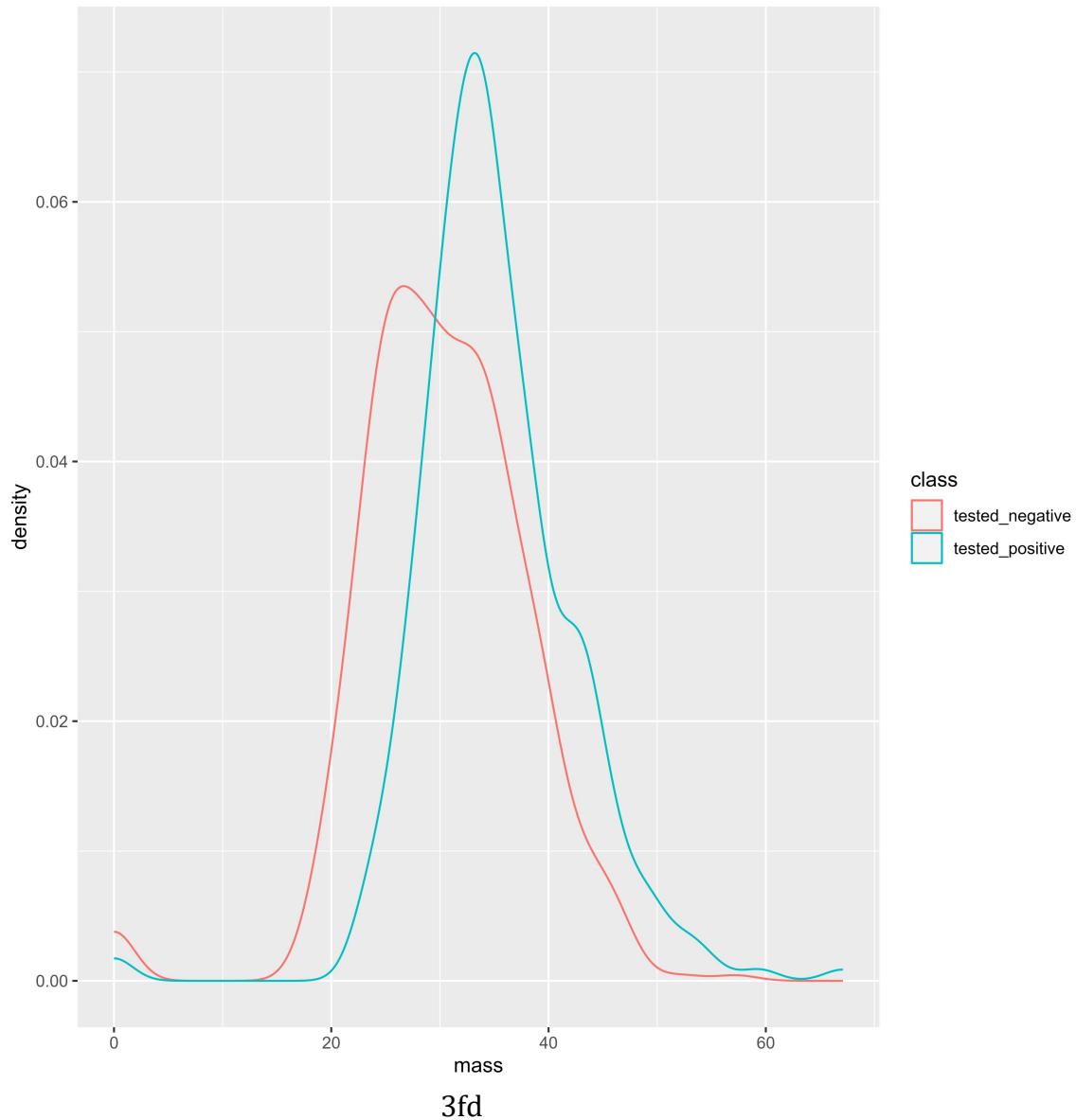
3c

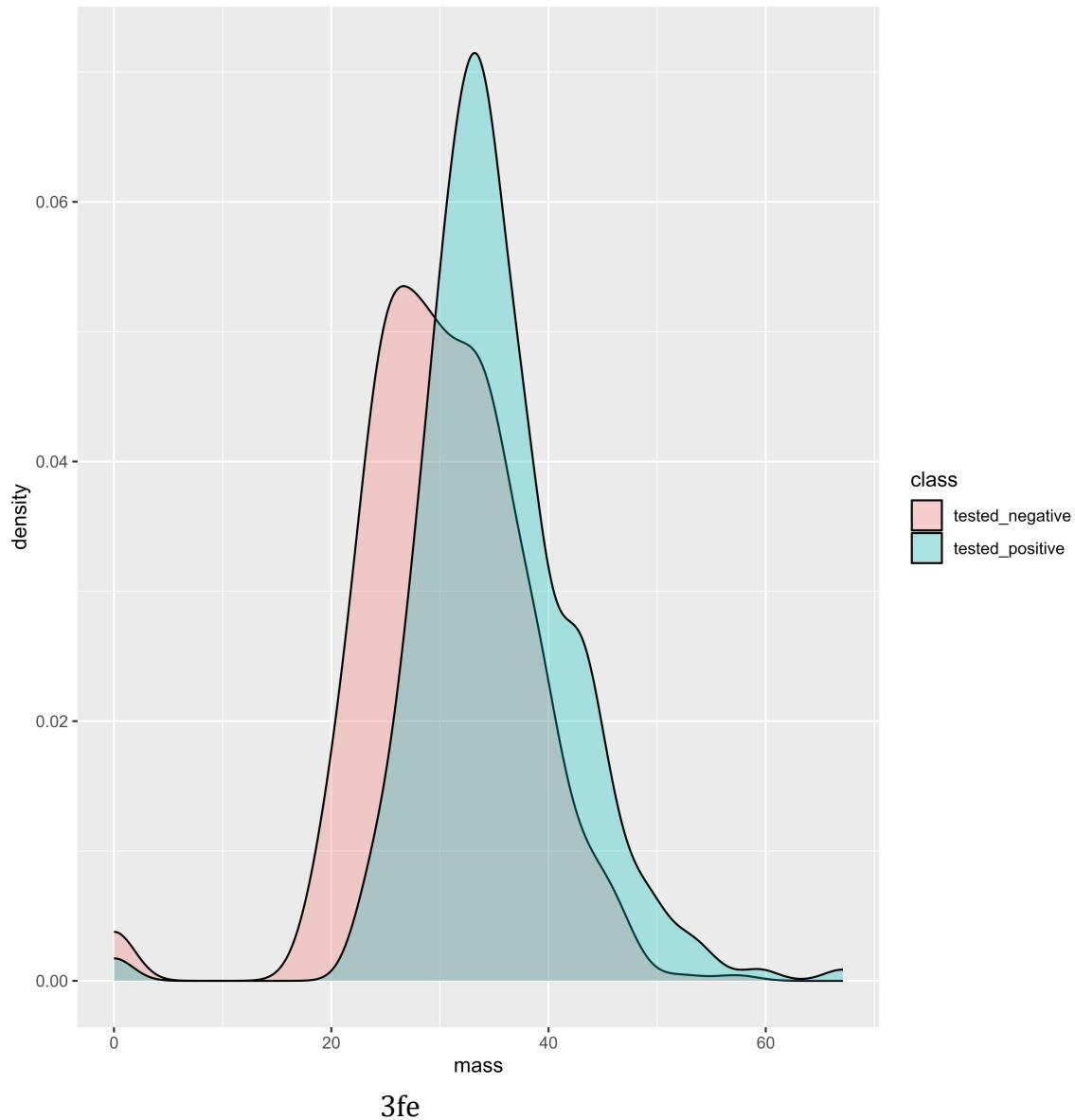












4. Read the titanic.csv file from DATA folder to a variable named passengers and perform the following steps and explain the operation very briefly (20 points):

```

data.file <- file.path('data', 'titanic.csv')
passengers <- read.csv(data.file, header = TRUE, sep = ',')
library(tidyr)
library(dplyr)
(a) passengers %>% drop_na() %>% summary()
First drops rows where any column contains a missing value from passengers. Then make the summary.
passengers %>% drop_na() %>% summary()

```

(b) passengers %>% filter(Sex == "male")

This code filters the passengers data to include only rows where the 'Sex' column is equal to "male".

