

In The Name of God

Statistical Inference Homework #2

Amin Asadi Sarijaloo
810196410

Problem 1

- a. False, because the coin doesn't have memory, and the tosses are independent. Thus the chance for every toss to become head or tail is equal to 50%.
- b. False, because if two random variables A and B are independent, we have that $P(A | B) = P(A)$ which means that knowing B, gives no information about A, but for mutual exclusiveness, we must have $P(A \cap B) = 0$ which means knowing that B has happened, A can not happen, hence they are dependent(except when at least one of the events is impossible).
- c. False, the general additions rule states that: The probability of either X or Y occurring is: $P(X \cup Y) = P(X) + P(Y) - P(A \cap B)$
- d. False, we have that :
- $P(A | B) = P(A \cap B) / P(B)$
 - $P(B | A) = P(A \cap B) / P(A)$
- So in general $\rightarrow P(A | B) \neq P(B | A)$
- e. True

- f. False, $P(A \text{ and } B) = 0$ means that A and B are mutually exclusive. (One of them having zero probability is not necessarily true).
- g. False, the mean and variance of Poisson Distributions are both equal to λ which is the expected number of the occurrences of a particular event in unit measurement.
- h. True
- i. True

Problem 2

X: being infected to COVID-19

Y: speaking other languages than French

$$P(X) = 14.6\%$$

$$P(Y) = 20.7\%$$

$$P(X \cap Y) = 4.2\%$$

- a. No they aren't, because $P(X \text{ and } Y) = 4.2\% \neq 0$
- b. $P(X \cap \sim Y) = P(X) - P(X \cap Y) = 14.6\% - 4.2\% = 12.4\%$
- c. $P(X \cup Y) = P(X) + P(Y) - P(X \cap Y) = 14.6\% + 20.7\% - 4.2\% = 31.1\%$
- d. $P(\sim X \cap \sim Y) = 100\% - P(X \cup Y) = 68.9\%$
- e.
- $P(X \cap Y) = 4.2\%$
 - $P(X) * P(Y) = 14.6\% * 20.7\% = 3.02\%$
 - $\rightarrow P(X \cap Y) \neq P(X) * P(Y)$

because the necessary condition for independence doesn't hold,
they are not independent.

Problem 3

We denote the random variable of points he won in each game of round i with x_i

Calculating $E(x_1)$:

He predicts the correct winner with probability $\frac{1}{2}$, that is +1 with probability $\frac{1}{2}$ and +0 with probability $\frac{1}{2}$, so

$$E(x_1) = (1 * \frac{1}{2} + 0 * \frac{1}{2}) = \frac{1}{2}$$

Calculating $E(x_i)$:

He predicts the correct winner with probability $(\frac{1}{2})^i$ because the predicted team must have won in all the $(i-1)$ previous rounds, so because he earns 2^{i-1} for each correct prediction:

$$E(x_i) = (2^{i-1} * (\frac{1}{2})^i + 0) = \frac{1}{2}$$

Thus the expected value of the total points he wins is =

$$(32+16+8+4+2+1) * \frac{1}{2} = 63 * \frac{1}{2} = 31.5$$

Problem 4

$$P(L^c) = 1 - P(L) = 0.9$$

$$P(M_1 | L^c) = 1 - P(M_1^c | L^c) = 0.1$$

a.

$$P(L|M_1) = \frac{P(M_1 | L) P(L)}{P(M_1 | L) P(L) + P(M_1 | L^c) P(L^c)} = \frac{0.9 * 0.1}{0.9 * 0.1 + 0.1 * 0.9} = 0.5$$

$$\begin{aligned} \text{b. } P(L|M_1 \cap M_2) &= \frac{P(M_1 \cap M_2|L) P(L)}{P(M_1 \cap M_2|L) P(L) + P(M_1 \cap M_2|L^c) P(L^c)} = \\ &= \frac{P(M_1|L) P(M_2|L) P(L)}{P(M_1|L) P(M_2|L) P(L) + P(M_1|L) P(M_2|L^c) P(L^c)} = \frac{0.9^2 * 0.1}{0.9^2 * 0.1 + 0.1^2 * 0.9} = 0.9 \end{aligned}$$

