

STAT-S 520 Data Analysis Project

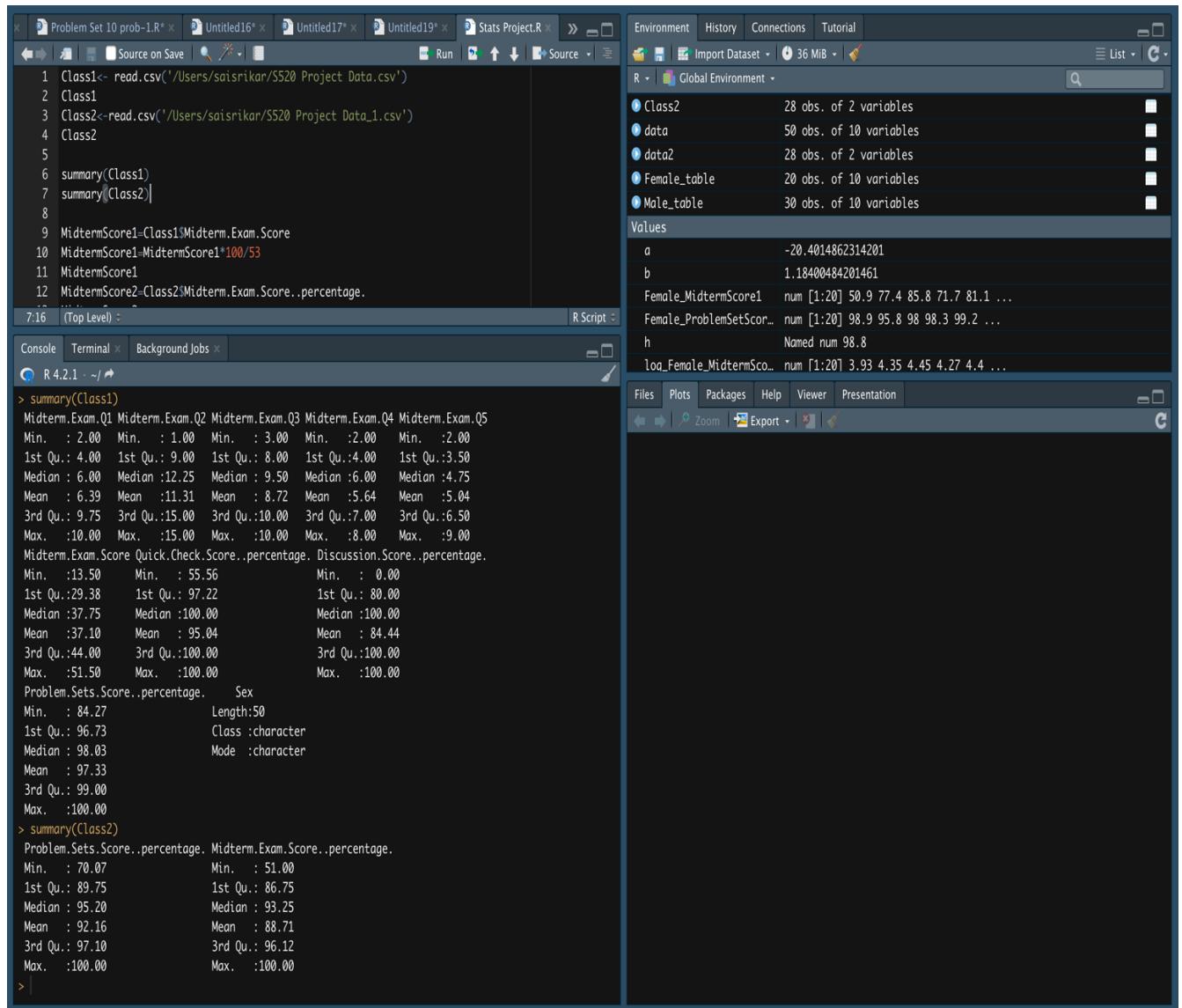
SAI SRIKAR GANDHE¹, MANASA GUDISE¹

1. **Research question:** Compare these two classes' performance in the first half of Fall semester.

Topic: Descriptive Statistics

Purpose of the Topic: To compare two classes performance and see which class's performance is better which may help the instructor to decide which one of the classes students need more attention and help.

Let's consider the midterm scores of both classes to see which class performance is better. Let's see the summary of both classes



The screenshot shows the RStudio interface with the following details:

- Script Editor:** Displays R code for reading CSV files, creating objects for Class1 and Class2, calculating MidtermScore1 and MidtermScore2, and printing the summary of Class1 and Class2.
- Console:** Shows the output of the R code, including the summary statistics for Class1 and Class2.
- Environment View:** Shows the global environment with various objects and their characteristics.
- Plots:** No plots are present in the interface.

```
1 Class1<- read.csv('/Users/saisrikan/S520 Project Data.csv')
2 Class1
3 Class2<-read.csv('/Users/saisrikan/S520 Project Data_1.csv')
4 Class2
5
6 summary(Class1)
7 summary(Class2)
8
9 MidtermScore1=Class1$Midterm.Exam.Score
10 MidtermScore1=MidtermScore1*100/53
11 MidtermScore1
12 MidtermScore2=Class2$Midterm.Exam.Score..percentage.
> summary(Class1)
Midterm.Exam.Q1 Midterm.Exam.Q2 Midterm.Exam.Q3 Midterm.Exam.Q4 Midterm.Exam.Q5
Min. : 2.00 Min. : 1.00 Min. : 3.00 Min. : 2.00 Min. : 2.00
1st Qu.: 4.00 1st Qu.: 9.00 1st Qu.: 8.00 1st Qu.: 4.00 1st Qu.: 3.50
Median : 6.00 Median :12.25 Median : 9.50 Median : 6.00 Median : 4.75
Mean : 6.39 Mean :11.31 Mean : 8.72 Mean : 5.64 Mean : 5.04
3rd Qu.: 9.75 3rd Qu.:15.00 3rd Qu.:10.00 3rd Qu.: 7.00 3rd Qu.: 6.50
Max. :10.00 Max. :15.00 Max. :10.00 Max. : 8.00 Max. : 9.00
Midterm.Exam.Score.Quick.Check.Score..percentage. Discussion.Score..percentage.
Min. :13.50 Min. : 55.56 Min. : 0.00
1st Qu.:29.38 1st Qu.: 97.22 1st Qu.: 80.00
Median :37.75 Median :100.00 Median :100.00
Mean :37.10 Mean : 95.04 Mean : 84.44
3rd Qu.:44.00 3rd Qu.:100.00 3rd Qu.:100.00
Max. :51.50 Max. :100.00 Max. :100.00
Problem.Sets.Score..percentage. Sex
Min. : 84.27 Length:50
1st Qu.: 96.73 Class :character
Median : 98.03 Mode :character
Mean : 97.33
3rd Qu.: 99.00
Max. :100.00
> summary(Class2)
Problem.Sets.Score..percentage. Midterm.Exam.Score..percentage.
Min. : 70.07 Min. : 51.00
1st Qu.: 89.75 1st Qu.: 86.75
Median : 95.20 Median : 93.25
Mean : 92.16 Mean : 88.71
3rd Qu.: 97.10 3rd Qu.: 96.12
Max. :100.00 Max. :100.00
```

As we are considering midterm scores of both classes, let ‘MidtermScore1’ be the midterm score percentage of class1 and ‘MidtermScore2’ be the midterm score percentage of class2. Now, let’s look at the numerical summaries of midterm scores of class1 and class2

The screenshot shows the RStudio interface. The R console window displays R code and its output. The code includes reading datasets, calculating percentages, and summarizing data for Class1 and Class2. The Global Environment window shows various objects and their types and values. The File menu is visible at the top.

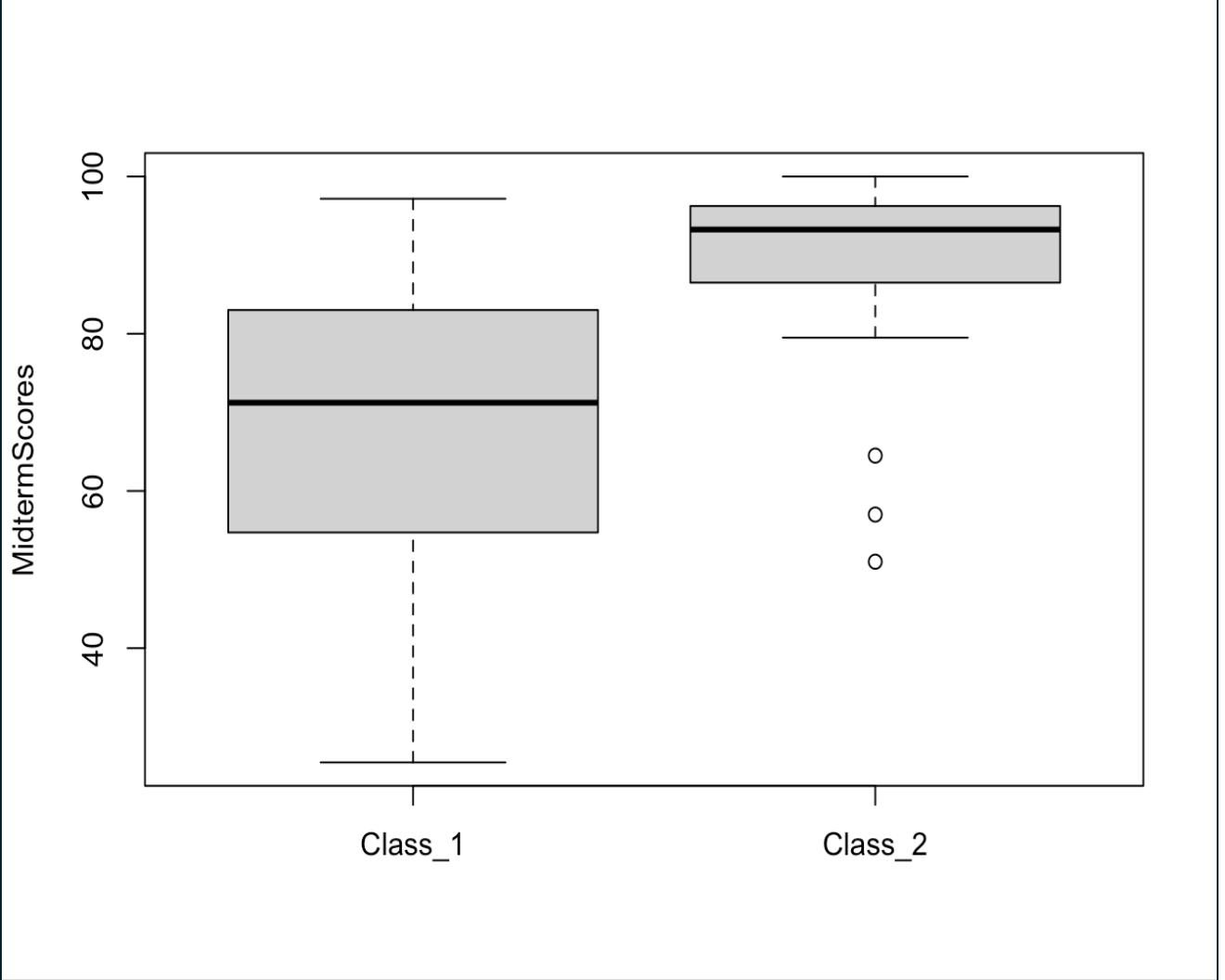
```

8
9 MidtermScore1<-Class1$Midterm.Exam.Score
10 MidtermScore1<-MidtermScore1*100/53
11 MidtermScore1
12 MidtermScore2<-Class2$Midterm.Exam.Score..percentage.
13 MidtermScore2
14
15 summary(MidtermScore1)
16 summary(MidtermScore2)
17
18 boxplot(MidtermScore1,MidtermScore2,names=c("Class_1","Class_2"),ylab="MidtermScores")
19
20 hist(density(MidtermScore2))
21
22 (Top Level) ->
R Script
Console Terminal Background Jobs
R 4.2.1 ->
Median : 98.03 Mode :character
Mean : 97.33
3rd Qu.: 99.00
Max. :100.00
> summary(Class2)
Problem.Sets.Score..percentage. Midterm.Exam.Score..percentage.
Min. : 70.07 Min. : 51.00
1st Qu.: 89.75 1st Qu.: 86.75
Median : 95.20 Median : 93.25
Mean : 92.16 Mean : 88.71
3rd Qu.: 97.10 3rd Qu.: 96.12
Max. :100.00 Max. :100.00
> MidtermScore1<-Class1$Midterm.Exam.Score
> MidtermScore1<-MidtermScore1*100/53
> MidtermScore1
[1] 62.26415 50.94340 57.54717 81.13208 83.01887 77.35849 72.64151 77.35849 85.84906 71.69811
[11] 81.13208 83.01887 89.62264 83.01887 91.50943 93.39623 86.79245 79.24528 94.33962 84.90566
[21] 95.28302 96.22642 92.45283 94.33962 97.16981 52.83019 61.32075 64.15094 60.37736 66.98113
[31] 81.13208 69.81132 81.13208 52.83019 80.18868 54.71698 59.43396 70.75472 69.81132 50.00000
[41] 54.71698 60.37736 52.83019 33.01887 57.54717 53.77358 49.05660 33.96226 25.47170 41.50943
> MidtermScore2<-Class2$Midterm.Exam.Score..percentage.
> MidtermScore2
[1] 57.0 96.0 87.0 86.0 99.5 93.0 100.0 96.5 94.5 90.0 93.5 88.5 51.0 93.5 85.0
[16] 94.0 83.0 91.0 96.5 97.0 87.0 95.0 96.0 99.0 79.5 96.0 97.5 91.5
> summary(MidtermScore1)
Min. 1st Qu. Median Mean 3rd Qu. Max.
25.47 55.42 71.23 70.00 83.02 97.17
> summary(MidtermScore2)
Min. 1st Qu. Median Mean 3rd Qu. Max.
51.00 86.75 93.25 88.71 96.12 100.00
>

```

From above summary statistics we can see that the mean and median of class2 midterm scores is higher than the class1 midterm scores. But we cannot say that for sure by just looking at numerical summary statistics. So, let’s further investigate. Let’s draw some boxplots for comparisons of midterm scores of class1 and class2. Below is a box plots of midterm scores of both classes