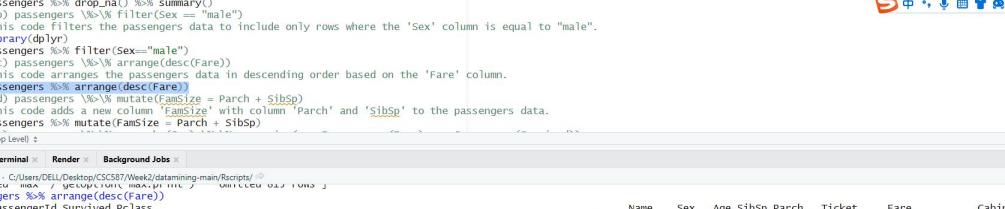
```
passengers %>% filter(Sex=="male")
```

```
① CSC587.W24_HW1_JIAN_XUR ❁ Go to/file/function ❁ Addins ❁
② Source on Save ❁ Run ❁ Help ❁ Source ❁
③ library(tidyverse)
④ passengers %>% drop_na() %>% summary()
⑤ #(b) passengers %>% filter(Sex == "male")
⑥ #This code filters the passengers data to include only rows where the 'Sex' column is equal to "male".
⑦ library(dplyr)
⑧ passengers %>% filter(Sex=="male")
⑨ #(C) passengers %>% arrange(desc(Fare))
⑩ #This code arranges the passengers data in descending order based on the 'Fare' column.
⑪ passengers %>% arrange(desc(Fare))
⑫ #(D) passengers %>% mutate(FamSize = Parch + SibSp)
⑬ #This code adds a new column 'FamSize' with column 'Parch' and 'SibSp' to the passengers data.
⑭ passengers %>% mutate(FamSize = Parch + SibSp)
⑮
⑯ (Top Level) ❁ Script ❁
```

(c) passengers %>% arrange(desc(Fare))

This code arranges the passengers data in descending order based on the 'Fare' column.

```
passengers %>% arrange(desc(Fare))
```



```
## Go to file function Addins Project: (None)
Source on Save Run Source Help
library(dplyr)
#> [1] "library(dplyr)" " "
77 passengers %>% drop_na() %>% summary()
78 #(#b) passengers %>% filter(Sex == "male")
79 #This code filters the passengers data to include only rows where the 'Sex' column is equal to "male".
80 library(dplyr)
81 passengers %>% filter(Sex=="male")
82 #(#c) passengers %>% arrange(desc(Fare))
83 #This code arranges the passengers data in descending order based on the 'Fare' column.
84 passengers %>% arrange(desc(Fare))
85 #(#d) passengers %>% mutate(FamSize = Parch + SibSp)
86 #This code adds a new column 'FamSize' with column 'Parch' and 'SibSp' to the passengers data.
87 passengers %>% mutate(FamSize = Parch + SibSp)
84:35 [Top Level] R Script
```

Console Terminal Render Background Jobs

R 4.3.2 | C:\Users\DELL\Desktop\SCS587\Week2\datamining-main\Scripts\R |

> passengers %>% arrange(desc(Fare))

	PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Cabin		
1	258	259	1	Ward, Miss. Anna	female	35.00	0	0 PC	17755	512.3292	B51 B53 B55		
2	679	680	1	1	Cardeza, Mr. Thomas Drake Martinez	male	36.00	0	1 PC	17755	512.3292	B101	
3	737	738	1	1	Lesurer, Mr. Gustave J	male	35.00	0	0 PC	17755	512.3292		
4	27	28	0	1	Fortune, Mr. Charles Alexander	male	19.00	3	2	19950	263.0000	C23 C25 C27	
5	38	39	1	1	Fortune, Miss. Alice Elizabeth	female	24.00	3	2	19950	263.0000	C23 C25 C27	
6	341	342	1	1	Fortune, Mr. Mark	male	64.00	1	4	19950	263.0000	C23 C25 C27	
7	438	439	0	1	Ryerson, Miss. Emily Borie	female	18.00	0	2	2 PC	17608	262.3750 B57	B59 B63 B66
8	311	312	1	1	Lesurer, Mr. Gustave J	male	35.00	0	0 PC	17608	262.3750 B57	B59 B63 B66	
9	742	743	1	1	Ryerson, Miss. Susan Parker "Suzette"	female	21.00	2	2 PC	17608	262.3750 B57	B59 B63 B66	
10	118	119	0	1	Baxter, Mrs. James (Helene Delaudeniere Chaput)	female	50.00	0	1 PC	17558	247.5208	B58 B60	
11	299	300	1	1	Baxter, Mr. Quige Edmond	male	24.00	0	0 PC	17558	247.5208	B58 B60	
12	380	381	1	?	Bidlo, Miss. Rosalie	female	42.00	0	0 PC	17557	227.5250		
13	557	558	0	1	Robbins, Mr. Hector	male	30.00	0	0 PC	17557	227.5250		
14	200	701	1	1	Astor, Mrs. John Jacob (Madeline Talmadge Foxe)	female	18.00	1	0 PC	17557	227.5250	C62 C64	
15	716	717	1	1	Endres, Miss. Caroline Louise	female	38.00	0	0 PC	17557	227.5250	C45	
16	527	528	0	1	Farthing, Mr. John	male	NA	0	0 PC	17483	221.7792	C95	
17	377	378	0	1	Widener, Mr. Harry Elkins	male	27.00	0	2	113503	211.5000	C82	
18	689	690	1	1	Madill, Miss. Georgette Alexandra	female	15.00	0	1	24160	211.3375	B5	
19	730	731	1	?	Allen, Miss. Elisabeth Walton	female	29.00	0	0	24160	211.3375	B5	
20	79	780	1	1	Evans, Mr. Frank	male	32.00	0	1	24160	211.3375	B3	
21	318	319	1	1	Wick, Miss. Mary Natalie	female	31.00	0	2	30326	164.0667	C7	
22	856	857	1	1	Wick, Mrs. George Donnick (Mary Hitchcock)	female	45.00	1	1	36928	164.8667		
23	268	269	1	1	Graham, Mrs. William Thompson (Edith Jenkins)	female	58.00	0	1 PC	17582	153.4625	C125	
24	332	333	0	1	Graham, Mr. George Edward	male	38.00	0	1 PC	17582	153.4625	C91	