- c. Yes, they are equivalent. We can use both methods for updating probabilities and they give the same results. This is because in both methods we use the same evidence to update our probabilities. The order of applying evidence does not matter.
-

Problem 5

$P(\text{a pill is damaged}) = 0.005$

$P(\text{a pill is not damaged}) = 1 - 0.005 = 0.995$

a.

This has a binomial distribution.

Each pill of 100 pills in a box is healthy independent from other pills with a probability of 0.995 and damaged with a probability of 0.005.

$P(\text{box having no damaged pills}) =$

$$\text{choose}(100, 0) \times 0.005^0 \times 0.995^{100} \times 100\% = 60.5 \%$$

b.

$P(\text{Having 2 or more damaged pills}) =$

$$\begin{aligned} &1 - \text{choose}(100, 0) \times 0.005^0 \times 0.995^{100} - \\ &\text{choose}(100, 1) \times 0.005^1 \times 0.995^{99} = 8.9 \% \end{aligned}$$

c. In the image below the result of the R-script is the same as the result of parts a and b:

- Calculation with “dbinom” function of R:

```
1 ##### CALCULATION #####  
2 cat('part one calculation percentage:', dbinom(0, 100, 0.005) * 100, '%\n')  
3  
4 cat('part two simulation percentage:', (1 - dbinom(0, 100, 0.005) - dbinom(1, 100, 0.005)) * 100, '%\n')  
5 |  
6
```

- Simulation: I created 1000 boxes with 100 pills in each. Each pill was damaged with a probability of 0.5%. Then I calculated the average number of boxes having no damaged pills and the average number of boxes having at least 2 damaged pills among all 1000 boxes:

```

7 ▾ ##### SIMULATION #####
8 ▾ create_one_box <- function() {
9   boxes = list(sample(c(0, 1), size = 100, replace = TRUE, prob = c(0.005, 0.995)));
10  return(boxes)
11 ^ }
12
13 ▾ create_boxes <- function() {
14   num_boxes <- 1000;
15   boxes <- c()
16   for (i in 1:num_boxes)
17     boxes <- c(boxes, create_one_box())
18   return(boxes)
19 ^ }
20
21 ##### PART a : Finding percentage of boxes which have no damaged pill
22 num_epochs <- 100
23 all_healthy_percentages <- c()
24 ▾ for (i in 1:num_epochs) {
25   boxes <- create_boxes()
26   num_healthy_boxes = sum(sapply(boxes, function (x) sum(unlist(x)) == 100))
27   percentage = num_healthy_boxes / num_boxes * 100
28   all_healthy_percentages <- c(all_healthy_percentages, percentage)
29 ^ }
30
31 cat('part one simulation percentage:', mean(all_healthy_percentages), '\n')
32
33
34 atleast_two_damaged_percentages <- c()
35 ▾ for (i in 1:num_epochs) {
36   boxes <- create_boxes()
37   num_healthy_boxes = sum(sapply(boxes, function (x) sum(unlist(x)) <= 98))
38   percentage = num_healthy_boxes / num_boxes * 100
39   atleast_two_damaged_percentages <- c(atleast_two_damaged_percentages, percentage)
40 ^ }
41
42 cat('part two simulation percentage:', mean(atleast_two_damaged_percentages), '\n')
43

```

Result of R codes:

As we can see the result of both the calculation code and simulation code are equal to what we found in parts a and b:

```

> source('~\Desktop\uni\Spring 00\SI\hws\hw2\p5.r')
part one calculation percentage: 60.57704 %
part two simulation percentage: 8.982231 %
part one simulation percentage: 60.589 %
part two simulation percentage: 8.941 %
> |

```

Problem 6

Since every passenger has two possibilities(showing up or not), each of them have a Bernoulli distribution.

a.