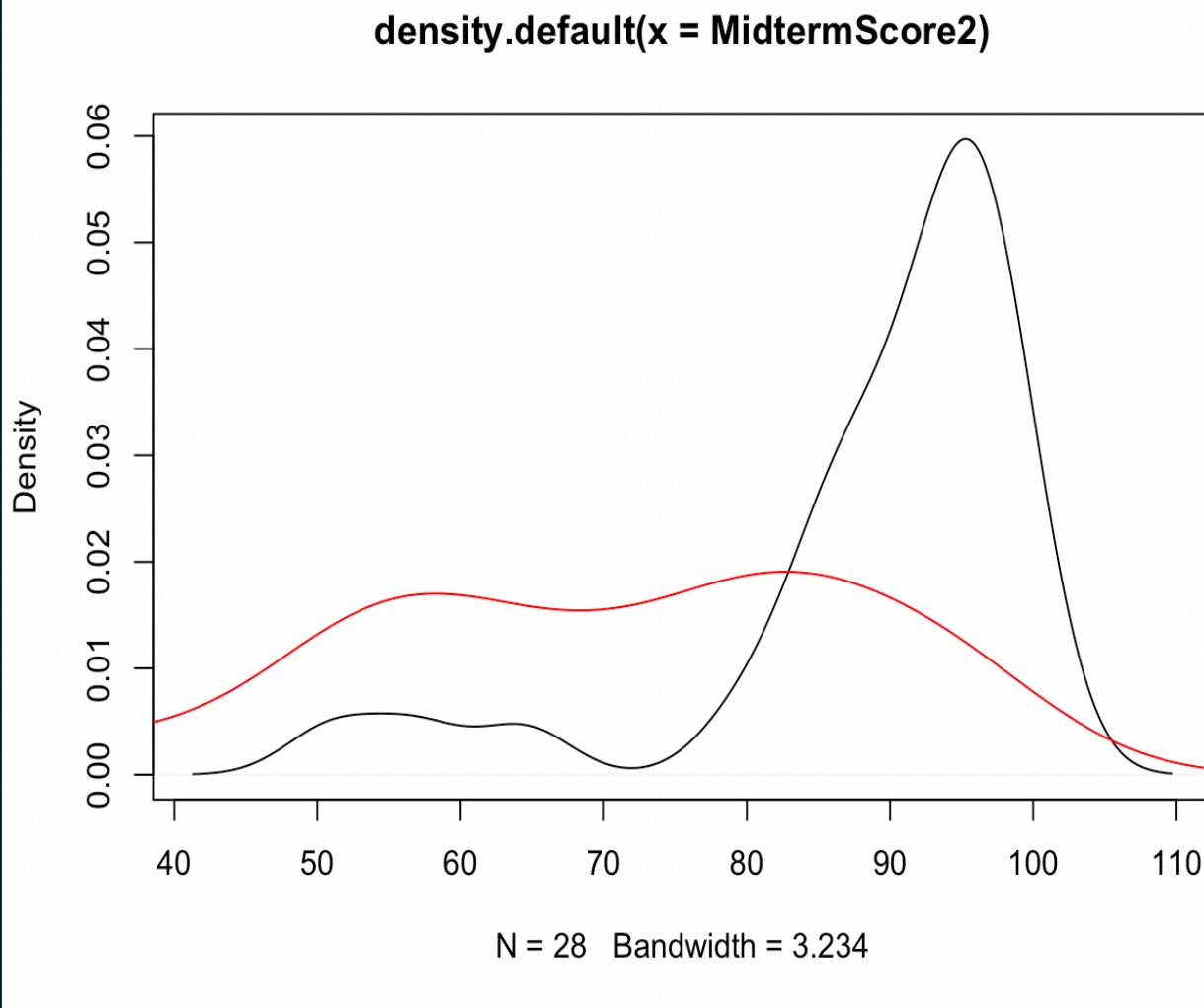
```
boxplot(MidtermScore1,MidtermScore2,names=c("Class_1","Class_2"),ylab="MidtermScores")
```



From the above boxplot, we can see that the median of the midterm scores of class1 lie around 71 and the median of the midterm scores of class2 lie around 93. We can also see that there are only three students in class2 whose midterm scores are less than 70%. Most of the students' scores in class2 lie between 86% and 100% whereas most of the students score in class1 lie between 55% and 97%.

Let's draw density plot of midterm scores of both classes

```
plot(density(MidtermScore2))
lines(density(MidtermScore1), col="red")
```



We can see that the distribution of midterm scores of class1 (red) is shifted to the left compared to midterm scores of class 2(black).

To find out more we can do a hypothesis test.

Hypothesis test:

Here we are considering midterm scores of both classes so it's a two-sample location problem and the parameter of interest is ' Δ ' where $\Delta = \mu_2 - \mu_1$. Let μ_1 be population mean of the midterm scores of class1 and μ_2 be the population mean of midterm scores of class2. So, we want to see if the population mean of the class2 is greater than mean of class2. So, this becomes our alternative hypothesis i.e, $\mu_2 - \mu_1 > 0$ which says it is a right tailed test. Null hypothesis is exactly opposite to alternative hypothesis.

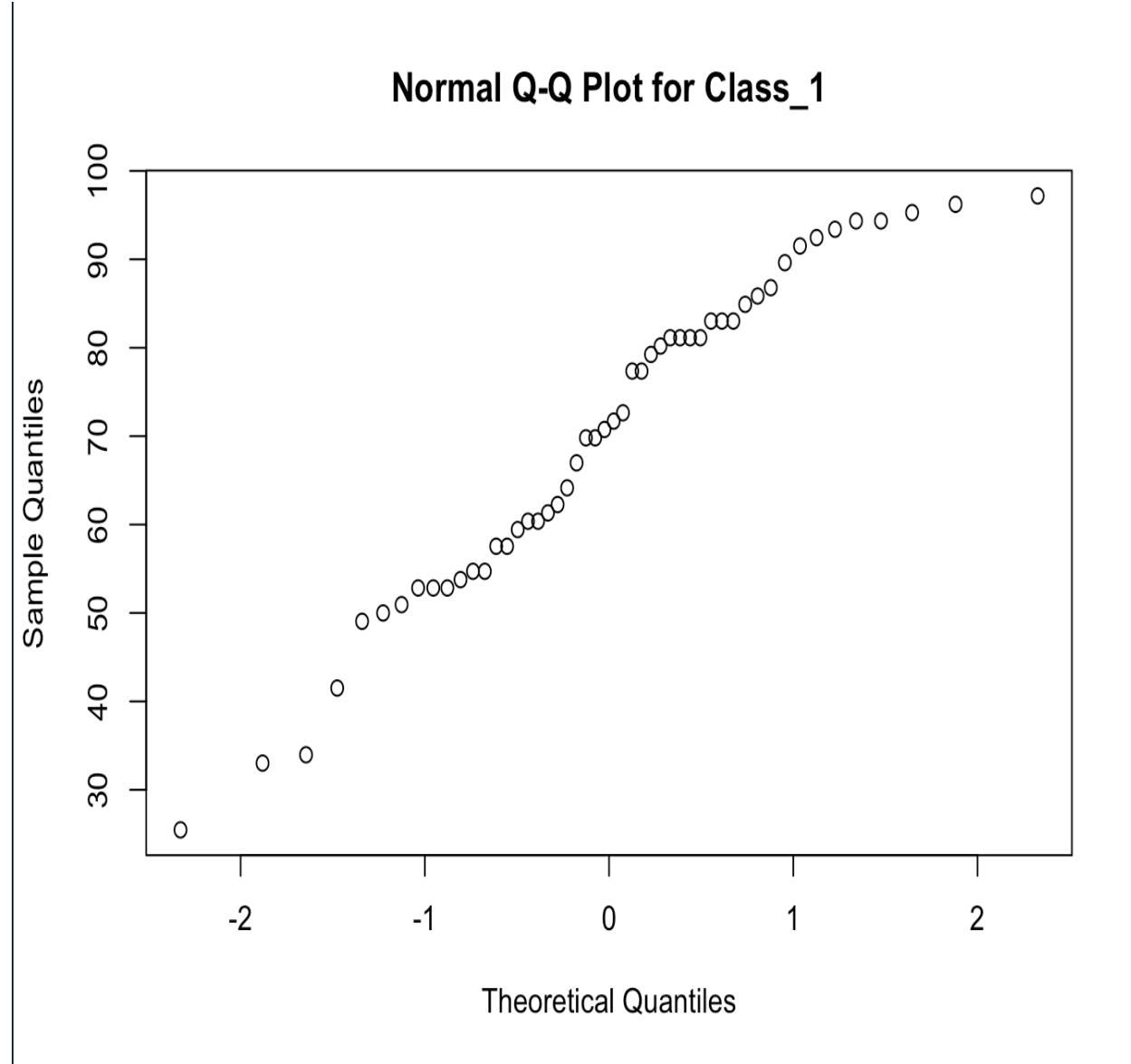
$$H_0: \Delta \leq 0$$

$$H_1: \Delta > 0$$

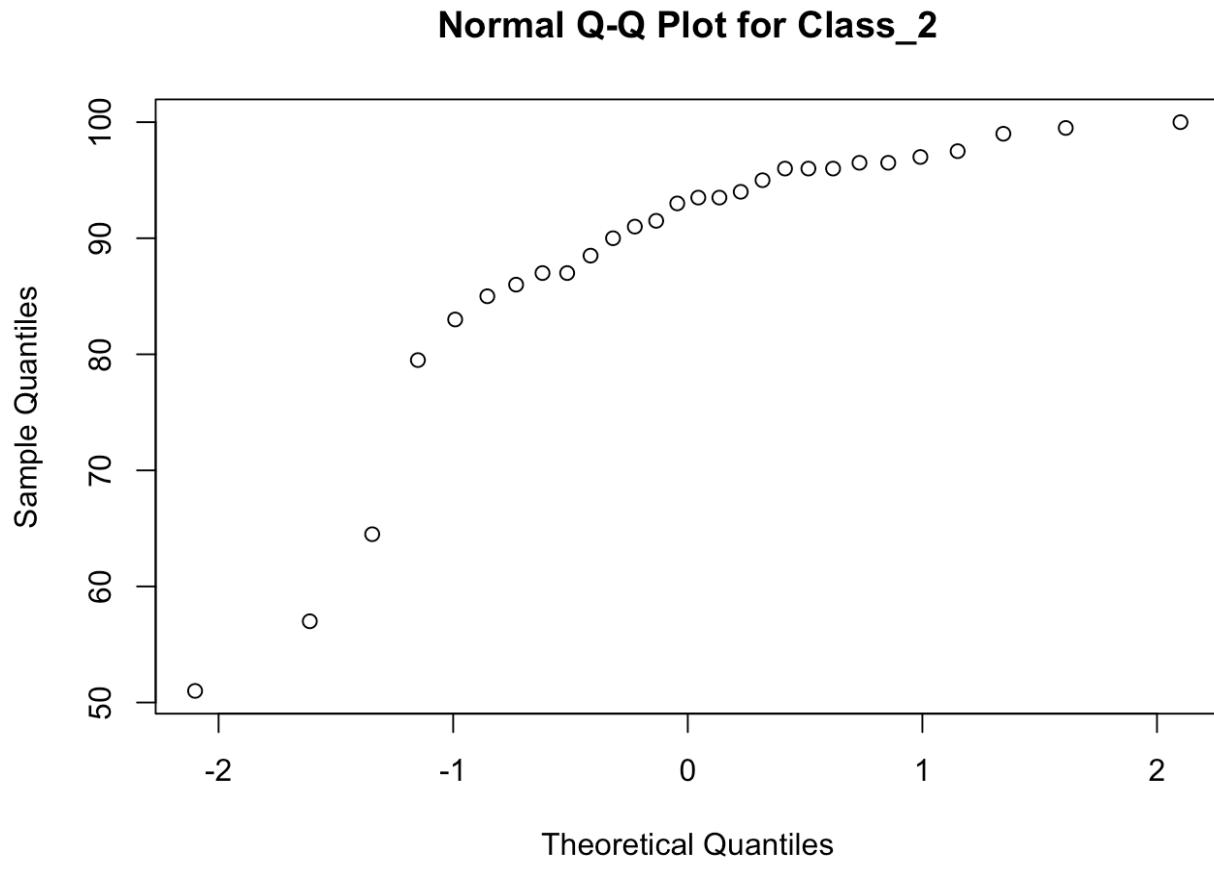
Now, let's draw QQ plots of midterm scores of both classes to see if the distribution is normal.

QQ plot of class1 midterm scores:

```
qqnorm(MidtermScore1,main='Normal Q-Q Plot for Class_1')
qqnorm(MidtermScore2,main='Normal Q-Q Plot for Class_2')
```



QQ plot of class2 midterm scores:



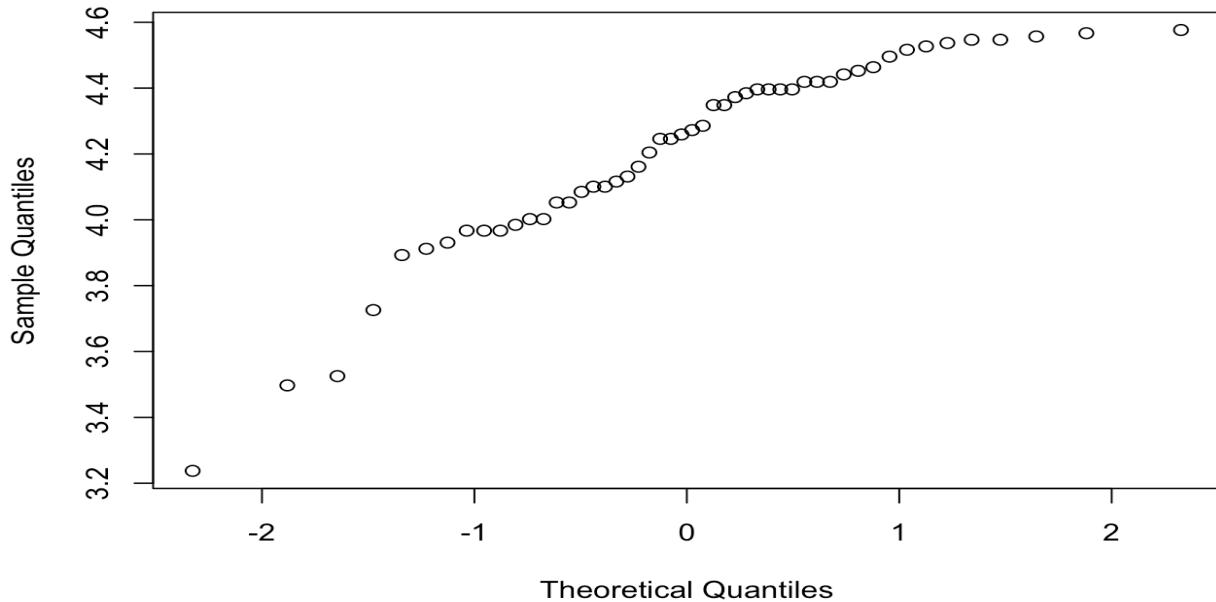
From above plots we can see that the QQ plot for class1 midterm scores looks almost like a straight line and the QQ plot for class2 midterm scores is more like curved line than a straight line. Now, let us transform the data and plot the graphs. Below is a QQ plot for log transformed midterm scores of class1.