(d) passengers %>% mutate(FamSize = Parch + SibSp)

This code adds a new column 'FamSize' with column 'Parch' and 'SibSp' to the passengers data.

```
passengers %>% mutate(FamSize = Parch + SibSp)
```

```
80 library(dplyr)
81 passengers %>% filter(Sex=="male")
82 #> #<--(c) passengers %>% arrange(desc(Fare))
83 #> #This code arranges the passengers data in descending order based on the 'Fare' column.
84 #> passengers %>% arrange(Fare)
85 #> #This code adds a new column 'FamSize' with column 'Parch' and 'SibSp' to the passengers data.
86 passengers %>% mutate(FamSize = Parch + SibSp)
87 #> #This code adds a new column 'FamSize' with column 'Parch' and 'SibSp'
88 #> #(e) passengers %>% group_by(Sex) %>% summarise(meanFare = mean(Fare), numSurv = sum(Survived))
89 #> #This code groups the passengers data by the 'Sex' column. Then it calculates the mean of 'Fare' and the sum of 'Survived' for each group ('Sex').
90 passengers %>% group_by(Sex) %>% summarise(meanFare = mean(Fare), numSurv = sum(Survived))
91 #> #3. By using quantile(), calculate 10th,30th,50th,60th percentiles of skin attribute of diabetes data. (10 points)
88:1 (Top Level) : 
```

R 4.3.2 C:\Users\DELL\Desktop\SCS878\Week2\datamining\main\Scripts.R

```
> #This code adds a new column 'FamSize' with column 'Parch' and 'SibSp' to the passengers data.
> passengers %>% mutate(FamSize = Parch + SibSp)
```

PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Cabin	Embarked	FamSize	
1	0	1	3	Braund, Mr. Owen Harris	male	22.0	1	0	A/5 21171	7.2500	S	1	1
2	1	2	1	Cumings, Mrs. John Bradley (Florence Briggs Thayer)	female	38.0	1	0	PC 17599	71.2833	C85	C	1
3	2	3	1	Heikkinen, Miss. Laina	female	26.0	0	0	STON/O2. 3101282	7.9250	S	0	0
4	3	4	1	Futrelle, Mrs. Jacques Heikkinen (Lily May Langford)	female	35.0	1	0	113803	53.1300	CL123	S	1
5	4	5	0	Allen, Mr. William Henry	male	35.0	1	0	341470	8.0500	S	0	0
6	5	0	3	Moran, Mr. James	male	NA	NA	0	303077	8.4583	Q	0	0
7	6	7	0	McCarthy, Mr. Timothy J	male	54.0	0	0	17463	51.8625	E46	S	0
8	7	8	0	Palsson, Master. Gosta Leonard	male	2.0	3	1	349909	21.0750	S	4	0
9	8	9	1	Johnson, Mrs. Oscar W (Elisabeth Vilhelmina Berg)	female	27.0	0	2	347742	11.1333	S	2	0
10	9	10	1	Nasser, Mrs. Nicholas (Adele) Achem	female	14.0	1	0	237736	30.0708	C	1	0
11	10	11	1	Sandstrom, Miss. Marguerite Rut	female	4.0	1	1	PP 9549	16.7000	G6	S	2
12	11	12	1	Bonnell, Miss. Elizabeth	female	38.0	0	0	131783	26.5500	CL03	S	0
13	12	13	0	Saunders, Mr. Edward Wilding	male	32.0	0	0	A/5 23349	3.0000	S	0	0
14	13	14	0	Anderson, Mr. Anders Johnson	male	39.0	1	5	347082	31.2750	S	6	0
15	14	15	0	Vestrom, Miss. Hilda Andersdotter	female	14.0	0	0	350406	7.8542	S	0	0
16	15	16	1	Hewlett, Mrs. (Mary D Kingcome)	female	55.0	0	0	248706	16.0000	S	0	0
17	16	17	0	Rice, Master. Eugene	male	2.0	4	1	382652	29.1250	Q	5	0
18	17	18	1	Williams, Mr. Charles Eugene	male	NA	0	0	244373	13.0000	S	0	0
19	18	19	0	Vander Planke, Mrs. Julius (Emilia Maria Vandeplante)	female	31.0	1	0	34363	19.0000	S	1	0
20	19	20	1	Masselman, Mrs. Ernesto	female	NA	0	0	261320	7.2500	C	0	0
21	20	21	0	Dyney, Mr. Joseph J	male	35.0	0	0	239865	26.0000	S	0	0
22	21	22	1	Beesley, Mr. Lawrence	male	34.0	0	0	248688	13.0000	D56	S	0
23	22	23	1	McGowan, Miss. Anna "Annie"	female	15.0	0	0	310923	8.0929	Q	0	0
24	23	24	1	Slater, William Thomas	male	38.0	0	0	112788	35.1500	A6	S	0