Let T be the random variable indicating the total number of people who show up. The objective is to find the probability of T being more than 300(overbook). T has binomial distribution:

$$m = E(T) = Np = 324 \times 0.9 = 291.6$$

$$\sigma = \sigma(T) = \sqrt{Np(1-p)} = \sqrt{324 \times 0.9 \times 0.1} = 5.4$$

$$P(T \geq 301) = P(T \geq 301 - 0.5) = P\left(\frac{T-m}{\sigma} \geq \frac{301-0.5-m}{\sigma}\right)$$

$$1 - \Phi\left(\frac{301-0.5-324 \times 0.9}{\sqrt{324 \times 0.9 \times 0.1}}\right) = 0.04$$

b.

There are $324/2=162$ pairs. Each pair will travel with probability of $0.9 \times 0.9 = 0.81$ independent of other pairs.

mean of the pairs will be $162 \times 0.81 = 131.22$

mean of the people will be $2 \times 162 \times 0.81 = 262.44$

Std of pairs will be $\sqrt{162 \times 0.81 \times (1 - 0.81)}$

Std of people will be $\sqrt{2 \times 162 \times 0.81 \times (1 - 0.81)} = 7.06$

So the mean is decreased a lot but std is increased a little. So the probability is decreased comparing to the previous part.

Problem 7

a.

Shorthand for exams = $N(462, 119)$

Shorthand for projects = $N(584, 151)$

```
1 library(ggplot2)
2
3 #a
4 exam_mean <- 462;
5 exam_std <- 119;
6 project_mean <- 584;
7 project_std <- 151
8 # Shorthand for exams = N(462, 119)
9 # Shorthand for projects = N(584, 151)
10
11 exam_score <- 620;
12 project_score <- 670;
13
```

b.

```
14 #b
15 exam_z_score = (exam_score - exam_mean) / exam_std;
16 project_z_score = (project_score - project_mean) / project_std;
17
18 cat('exam z-score is:', exam_z_score, '\n')
19 cat('project z-score is', project_z_score, '\n')
```