RCode:

```
# using log
log_MidtermScore1=log(MidtermScore1)
log_MidtermScore2=log(MidtermScore2)

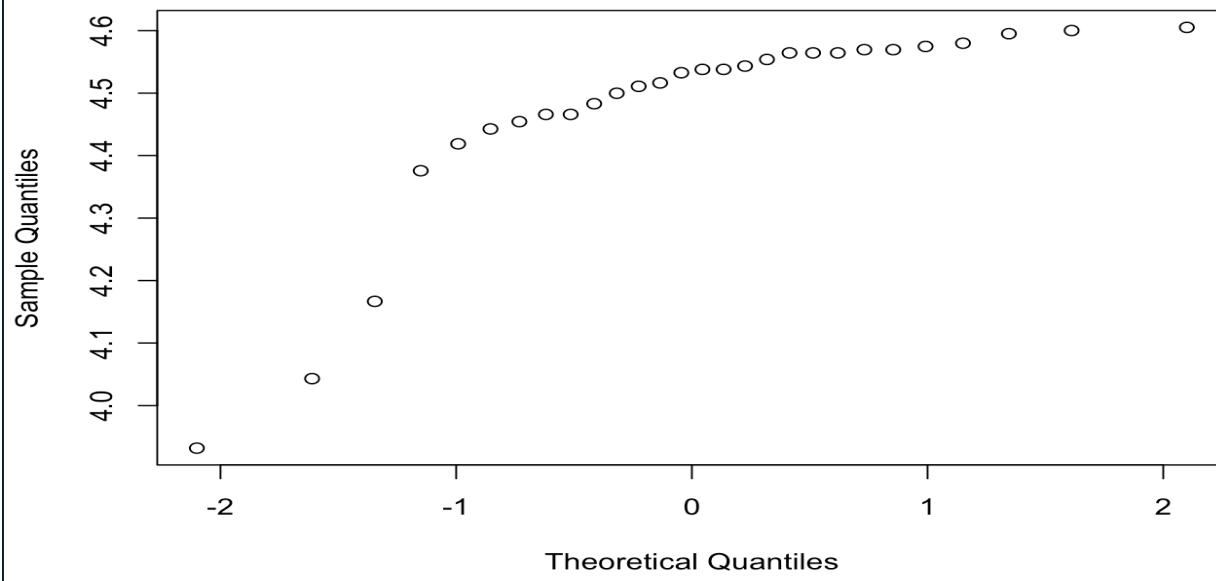
qqnorm(log_MidtermScore1,main='Normal Log Q-Q Plot for Class_1')
qqnorm(log_MidtermScore2,main='Normal Log Q-Q Plot for Class_2')
```

Normal Log Q-Q Plot for Class_1



QQ plot of Log transformed Midterm Scores of class2:

Normal Log Q-Q Plot for Class_2



The log transformation made the QQ plot of class1 more curved so let's stick with the untransformed data. Assuming that samples are IID samples from two independent normal populations, we can do Welch's two- sample t-test.

```
37
38 t.test(MidtermScore2, MidtermScore1, alternative = "greater")
39
40
41
```

38:59 (Top Level) R Script

Console Terminal Background Jobs

R 4.2.1 · ~/

```
> boxplot(MidtermScore1, MidtermScore2, names=c("Class_1", "Class_2"), ylab="MidtermScores")
> plot(density(MidtermScore2))
> lines(density(MidtermScore1), col="red")
> qqnorm(MidtermScore1, main='Normal Q-Q Plot for Class_1')
> qqnorm(MidtermScore2, main='Normal Q-Q Plot for Class_2')
> log_MidtermScore1=log(MidtermScore1)
> log_MidtermScore2=log(MidtermScore2)
> qqnorm(log_MidtermScore1, main='Normal Log Q-Q Plot for Class_1')
> qqnorm(log_MidtermScore2, main='Normal Log Q-Q Plot for Class_2')
> t.test(MidtermScore2, MidtermScore1, alternative ="greater")
```

Welch Two Sample t-test

```
data: MidtermScore2 and MidtermScore1
t = 5.3926, df = 73.235, p-value = 4.065e-07
alternative hypothesis: true difference in means is greater than 0
95 percent confidence interval:
12.93295      Inf
sample estimates:
mean of x mean of y
88.71429 70.00000
```

> |

From above t-test we can see that the p-value is $4.065e-07$ which is very small, so it means we have evidence against the null hypothesis. So, we need to reject null hypothesis. So, our alternative hypothesis becomes true. That is class2 performance is better than class1 performance when taken midterm scores into account. From above, we can see that the 95% confidence interval is $(12.93295, \infty)$.

2. Research question: Does gender make a difference in student's performance on assignments and midterm exam in S520?

Topic: Inferential Statistics

Purpose of the Topic: To see if male students or female students perform better in assignments and midterm exam.

Firstly, let's see summary of midterm scores and problem set scores of female and male students

The screenshot shows the RStudio interface with the following details:

- Environment View:** Shows the global environment with two data frames:
 - Female_table**: 20 obs. of 10 variables
 - Male_table**: 30 obs. of 10 variables
- Values View:** Displays variable definitions and their values.
- Console View:** Shows the R script and its output. The script reads data from 'Class1' and performs the following steps:
 - Creates 'Female_table' for female students.
 - Calculates 'Female_MidtermScore1' by multiplying 'Midterm.Exam.Score' by 100/53.
 - Prints the summary of 'Female_MidtermScore1'.
 - Creates 'Female_ProblemSetScore1' by multiplying 'Problem.Sets.Score..percentage' by 100.
 - Prints the summary of 'Female_ProblemSetScore1'.
 - Creates 'Male_table' for male students.
 - Calculates 'Male_MidtermScore1' by multiplying 'Midterm.Exam.Score' by 100/53.
 - Prints the summary of 'Male_MidtermScore1'.
 - Creates 'Male_ProblemSetScore1' by multiplying 'Problem.Sets.Score..percentage' by 100.
 - Prints the summary of 'Male_ProblemSetScore1'.

The screenshot shows the RStudio interface with the following components:

- Top Bar:** Shows tabs for "Problem Set 10 prob-1.R" (active), "Untitled16.R", "Untitled17.R", "Untitled19.R", and "Stats Project.R".
- Source Editor:** Displays R script code for calculating average scores for female and male students.
- Console:** Shows the output of the R script, listing individual student scores followed by summary statistics for both male and female students.
- Environment View:** Shows the global environment with variables like Female_table, Male_table, and various score summaries.
- Bottom Navigation:** Includes tabs for "Files", "Plots", "Packages", "Help", "Viewer", and "Presentation".

```

56 #sum_scores_avg_fe<- sum(Female_table$Midterm.Exam.Score)/nrow(Female_table)
57
58 Male_table<-Class1[Class1$Sex=="Male",]
59 Male_table
60 Male_MidtermScore1 = Male_table$Midterm.Exam.Score
61 Male_MidtermScore1=Male_MidtermScore1*100/53
62 summary(Male_MidtermScore1)
63
64 Male_ProblemSetScore1=Male_table$Problem.Sets.Score..percentage.
65 summary(Male_ProblemSetScore1)
66 |
67 #sum_scores_avg_ma<- sum(Male_table$Midterm.Exam.Score)/nrow(Male_table)
68
69 (Untitled) :

```

Console output (partial):

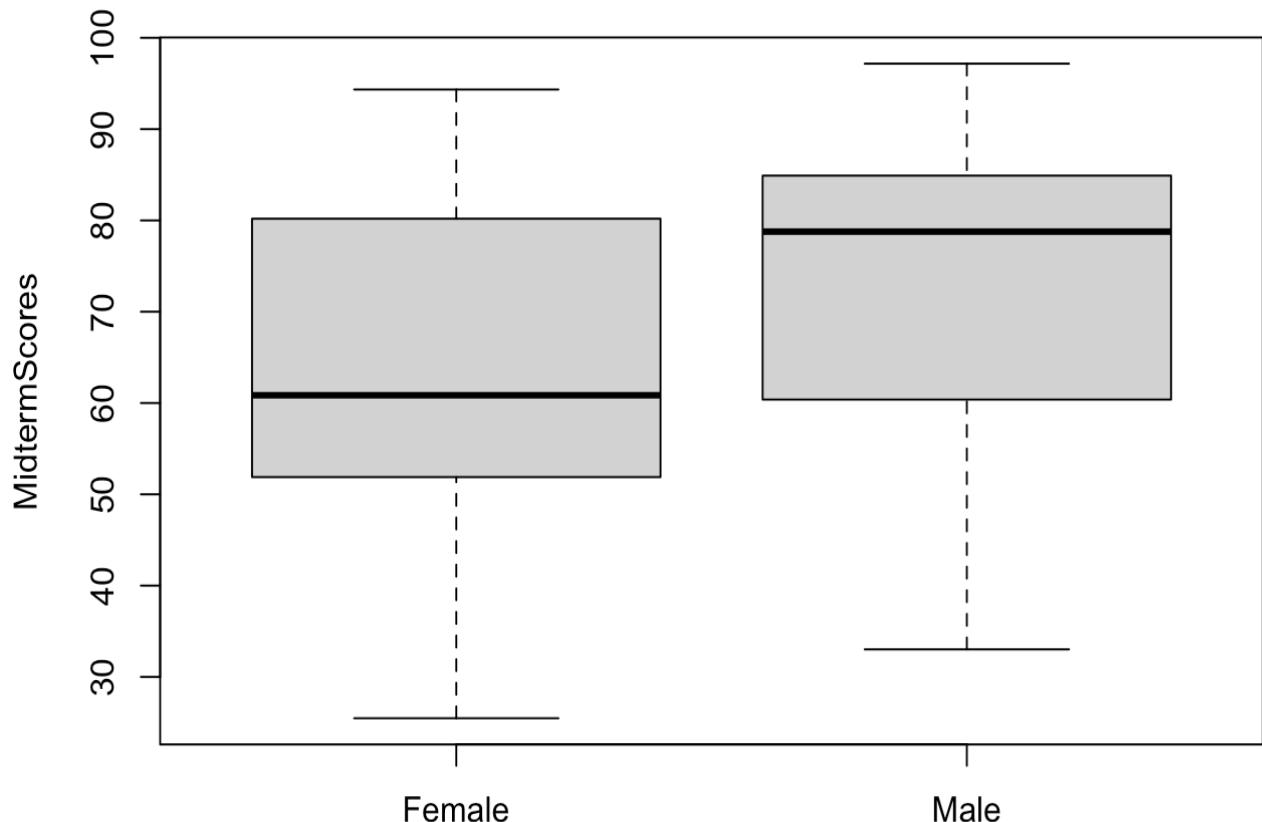
```

20      97.00 Male
21      98.00 Male
22      96.07 Male
23      96.73 Male
24      97.00 Male
25      99.13 Male
26      99.73 Male
27      84.27 Male
31      98.47 Male
32      99.07 Male
33      91.53 Male
35      98.67 Male
37      98.67 Male
38      97.73 Male
39      97.07 Male
41      96.33 Male
44      99.00 Male
45      99.53 Male
46      97.40 Male
48      98.33 Male
> Male_MidtermScore1 = Male_table$Midterm.Exam.Score
> Male_MidtermScore1=Male_MidtermScore1*100/53
> summary(Male_MidtermScore1)
   Min. 1st Qu. Median  Mean 3rd Qu.  Max.
 33.02   60.85   78.77   73.71   84.43   97.17
>
> Male_ProblemSetScore1=Male_table$Problem.Sets.Score..percentage.
> summary(Male_ProblemSetScore1)
   Min. 1st Qu. Median  Mean 3rd Qu.  Max.
 84.27   96.43   97.86   97.04   99.00   99.87
> |

```

Let's draw some boxplots for comparisons of midterm scores of female and male students of class1. Below is a box plots of midterm scores of both classes

```
boxplot(Female_MidtermScore1,Male_MidtermScore1,names=c("Female","Male"),ylab="MidtermScores")
```

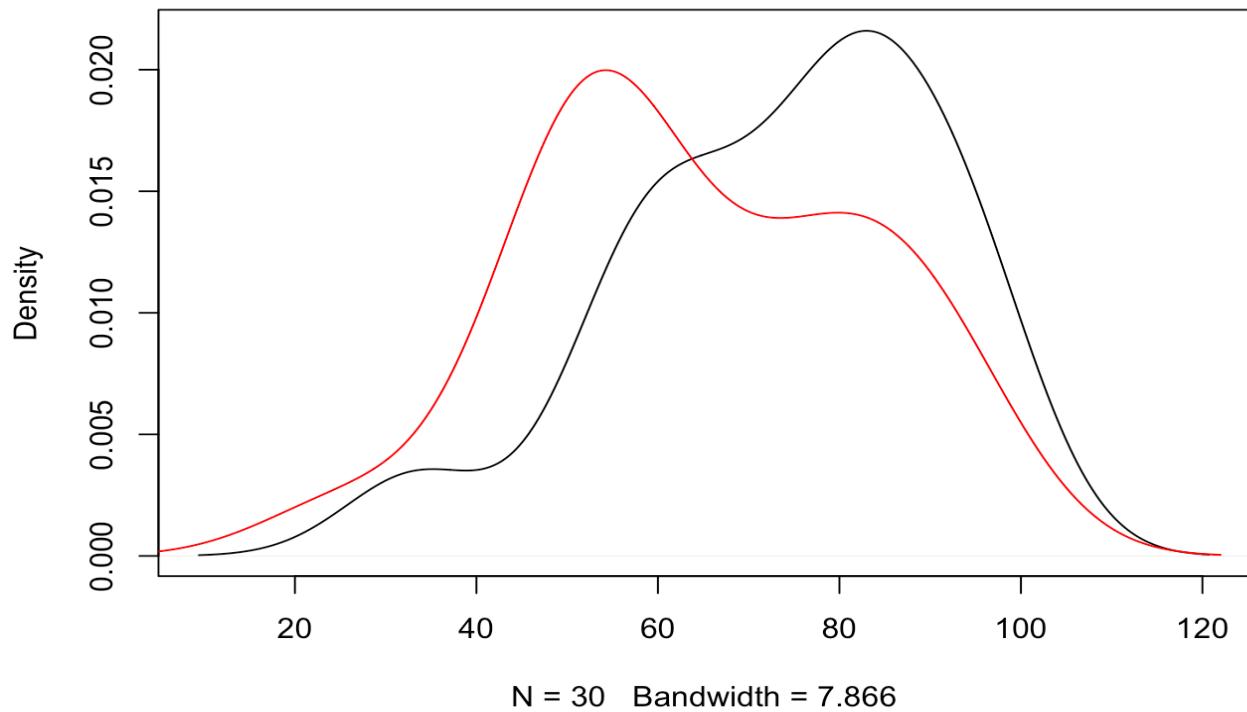


From the above boxplot, we can see that the median of the female midterm scores lie around 60 and the median of the midterm scores of male students lie around 78. 50% of the female students scores lie between 52 and 80 whereas male student scores lie between 60 and 84.

Let's draw density plots of midterm scores of male and female students.

```
plot(density(Male_MidtermScore1))
lines(density(Female_MidtermScore1), col="red")
```

density.default(x = Male_MidtermScore1)



We can see that the distribution of female midterm scores of class1 (red) is shifted to the left compared to male midterm scores of class1(black).