(e) passengers %>% group_by(Sex) %>% summarise(meanFare = mean(Fare), numSurv = sum(Survived))

This code groups the passengers data by the 'Sex' column. Then It then calculates the mean of 'Fare' and the sum of 'Survived' for each group ('Sex').

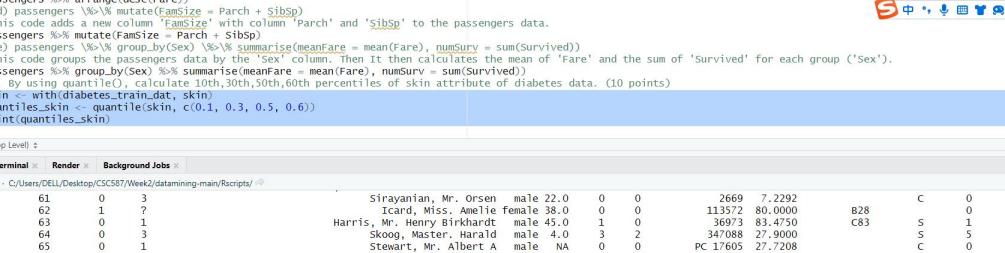
```
passengers %>% group_by(Sex) %>% summarise(meanFare = mean(Fare),  
numSurv = sum(Survived))
```

5. By using quantile(), calculate 10th,30th,50th,60th percentiles of skin attribute of diabetes data. (10 points)

```
skin <- with(diabetes_train_dat, skin)
```

```
quantiles_skin <- quantile(skin, c(0.1, 0.3, 0.5, 0.6))
```

```
print(quantiles_skin)
```



The screenshot shows a Shiny application running in RStudio. The application has a header "W24_HW1_JIAN_XU.Rmd" and a sidebar with tabs for "Source on Save", "Run", "Source", and "Project (None)". The main panel displays the R code for the Shiny app, which includes data manipulation with dplyr and ggplot2, and a conditional statement for the UI based on the number of passengers.

```
library(shiny)
library(dplyr)
library(ggplot2)

# This code reads the titanic dataset and creates a new column 'FamSize' = Parch + SibSp
# This code groups the passengers data by the 'Sex' column. Then It then calculates the mean of 'Fare' and the sum of 'Survived' for each group ('Sex').
# This code groups the passengers data by the 'Sex' column. Then It then calculates the mean of 'Fare' and the sum of 'Survived' for each group ('Sex').

#> #> passengers %>% group_by(Sex) %>% summarise(meanFare = mean(Fare), numSurv = sum(Survived))
#> #> This code groups the passengers data by the 'Sex' column. Then It then calculates the mean of 'Fare' and the sum of 'Survived' for each group ('Sex').
#> #> passengers %>% group_by(Sex) %>% summarise(meanFare = mean(Fare), numSurv = sum(Survived))
#> #> By using quantile(), calculate 10th,30th,50th,60th percentiles of skin attribute of diabetes data. (10 points)
#> skin <- with(diabetes_train_dat, skin)
#> quantiles_skin <- quantile(skin, c(0.1, 0.3, 0.5, 0.6))
#> print(quantiles_skin)
#> 95t
```

The console output shows the results of the quantile calculation:

Sex	meanFare	numSurv
<chr>	<dbl>	<int>
1 female	44.5	233

The UI output shows the conditional rendering of the plot based on the number of passengers:

```
if (n == 1) {
```

The plot area shows a scatter plot of Age vs. Fare for female passengers.