Result:

```
exam z-score is: 1.327731
project z-score is 0.5695364
```

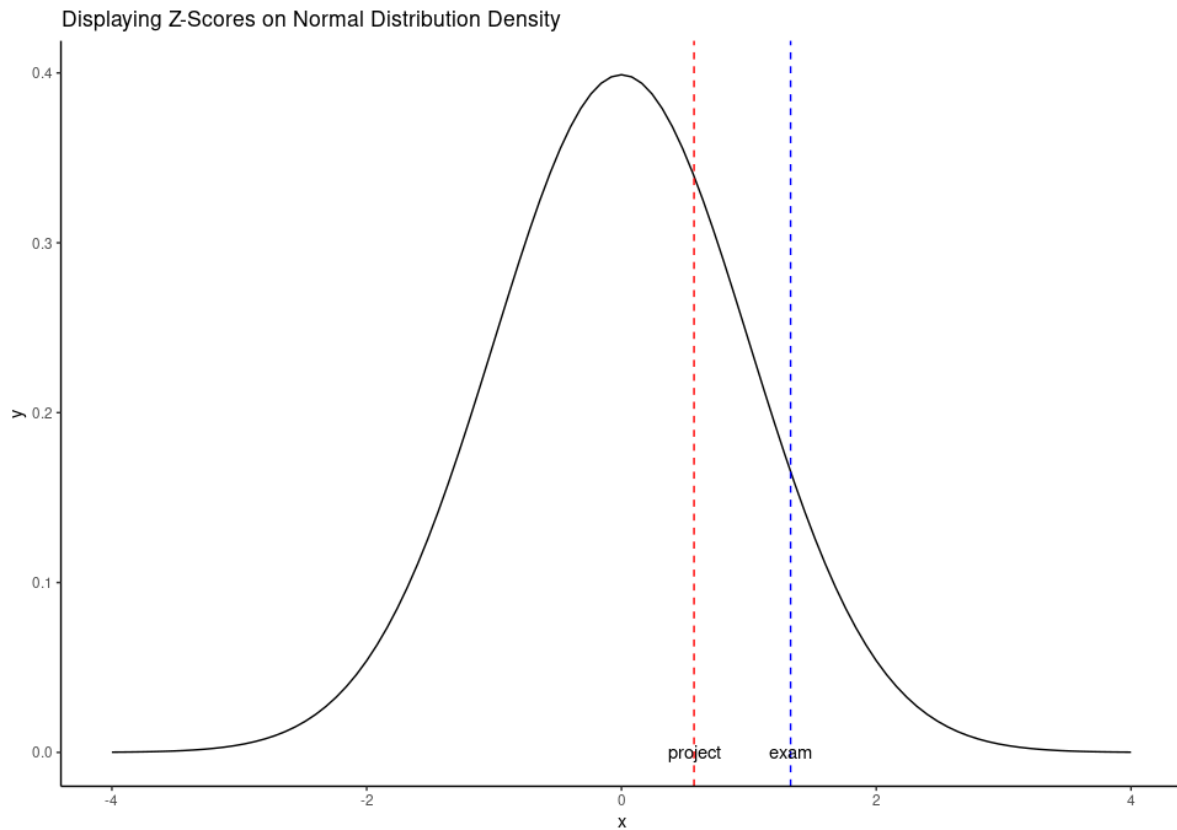
c.


```

21 #c
22 p1 <- ggplot(data.frame(x = seq(-4, 4, length=100)), aes(x = x)) +
23
24     stat_function(fun = dnorm) +
25
26     geom_vline(xintercept = c(exam_z_score, project_z_score),
27                 colour=c('blue', 'red'),
28                 linetype = "dashed") +
29     geom_text(aes(x = c(exam_z_score, project_z_score),
30                     y = 0,|
31                     label = c('exam', 'project')),
32               data=data.frame(c(exam_z_score, project_z_score)))
33 print(p1)
34

```

Result:



d.

a z-score tells us that how far is a datapoint from the mean, i.e. It measures the number of standard deviations an observation falls below or above the mean. So her exam score is 1.32 standard deviations far from the mean and her project score is 0.56 standard deviations far from the mean.

e.

z-score of her exam is more than z-score of her project so she performed better in exam.

f.

```
36 #f
37 exam_percentile_score <- pnorm(exam_z_score, 0, 1)
38 project_percentile_score <- pnorm(project_z_score, 0, 1)
39 cat('percentile score of exam:', exam_percentile_score , '\n')
40 cat('percentile score of project ', project_percentile_score, '\n')
41
```

Result:

```
percentile score of exam: 0.9078665
percentile score of project 0.7155039
```

g.

```
42 #g
43 cat(100 - exam_percentile_score * 100, "students performed better than her in exam\n")
44
```

Result:

```
9.213348 % of students performed better than her in exam
```

h.

```
45 #h
46 cat(100 - project_percentile_score * 100, "students performed better than her in project\n")
47
```