Hypothesis test:

Here we are considering midterm scores of both male and female students. so it's a two-sample location problem and the parameter of interest is ' Δ ' where $\Delta = \mu_2 - \mu_1$. Let μ_1 be population mean of the midterm scores of male students and μ_2 be the population mean of midterm scores of female students. So, we want to see if the population mean of the female students is equal to mean of male students. So, this becomes our alternative hypothesis i.e, $\mu_2 - \mu_1 = 0$ which says it is a two-tailed test. Null hypothesis is exactly opposite to alternative hypothesis. Assuming the significance level (α) is 0.05.

$$H_0: \Delta \neq 0$$

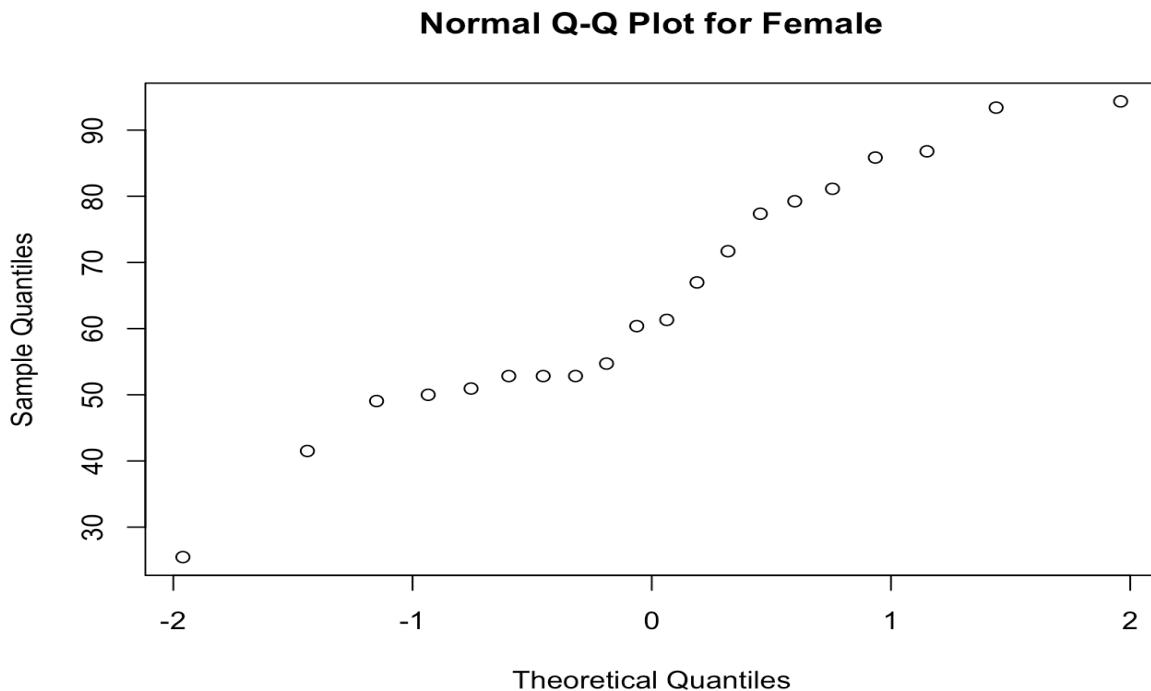
$$H_1: \Delta = 0$$

Now, let's draw QQ plots of midterm scores of both genders to see if the distribution is normal.

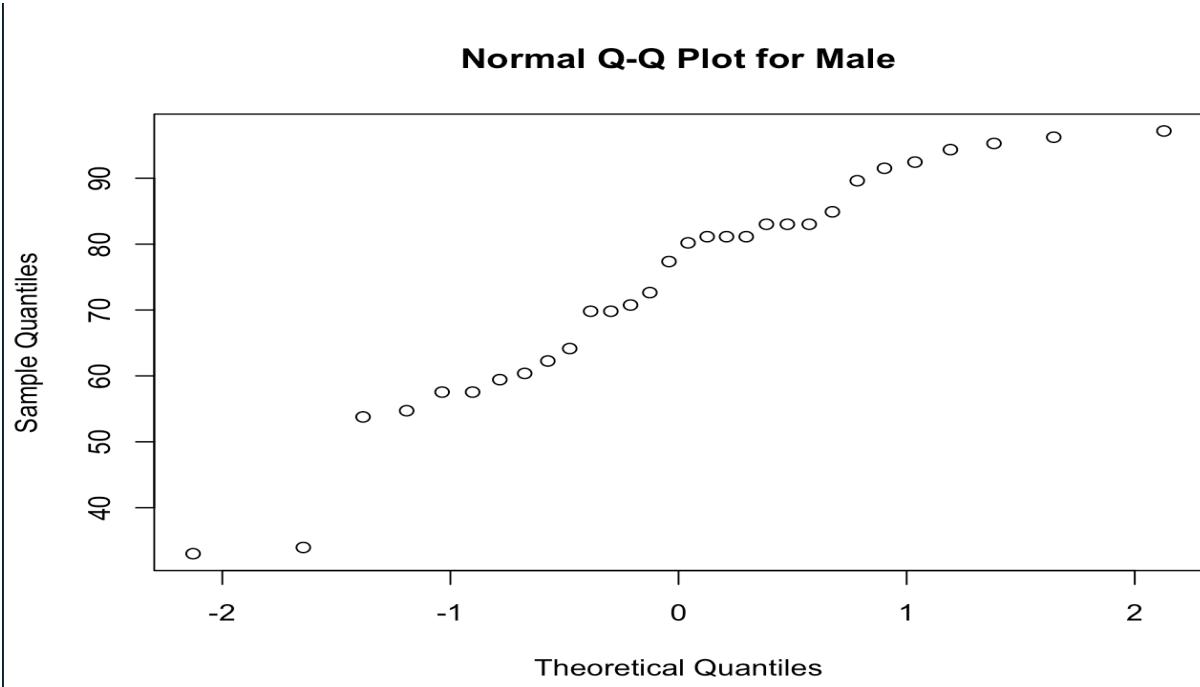
QQ plot of female and male student's midterm scores:

```
qqnorm(Female_MidtermScore1,main='Normal Q-Q Plot for Female')
qqnorm(Male_MidtermScore1,main='Normal Q-Q Plot for Male')
```

QQ plot of female student's midterm score:



QQ plot of male student's midterm score:



The QQ plot of male students Midterm score looks like a straight line except one outlier which is at 33. The QQ plot of female student scores also looks nearly straight except it is bent around 60.

Assuming that samples are IID samples from two independent normal populations, we can do Welch's two-sample t-test.

```
92
93 #Welch test
94
95 t.test(Male_MidtermScore1,Female_MidtermScore1)
96
97
98 # Problemset score male and female
99
```

95:1 # (Untitled) R Script

Console Terminal × Background Jobs × R 4.2.1 · ~/ ↗

```
Min. 1st Qu. Median Mean 3rd Qu. Max.
33.02 60.85 78.77 73.71 84.43 97.17
>
> Male_ProblemSetScore1=Male_table$Problem.Sets.Score..percentage.
> summary(Male_ProblemSetScore1)
  Min. 1st Qu. Median Mean 3rd Qu. Max.
84.27 96.43 97.86 97.04 99.00 99.87
> boxplot(Female_MidtermScore1, Male_MidtermScore1, names=c("Female", "Male"), ylab="MidtermScores")
> plot(density(Male_MidtermScore1))
> lines(density(Female_MidtermScore1), col="red")
> qqnorm(Female_MidtermScore1, main='Normal Q-Q Plot for Female')
> qqnorm(Male_MidtermScore1, main='Normal Q-Q Plot for Male')
> # using log
> log_Female_MidtermScore1=log(Female_MidtermScore1)
> log_Male_MidtermScore1=log(Male_MidtermScore1)
> qqnorm(log_Female_MidtermScore1, main='Normal Log Q-Q Plot for Female')
> qqnorm(log_Male_MidtermScore1, main='Normal Log Q-Q Plot for Male')
> t.test(Male_MidtermScore1, Female_MidtermScore1)

  Welch Two Sample t-test

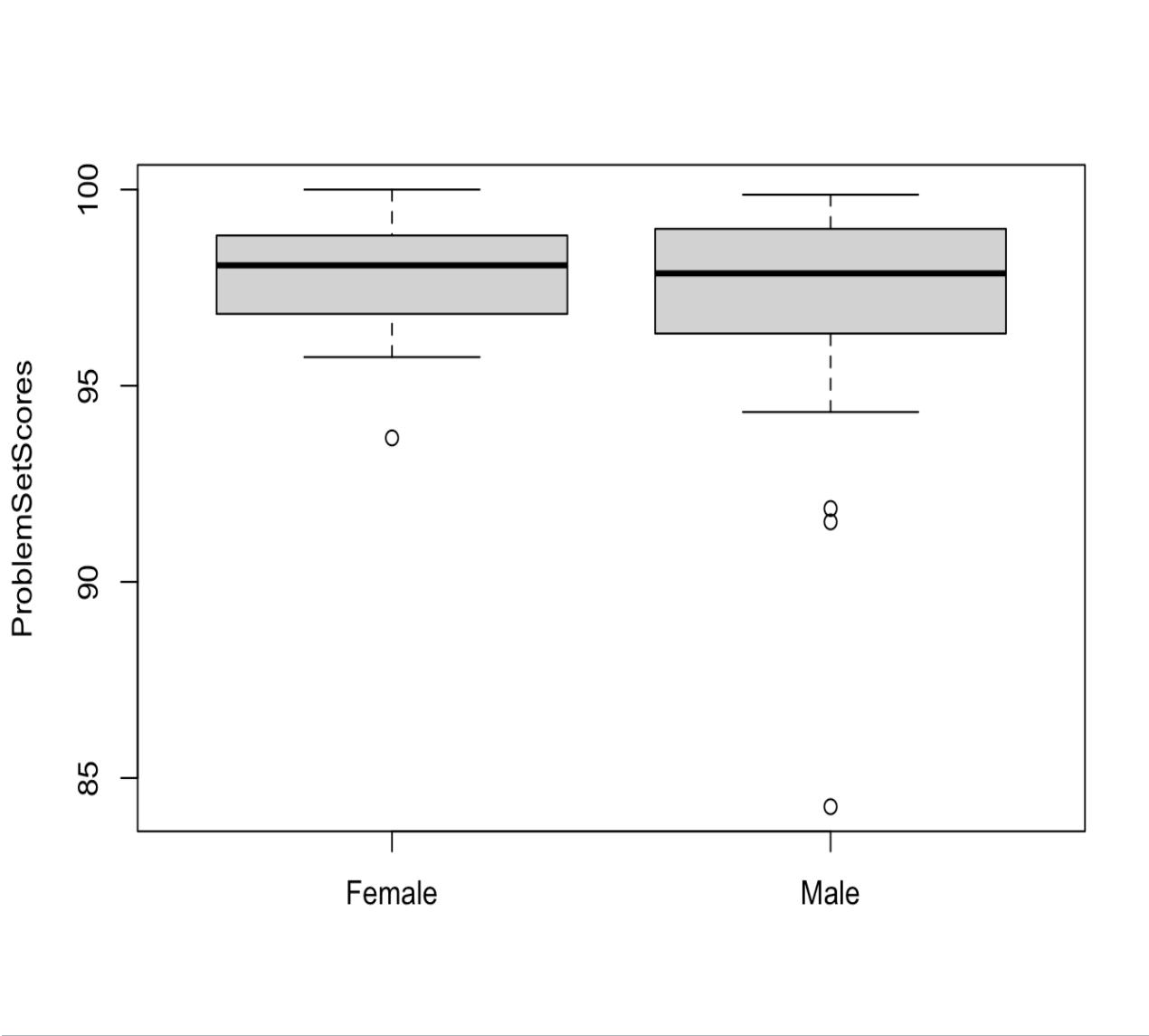
data: Male_MidtermScore1 and Female_MidtermScore1
t = 1.7747, df = 38.628, p-value = 0.08383
alternative hypothesis: true difference in means is not equal to 0
95 percent confidence interval:
-1.29962 19.85308
sample estimates:
mean of x mean of y
73.71069 64.43396
> |
```

From above t-test we can see that the p-value is 0.08383 which is greater than significance level (0.05) small, so it means we do not have any evidence against the null hypothesis. So, we cannot reject null hypothesis. So, our null hypothesis becomes true. That is gender does make difference in performance when taken midterm scores into account. From above, we can see that the 95% confidence interval is (-1.29962, 19.85308)

Now, Let's compare male and female students based on problem set scores.

Let's draw some boxplots for comparisons of problem set scores of female and male students of class1. Below is a box plots of problem set scores of both classes.

```
boxplot(Female_ProblemSetScore1,Male_ProblemSetScore1,names=c("Female","Male"),ylab="ProblemSetScores")
```

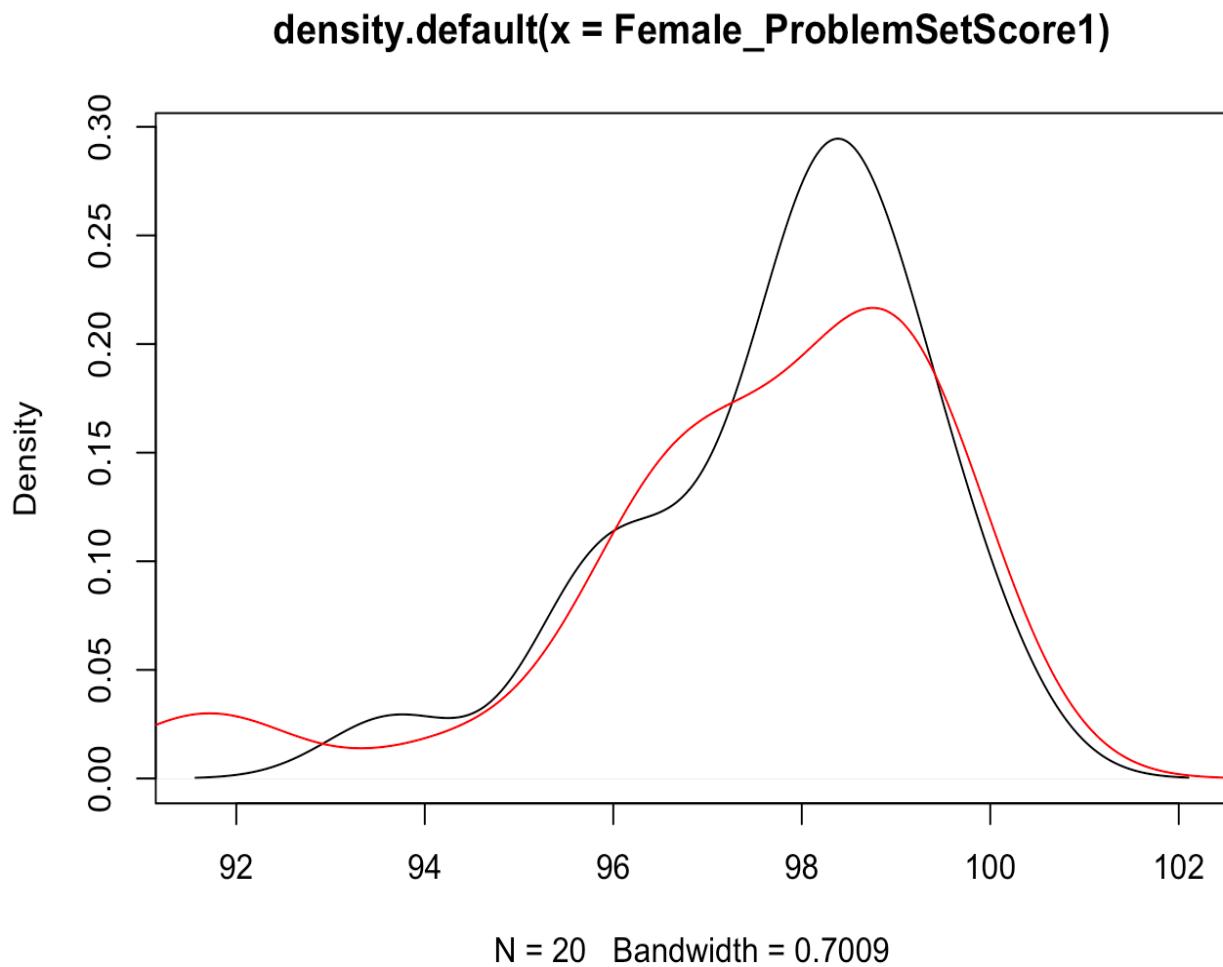


From above boxplot, we can see that the medians of problem set scores of both female and male students are nearly equal. We can also see that there is a outlier in female students scores which is around 93 .There are three outliers in male students scores in which two outliers lie between 90 and 92 and one outlier right below 85.

Let's draw density plot of problem set scores of male and female students.

R Code:

```
plot(density(Female_ProblemSetScore1))
lines(density(Male_ProblemSetScore1), col="red")
```



We can see that the distribution of female problem set scores of class1 (black) is nearly same to the male problem set scores of class1(red).

Hypothesis test:

Here we are considering problem set scores of both male and female students. so it's a two-sample location problem and the parameter of interest is ' Δ ' where $\Delta = \mu_2 - \mu_1$. Let μ_1 be population mean of the problem set scores of male students and μ_2 be the population mean of problem set scores of female students. So, we want to see if the population mean of the problem set scores of female students is equal to mean of problem set scores of male

students. So, this becomes our alternative hypothesis i.e, $\mu_2 - \mu_1 = 0$ which says it is a two-tailed test. Null hypothesis is exactly opposite to alternative hypothesis.

$$H_0: \Delta \neq 0$$

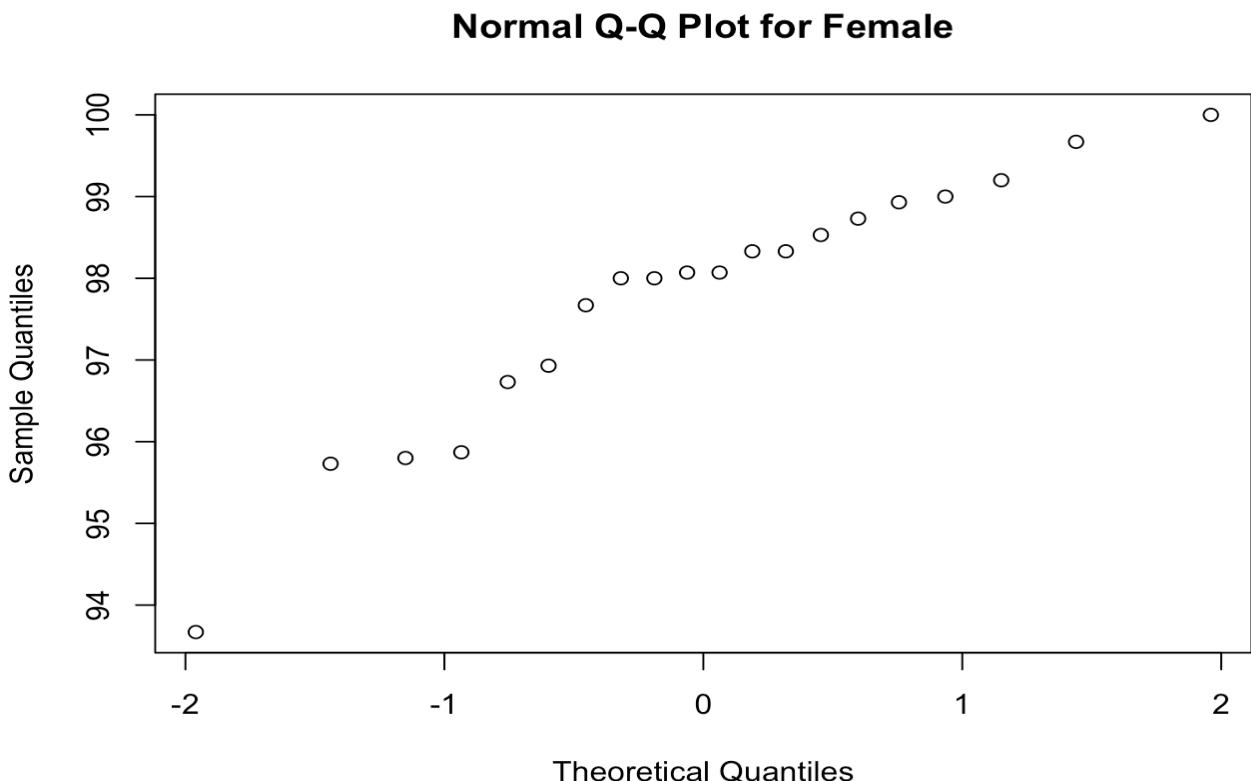
$$H_1: \Delta = 0$$

Now, let's draw QQ plots of problem set scores of both genders to see if the distribution is normal.

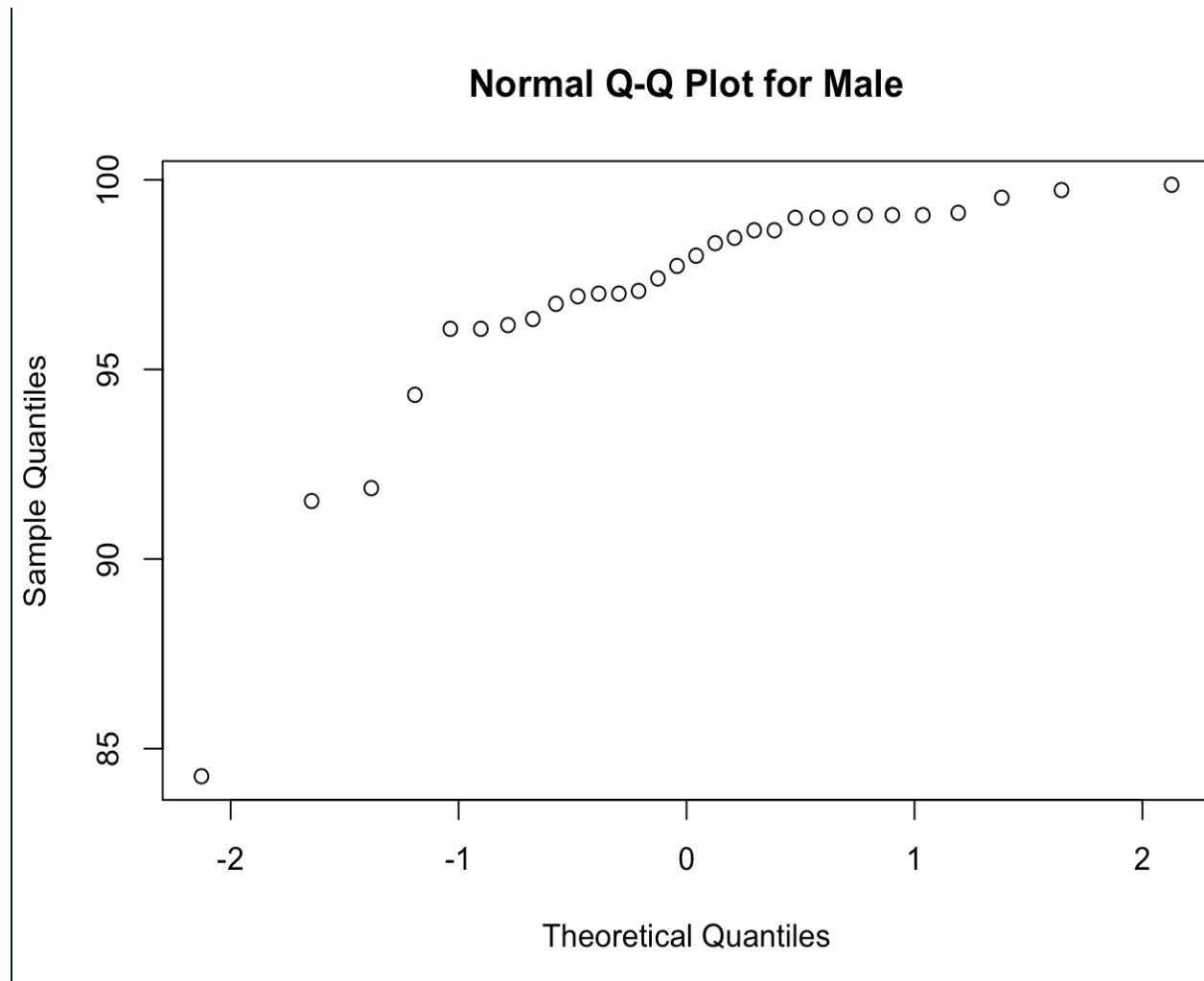
QQ plot of female and male student's problem set scores:

```
qqnorm(Female_ProblemSetScore1,main='Normal Q-Q Plot for Female')
qqnorm(Male_ProblemSetScore1,main='Normal Q-Q Plot for Male')
```

QQ plot of female student's problem set score:



QQ plot of male student's problem set score:



The QQ plot of female student's problem set score looks like a straight line but the QQ plot of male student's problem set score does not look like a straight line it looks like a curved line with outlier right below 85.

Assuming that samples are IID samples from two independent normal populations, we can do Welch's two- sample t-test.

The screenshot shows the RStudio interface with the following details:

- Top Panel (Code Editor):** Displays R code for statistical analysis. Lines 113-125 include plotting a Q-Q plot for males, calculating a delta value, performing a Welch t-test, and starting a regression section. Line 120 shows the t-test command: `t.test(Male_ProblemSetScore1, Female_ProblemSetScore1)`.
- Bottom Panel (Console):** Shows the R console output. It starts with sample estimates for means of x and y. Then it runs a series of commands to create boxplots, density plots, and Q-Q plots for both female and male problem set scores. It then performs log-transformed versions of these plots. Finally, it runs a Welch Two Sample t-test comparing Male_ProblemSetScore1 and Female_ProblemSetScore1.
- Output of the t-test:**

```

Welch Two Sample t-test

data: Male_ProblemSetScore1 and Female_ProblemSetScore1
t = -1.0752, df = 44.857, p-value = 0.2881
alternative hypothesis: true difference in means is not equal to 0
95 percent confidence interval:
-2.0861425 0.6341425
sample estimates:
mean of x mean of y
97.037   97.763

```

From above t-test we can see that the p-value is 0.2881 which is NOT small, so it means we do not have any evidence against the null hypothesis. So, we cannot reject null hypothesis. So, our null hypothesis becomes true. That is gender does make difference in performance when taken problem set scores into account. From above, we can see that the 95% confidence interval is (-2.0861425,0.6341425)

To conclude, gender does make difference in students' performance in assignments and midterm exam in S520.

3. Research question: Should we use students' assignment grades to predict their midterm exam scores? Is the conclusion the same for different programs (i.e., online vs residential)?

Topic: Regression and Prediction

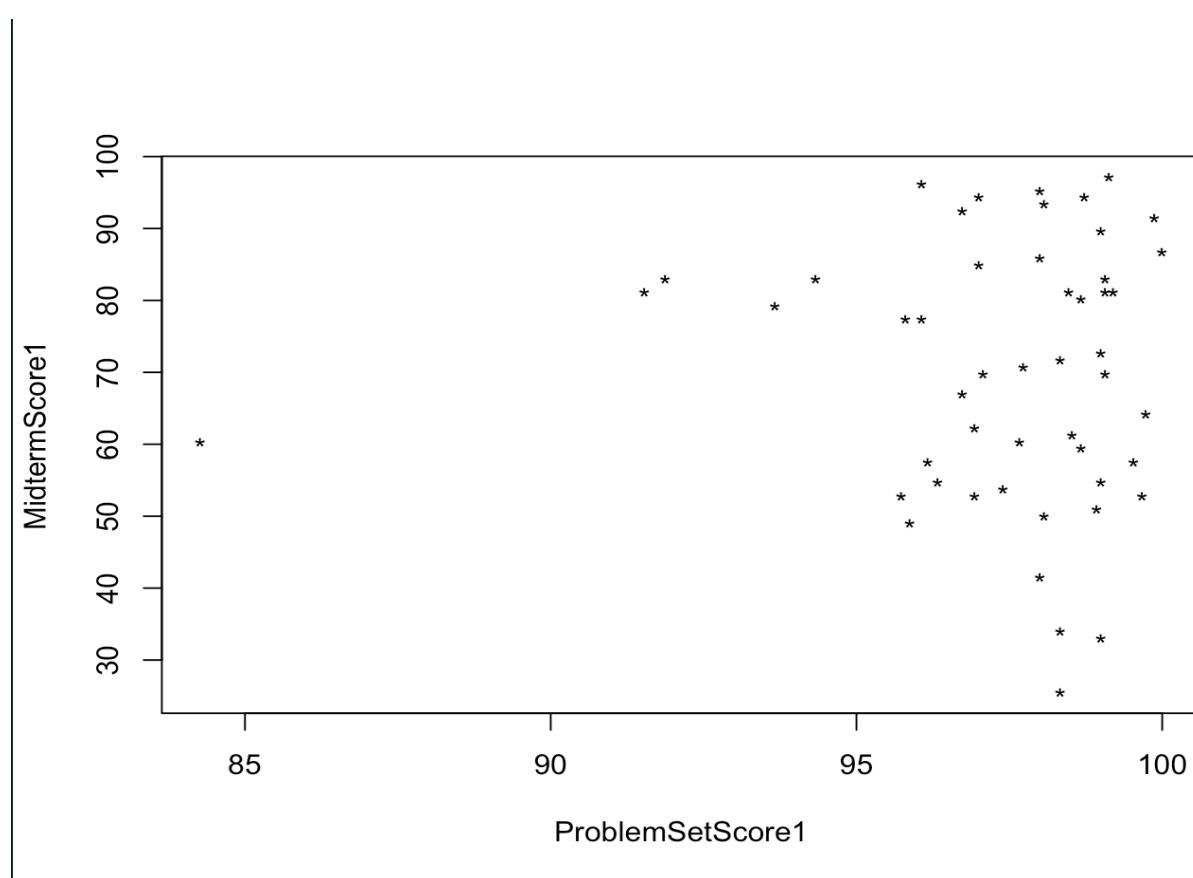
Purpose of the Topic: To check whether we can be able to predict student midterm exam score by using student's assignment grades for both online and residential.

Firstly, we have checked for Class 1.

For Class 1:

R Code:

```
plot(ProblemSetScore1, MidtermScore1, pch="*")
```



From the plot we can be able to observe that the data is less, and we cannot be able to say whether it is correlation high or low.

Next, we have found correlation between problem set score1 and mid-termscore1.

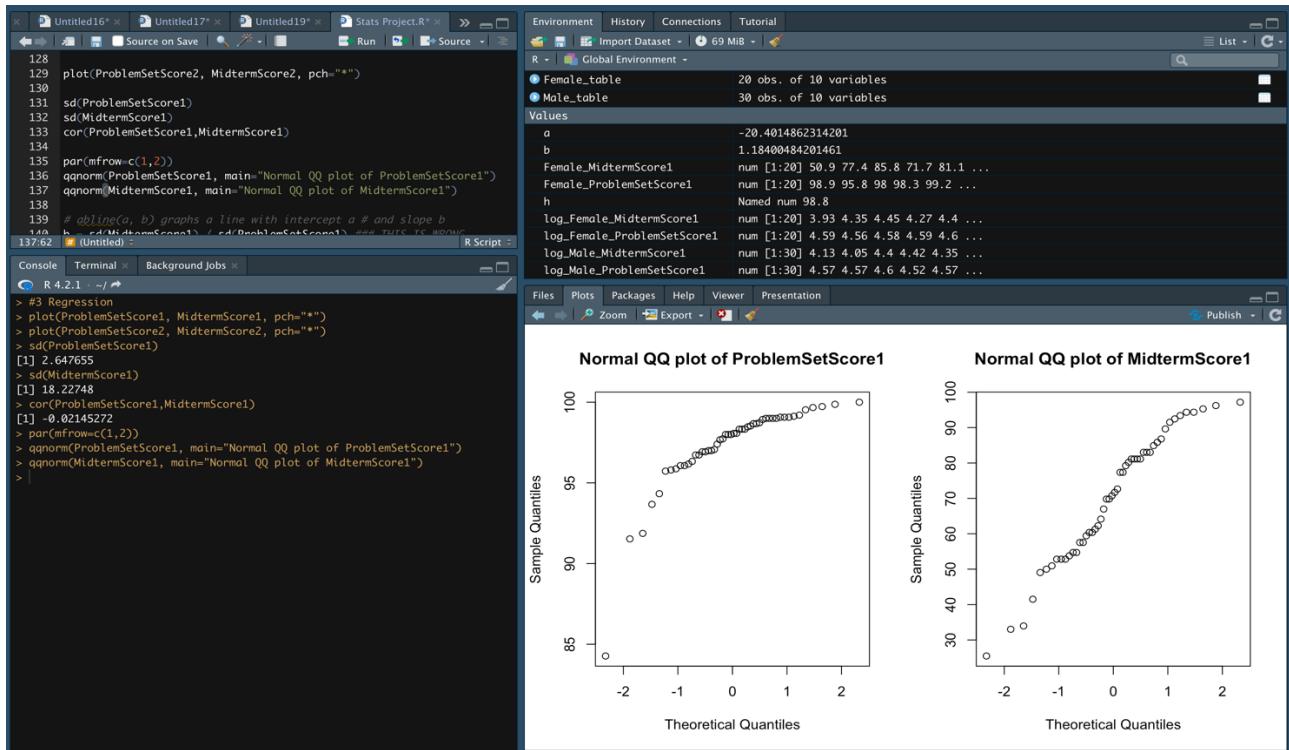
```

124
125 sd(ProblemSetScore1)
126 sd(MidtermScore1)
127 cor(ProblemSetScore1, MidtermScore1)
128
128:1 # (Untitled) R Script
Console Terminal Background Jobs
R 4.2.1 · ~/ ↗
> sd(ProblemSetScore1)
[1] 2.647655
> sd(MidtermScore1)
[1] 18.22748
> cor(ProblemSetScore1, MidtermScore1)
[1] -0.02145272
>

```

As we can see there is -0.02145272 correlation between them, and we can say that if one value increases then other value decreases as it is a negative correlation.

For checking that the data Is bivariate normal, we have drawn QQ plots:



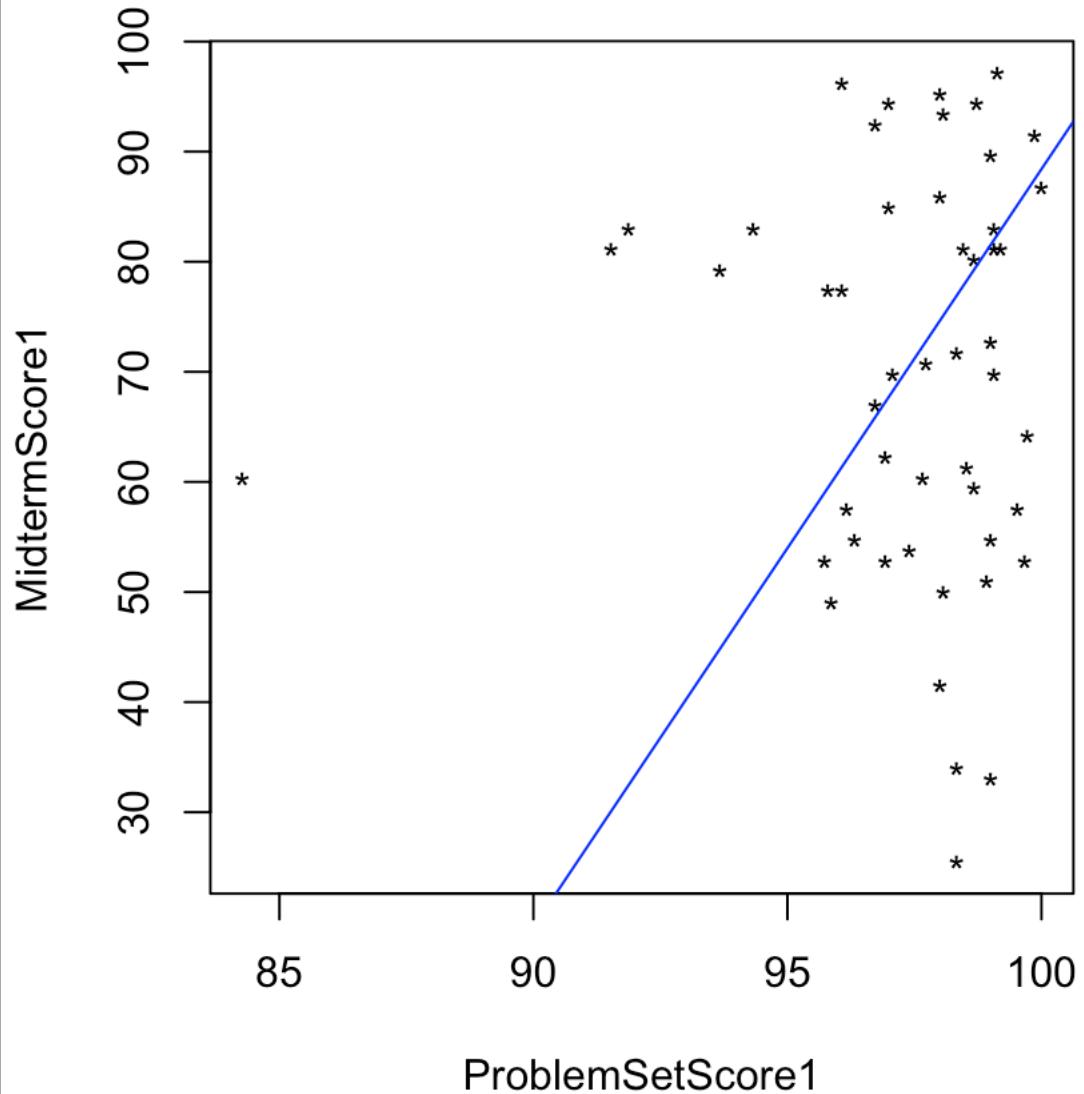
From above plots we can see that the QQ plot for class1 midterm scores looks almost like a straight line and the QQ plot for class1 problem set scores is more like curved line than a straight line. I cannot conclude that the data is very close to bivariate normal.

To find the prediction first we need to find slope so for that first we have done

SD-Line:

R Code:

```
# abline(a, b) graphs a line with intercept a # and slope b
b = sd(MidtermScore1) / sd(ProblemSetScore1) ### THIS IS WRONG
a = mean(MidtermScore1) - b * mean(ProblemSetScore1)
plot(ProblemSetScore1, MidtermScore1, pch="*")
abline(a, b, col="blue")
```

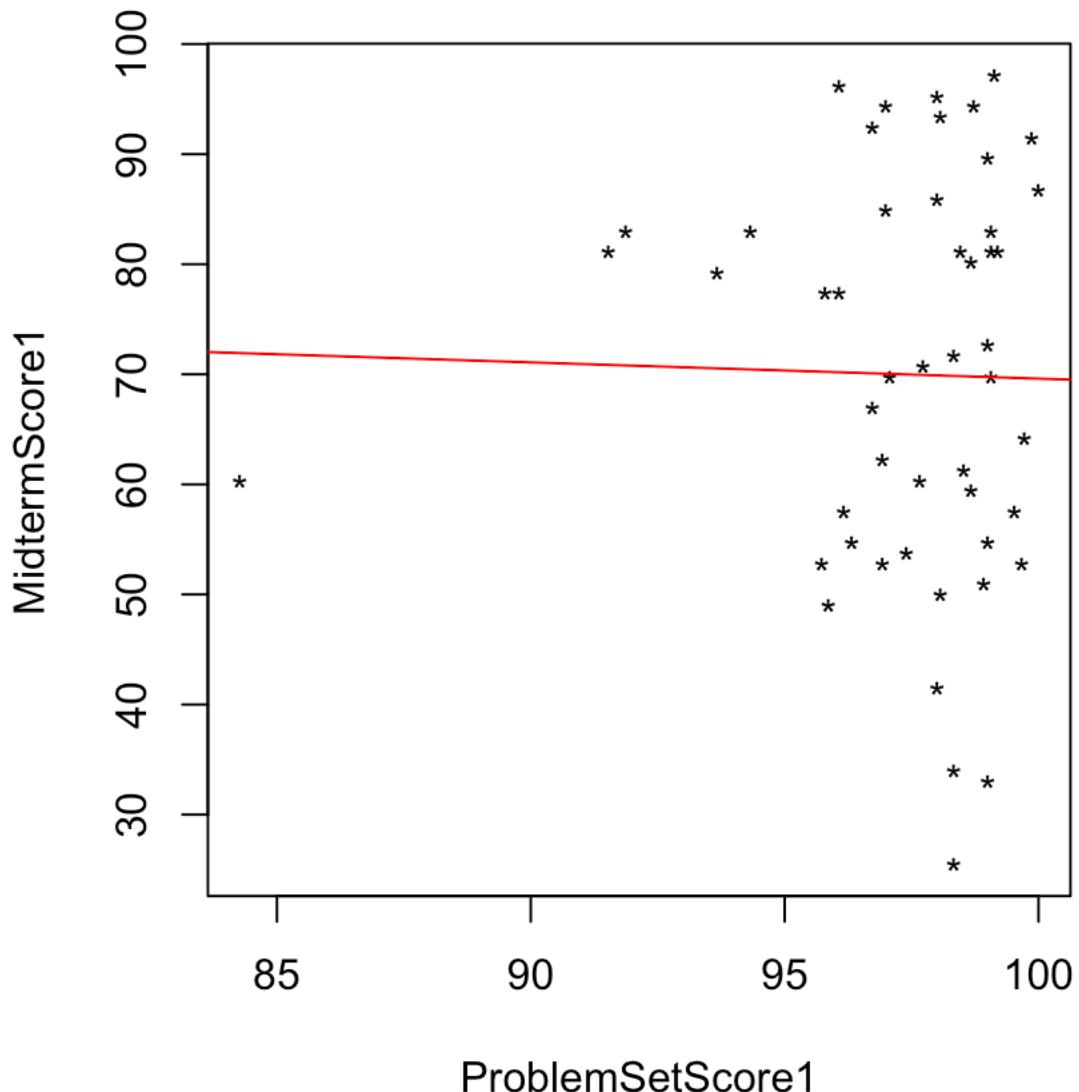


The blue line is called the SD line. It describes the data well but is terrible at prediction. The correct answer is not at all obvious until you do some calculus or linear algebra.

Correlation * SD of weight / SD of height:

R Code:

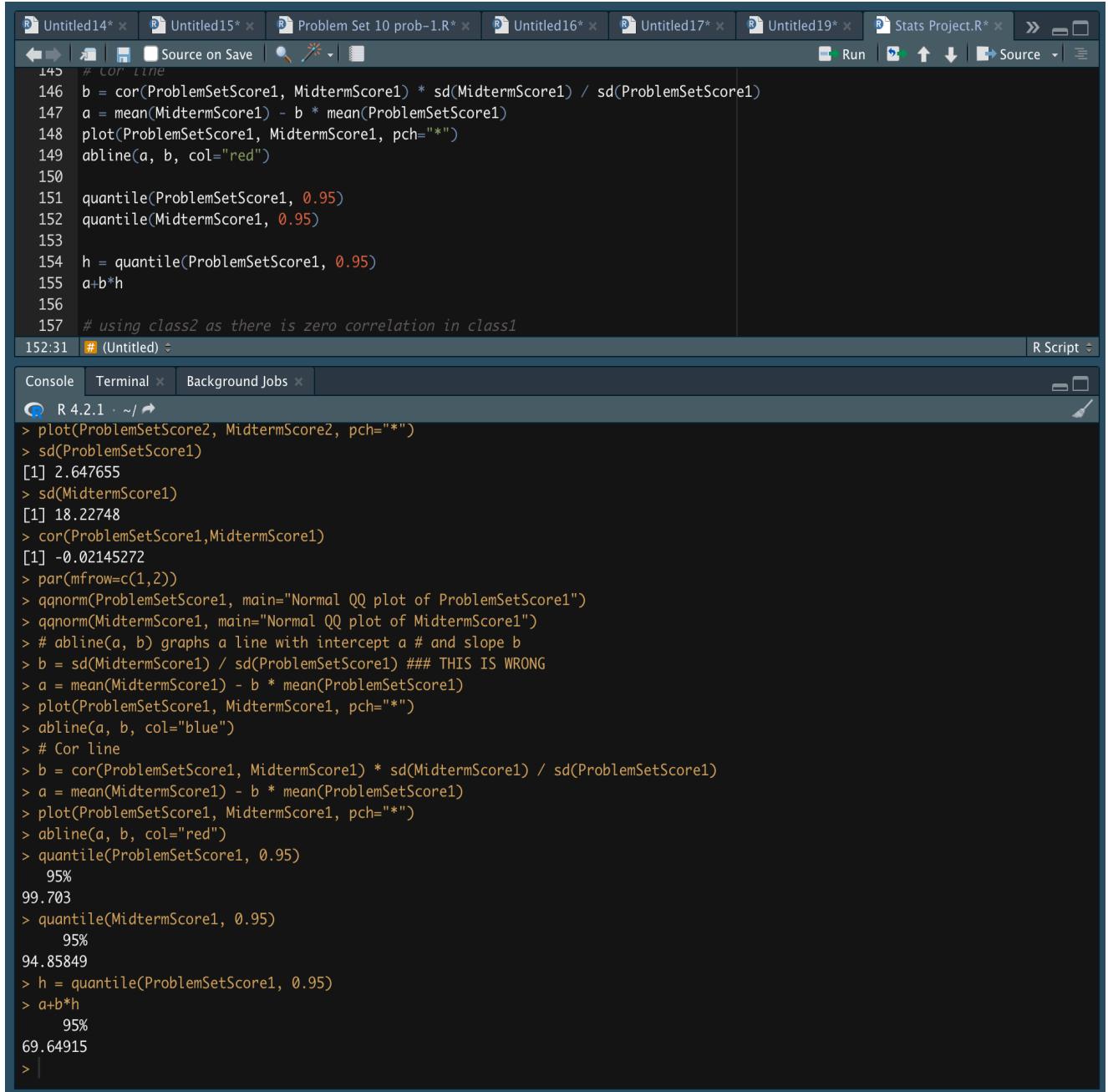
```
# Cor line
b = cor(ProblemSetScore1, MidtermScore1) * sd(MidtermScore1) / sd(ProblemSetScore1)
a = mean(MidtermScore1) - b * mean(ProblemSetScore1)
plot(ProblemSetScore1, MidtermScore1, pch="*")
abline(a, b, col="red")
```



From finding slope by using correlation is the best way to predict rather than SD Line. It seems that's there is no such linear relationship between problem set score 1 and midterm score1.

Next, we have found Quantile:

R Code:



The screenshot shows the RStudio interface with the following details:

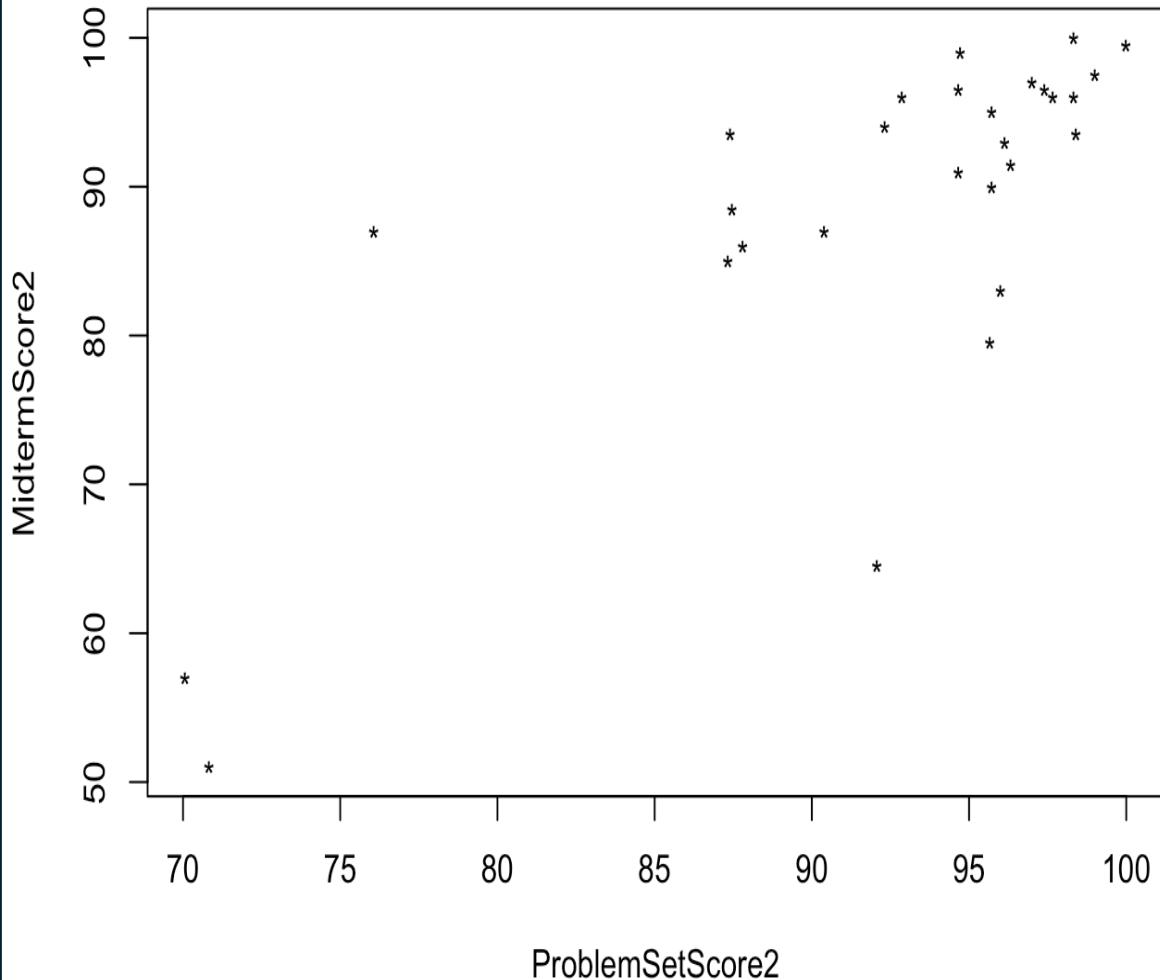
- Code Editor:** Shows R code for calculating a regression line and quantiles. The code includes comments like "# Cor line" and "# using class2 as there is zero correlation in class1".
- Console:** Shows the R session output. It starts with a warning about a non-existent file "R 4.2.1 · ~/". Then it performs calculations:
 - sd(ProblemSetScore1) = 2.647655
 - sd(MidtermScore1) = 18.22748
 - cor(ProblemSetScore1, MidtermScore1) = -0.02145272
 - par(mfrow=c(1,2))
 - qqnorm(ProblemSetScore1, main="Normal QQ plot of ProblemSetScore1")
 - qqnorm(MidtermScore1, main="Normal QQ plot of MidtermScore1")
 - # abline(a, b) graphs a line with intercept a # and slope b
 - b = sd(MidtermScore1) / sd(ProblemSetScore1) ## THIS IS WRONG
 - a = mean(MidtermScore1) - b * mean(ProblemSetScore1)
 - plot(ProblemSetScore1, MidtermScore1, pch="*")
 - abline(a, b, col="blue")
 - # Cor line
 - b = cor(ProblemSetScore1, MidtermScore1) * sd(MidtermScore1) / sd(ProblemSetScore1)
 - a = mean(MidtermScore1) - b * mean(ProblemSetScore1)
 - plot(ProblemSetScore1, MidtermScore1, pch="*")
 - abline(a, b, col="red")
 - quantile(ProblemSetScore1, 0.95)
95%
99.703
 - quantile(MidtermScore1, 0.95)
95%
94.85849
 - h = quantile(ProblemSetScore1, 0.95)
 - a+b*h
95%
69.64915

This is much less than the 95th percentile of MidtermScore1! So, we're saying that we do not predict midterm scores from problem set grades for Class1.

For Class 2:

R Code:

```
plot(ProblemSetScore2, MidtermScore2, pch="*")
```



From the plot we can be able to observe that the data is less, and we cannot be able to say whether it is correlation high or low.

Next, we have found correlation between problem set score2 and mid-termscore2.

The screenshot shows the RStudio interface. The top panel displays an R script with the following code:

```
150 a+b*h
151
152 # using class2 as there is zero correlation in class1
153 sd(ProblemSetScore2)
154 sd(MidtermScore2)
155 cor(ProblemSetScore2,MidtermScore2)
156
157 par(mfrow=c(1,2))
158 qqnorm(ProblemSetScore2, main="Normal QQ plot of ProblemSetScore2")
159 qqnorm(MidtermScore2, main="Normal QQ plot of MidtermScore2")
160
161 # abline(a, b) graphs a line with intercept a # and slope b
162 b = sd(MidtermScore2) / sd(ProblemSetScore2) ### THIS IS WRONG
163 a = mean(MidtermScore2) - b * mean(ProblemSetScore2)
```

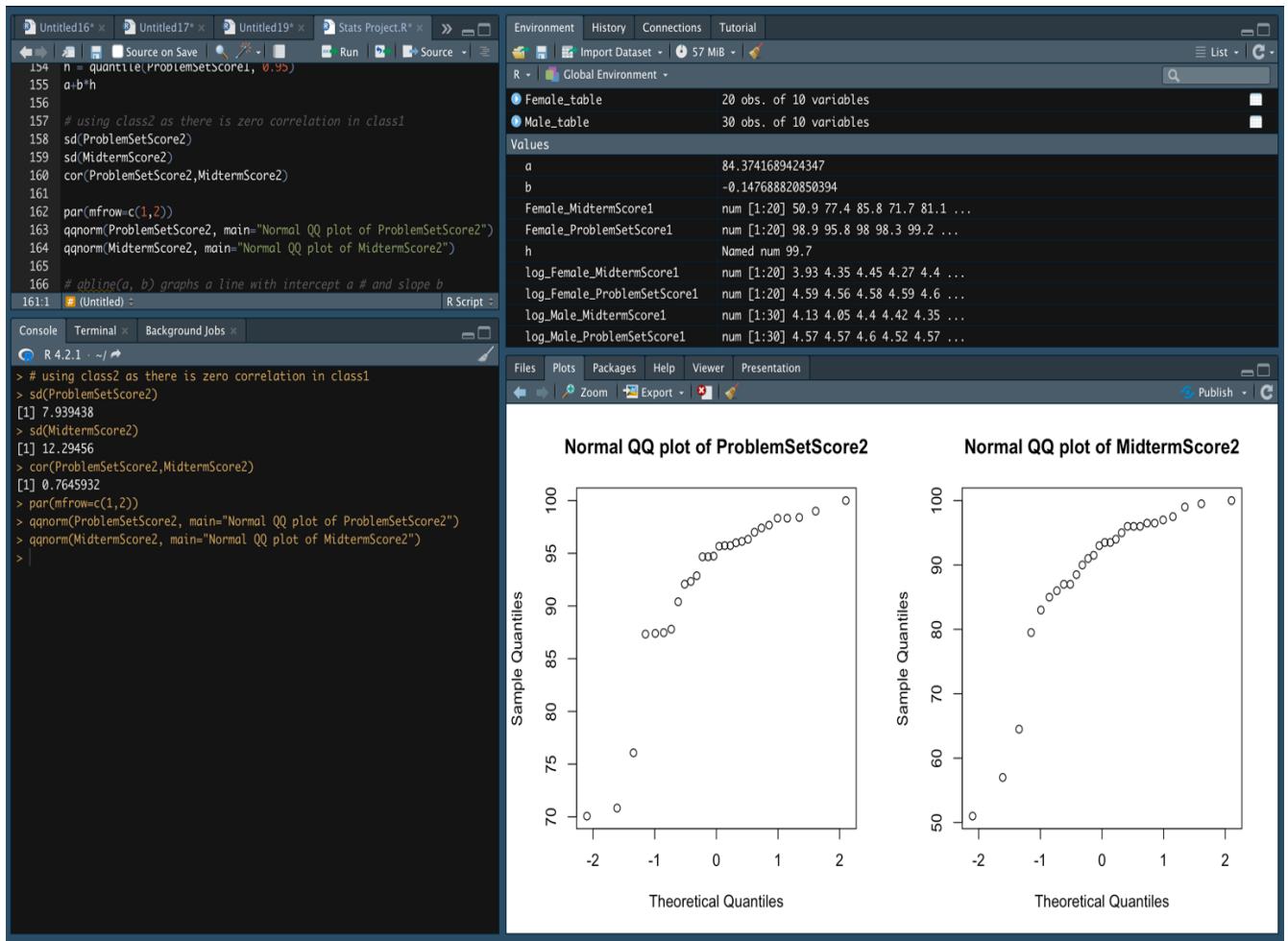
The status bar at the bottom left indicates "155:36 # (Untitled)". The bottom panel shows the R Console output:

```
R 4.2.1 · ~/r
> sd(ProblemSetScore2)
[1] 7.939438
> sd(MidtermScore2)
[1] 12.29456
> cor(ProblemSetScore2,MidtermScore2)
[1] 0.7645932
>
```

As we can see that there is 0.7645932 correlation between Problem set score 2 and midterm score 2, and we can say that it has good correlation between them.

For checking that the data Is bivariate normal, we have drawn QQ plots:

```
par(mfrow=c(1,2))
qqnorm(ProblemSetScore2, main="Normal QQ plot of ProblemSetScore2")
qqnorm(MidtermScore2, main="Normal QQ plot of MidtermScore2")
```



From above plots we can see that the QQ plots for class2 midterm scores and class2 problem set scores are more like curved line than a straight line. We cannot conclude that the data is very close to bivariate normal.

To find the prediction first we need to find slope so for that first we have done

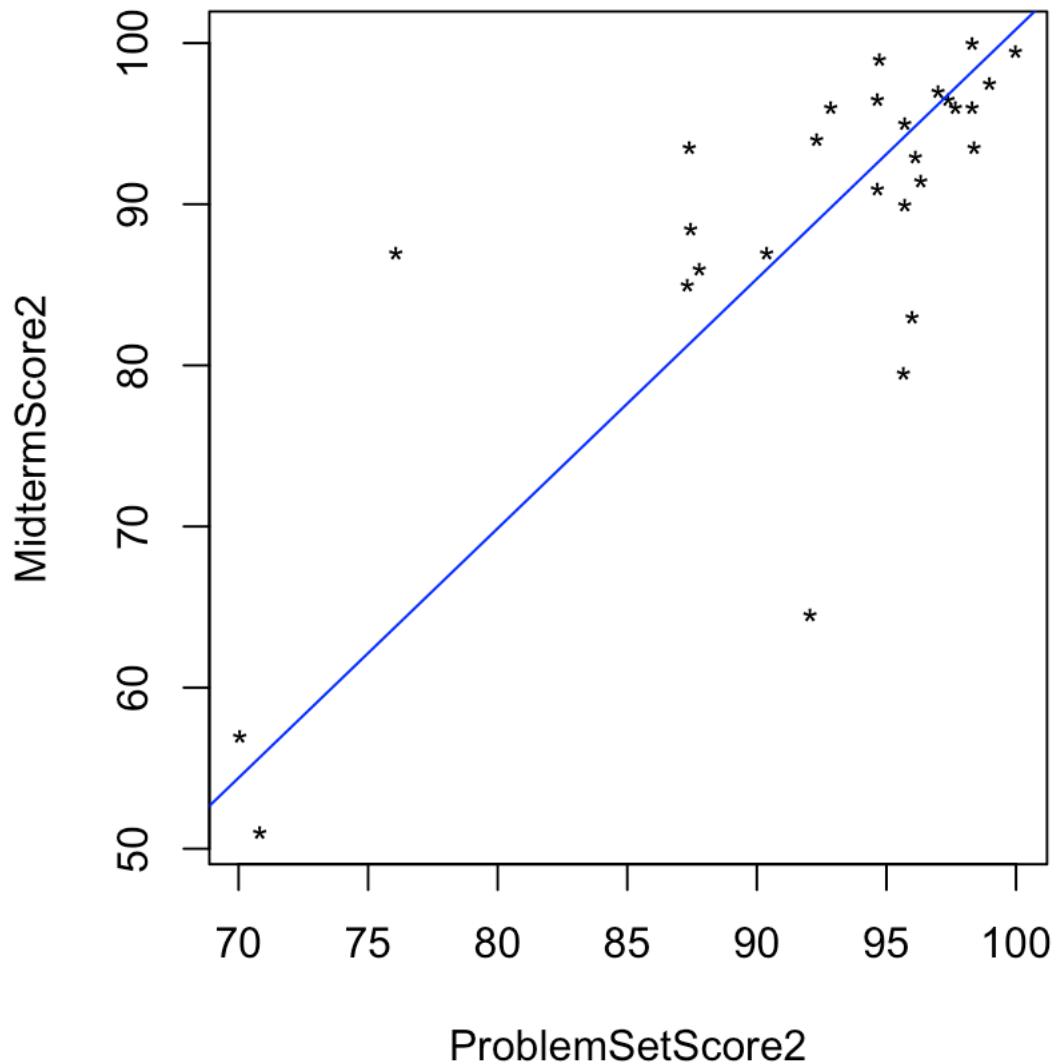
SD-Line:

R Code:

```

# abline(a, b) graphs a line with intercept a # and slope b
b = sd(MidtermScore2) / sd(ProblemSetScore2) ### THIS IS WRONG
a = mean(MidtermScore2) - b * mean(ProblemSetScore2)
plot(ProblemSetScore2, MidtermScore2, pch="*")
abline(a, b, col="blue")

```

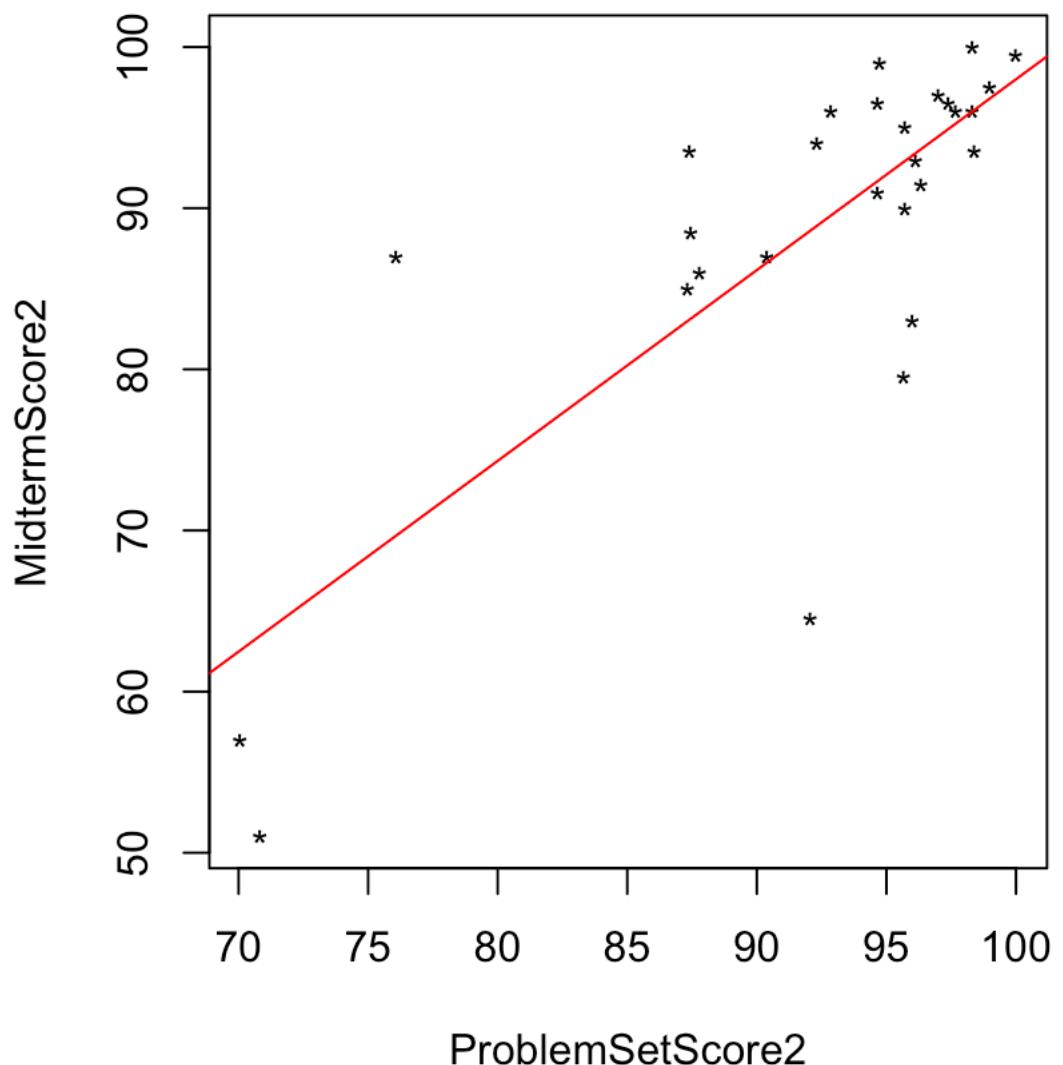


The blue line is called the SD line. It describes the data well but is terrible at prediction. The correct answer is not at all obvious until you do some calculus or linear algebra.

Correlation * SD of weight / SD of height:

R Code:

```
# Cor line
b = cor(ProblemSetScore2, MidtermScore2) * sd(MidtermScore2) / sd(ProblemSetScore2)
a = mean(MidtermScore2) - b * mean(ProblemSetScore2)
plot(ProblemSetScore2, MidtermScore2, pch="*")
abline(a, b, col="red")
```



From finding slope by using correlation is the best way to predict rather than SD Line.

Next, we have found Quantile:

R Code:

The screenshot shows the RStudio interface. The top bar has tabs for "Problem Set 10 prob-1.R*", "Untitled16*", "Untitled17*", "Untitled19*", "Stats Project.R*", and "Run". Below the tabs are icons for back, forward, save, and search. The main area contains R code:

```

172 # Cor line
173 b = cor(ProblemSetScore2, MidtermScore2) * sd(MidtermScore2) / sd(ProblemSetScore2)
174 a = mean(MidtermScore2) - b * mean(ProblemSetScore2)
175 plot(ProblemSetScore2, MidtermScore2, pch="*")
176 abline(a, b, col="red")
177
178 quantile(ProblemSetScore2, 0.95)
179 quantile(MidtermScore2, 0.95)
180
181 h = quantile(ProblemSetScore2, 0.95)
182 a+b*h
183 # This value is lesser than the 95percentile of the midtermscore2.
184

```

The status bar at the bottom left says "183:67 # (Untitled)". The bottom panel shows the R console with the following history:

```

Console Terminal × Background Jobs ×
R 4.2.1 · ~/ ↻
> par(mfrow=c(1,2))
> qqnorm(ProblemSetScore2, main="Normal QQ plot of ProblemSetScore2")
> qqnorm(MidtermScore2, main="Normal QQ plot of MidtermScore2")
> # abline(a, b) graphs a line with intercept a # and slope b
> b = sd(MidtermScore2) / sd(ProblemSetScore2) ### THIS IS WRONG
> a = mean(MidtermScore2) - b * mean(ProblemSetScore2)
> plot(ProblemSetScore2, MidtermScore2, pch="*")
> abline(a, b, col="blue")
> # Cor line
> b = cor(ProblemSetScore2, MidtermScore2) * sd(MidtermScore2) / sd(ProblemSetScore2)
> a = mean(MidtermScore2) - b * mean(ProblemSetScore2)
> plot(ProblemSetScore2, MidtermScore2, pch=".")
> abline(a, b, col="red")
> # Cor line
> b = cor(ProblemSetScore2, MidtermScore2) * sd(MidtermScore2) / sd(ProblemSetScore2)
> a = mean(MidtermScore2) - b * mean(ProblemSetScore2)
> plot(ProblemSetScore2, MidtermScore2, pch="*")
> abline(a, b, col="red")
> quantile(ProblemSetScore2, 0.95)
  95%
98.79
> quantile(MidtermScore2, 0.95)
  95%
99.325
>
> h = quantile(ProblemSetScore2, 0.95)
> a+b*h
  95%
96.56635
> # This value is lesser than the 95percentile of the midtermscore2.
>

```

We can predict the midterm scores by using problem set scores for Class 2 but I would help more if we have more data.

Total R-Code for 1 , 2 , 3:

```
Class1<- read.csv('/Users/saisrikar/S520 Project Data.csv')
```

```
Class1
```

```
Class2<-read.csv('/Users/saisrikar/S520 Project Data_1.csv')
```

```
Class2
```

```
summary(Class1)
```

```
summary(Class2)
```

```
MidtermScore1=Class1$Midterm.Exam.Score
```

```
MidtermScore1=MidtermScore1*100/53
```

```
MidtermScore1
```

```
MidtermScore2=Class2$Midterm.Exam.Score..percentage.
```

```
MidtermScore2
```

```
summary(MidtermScore1)
```

```
summary(MidtermScore2)
```

```
boxplot(MidtermScore1, MidtermScore2, names=c("Class_1", "Class_2"), ylab="Midterm Scores")
```

```
plot(density(MidtermScore2))
```

```
lines(density(MidtermScore1), col="red")
```

```
qqnorm(MidtermScore1, main='Normal Q-Q Plot for Class_1')
```

```
qqnorm(MidtermScore2, main='Normal Q-Q Plot for Class_2')
```

```
# using log
```

```
log_MidtermScore1=log(MidtermScore1)
```

```
log_MidtermScore2=log(MidtermScore2)
```

```

qqnorm(log_MidtermScore1,main='Normal Log Q-Q Plot for Class_1')
qqnorm(log_MidtermScore2,main='Normal Log Q-Q Plot for Class_2')
#Welch test
t.test(MidtermScore2,MidtermScore1,alternative ="greater")
#####
#####

#2 Inferential

Female_table=Class1[Class1$Sex=="Female",]

Female_table

Female_MidtermScore1 = Female_table$Midterm.Exam.Score

Female_MidtermScore1=Female_MidtermScore1*100/53

summary(Female_MidtermScore1)

Female_ProblemSetScore1=Female_table$Problem.Sets.Score..percentage.

summary(Female_ProblemSetScore1)

Male_table=Class1[Class1$Sex=="Male",]

Male_table

Male_MidtermScore1 = Male_table$Midterm.Exam.Score

Male_MidtermScore1=Male_MidtermScore1*100/53

summary(Male_MidtermScore1)

Male_ProblemSetScore1=Male_table$Problem.Sets.Score..percentage.

summary(Male_ProblemSetScore1)

boxplot(Female_MidtermScore1, Male_MidtermScore1, names=c("Female","Male"), ylab="MidtermScores")

plot(density(Male_MidtermScore1))
lines(density(Female_MidtermScore1), col="red")

```

```
qqnorm(Female_MidtermScore1,main='Normal Q-Q Plot for Female')

qqnorm(Male_MidtermScore1,main='Normal Q-Q Plot for Male')

# using log

log_Female_MidtermScore1=log(Female_MidtermScore1)

log_Male_MidtermScore1=log(Male_MidtermScore1)

qqnorm(log_Female_MidtermScore1,main='Normal Log Q-Q Plot for Female')

qqnorm(log_Male_MidtermScore1,main='Normal Log Q-Q Plot for Male')

#Welch test

t.test(Male_MidtermScore1,Female_MidtermScore1)

# Problemset score male and female

boxplot(Female_ProblemSetScore1,Male_ProblemSetScore1,names=c("Female","Male"),
,ylab="ProblemSetScores")

plot(density(Female_ProblemSetScore1))

lines(density(Male_ProblemSetScore1), col="red")

qqnorm(Female_ProblemSetScore1,main='Normal Q-Q Plot for Female')

qqnorm(Male_ProblemSetScore1,main='Normal Q-Q Plot for Male')

# using log

log_Female_ProblemSetScore1=log(Female_ProblemSetScore1)

log_Male_ProblemSetScore1=log(Male_ProblemSetScore1)

qqnorm(log_Female_ProblemSetScore1,main='Normal Log Q-Q Plot for Female')

qqnorm(log_Male_ProblemSetScore1,main='Normal Log Q-Q Plot for Male')

#Welch test

t.test(Male_ProblemSetScore1,Female_ProblemSetScore1)

#####
#####
```

#3 Regression

```
plot(ProblemSetScore1, MidtermScore1, pch="*")  
plot(ProblemSetScore2, MidtermScore2, pch="*")  
sd(ProblemSetScore1)  
sd(MidtermScore1)  
cor(ProblemSetScore1,MidtermScore1)  
par(mfrow=c(1,2))  
qqnorm(ProblemSetScore1, main="Normal QQ plot of ProblemSetScore1")  
qqnorm(MidtermScore1, main="Normal QQ plot of MidtermScore1")  
# abline(a, b) graphs a line with intercept a # and slope b  
b = sd(MidtermScore1) / sd(ProblemSetScore1) ### THIS IS WRONG  
a = mean(MidtermScore1) - b * mean(ProblemSetScore1)  
plot(ProblemSetScore1, MidtermScore1, pch="*")  
abline(a, b, col="blue")  
# Cor line  
b = cor(ProblemSetScore1, MidtermScore1) * sd(MidtermScore1) / sd(ProblemSetScore1)  
a = mean(MidtermScore1) - b * mean(ProblemSetScore1)  
plot(ProblemSetScore1, MidtermScore1, pch="*")  
abline(a, b, col="red")  
quantile(ProblemSetScore1, 0.95)  
quantile(MidtermScore1, 0.95)  
h = quantile(ProblemSetScore1, 0.95)  
a+b*h
```

```
# using class2 as there is zero correlation in class1

sd(ProblemSetScore2)

sd(MidtermScore2)

cor(ProblemSetScore2,MidtermScore2)

par(mfrow=c(1,2))

qqnorm(ProblemSetScore2, main="Normal QQ plot of ProblemSetScore2")

qqnorm(MidtermScore2, main="Normal QQ plot of MidtermScore2")

# abline(a, b) graphs a line with intercept a # and slope b

b = sd(MidtermScore2) / sd(ProblemSetScore2) ### THIS IS WRONG

a = mean(MidtermScore2) - b * mean(ProblemSetScore2)

plot(ProblemSetScore2, MidtermScore2, pch="*")

abline(a, b, col="blue")

# Cor line

b = cor(ProblemSetScore2, MidtermScore2) * sd(MidtermScore2) / sd(ProblemSetScore2)

a = mean(MidtermScore2) - b * mean(ProblemSetScore2)

plot(ProblemSetScore2, MidtermScore2, pch="*")

abline(a, b, col="red")

quantile(ProblemSetScore2, 0.95)

quantile(MidtermScore2, 0.95)

h = quantile(ProblemSetScore2, 0.95)

a+b*h

# This value is lesser than the 95percentile of the midtermscore2.
```