Result:

```
28.44961 % of students performed better than her in project
```

i.

```

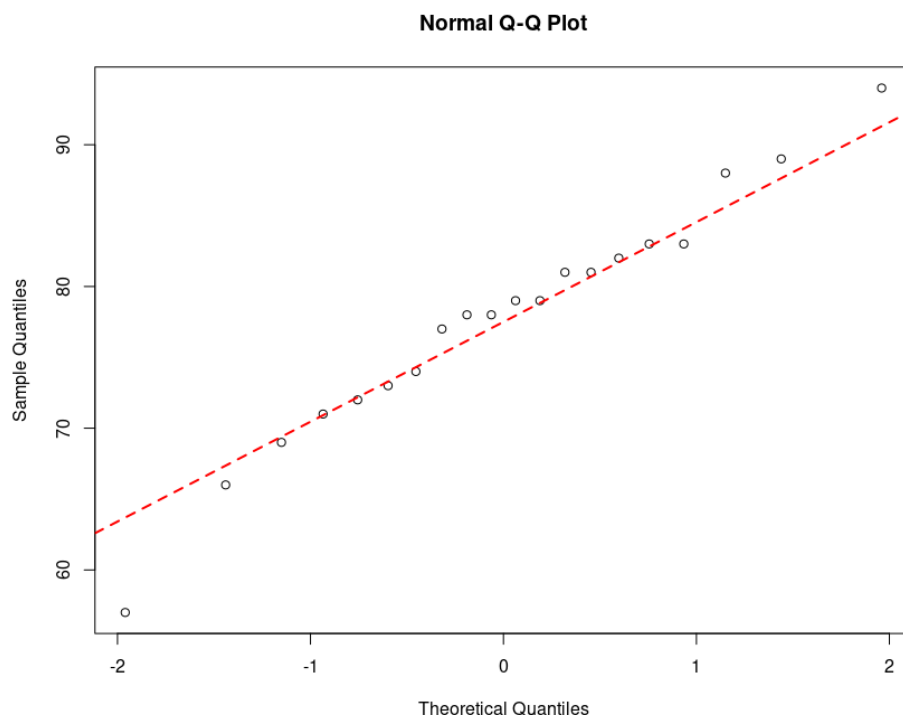
48 #i
49 scores <- c(57, 66, 69, 71, 72, 73, 74, 77, 78, 78, 79, 79, 81, 81, 82, 83, 83, 88, 89, 94)
50
51 ## QQplot
52 qqnorm(scores)
53 qqline(scores, col='red', lwd=2, lty=2)
54
55 ## Histogram
56 h <- hist(scores)
57 xfit <- seq(min(scores), max(scores), by=0.01)
58 yfit <- dnorm(xfit, mean = mean(scores), sd = sd(scores)) * 5 * length(scores)
59 lines(xfit, yfit, col="blue")
60
61 ## Box Plot
62 boxplot(scores)
63 title(main = "scores box plot", ylab = "score")

```

Result:

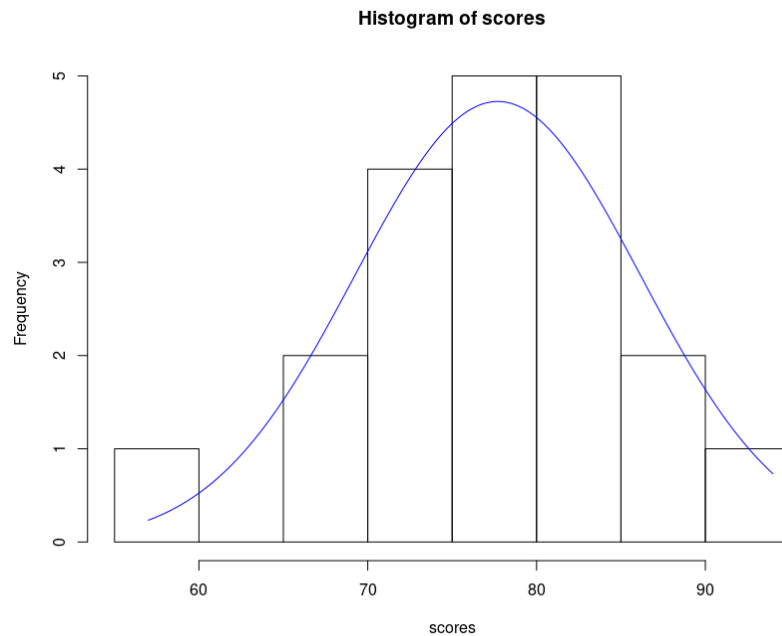
QQ Plot:

We can see that at some parts specially at first before first quartile the qqplot of scores is different from qqplot of normal distribution but approximately they are similar.



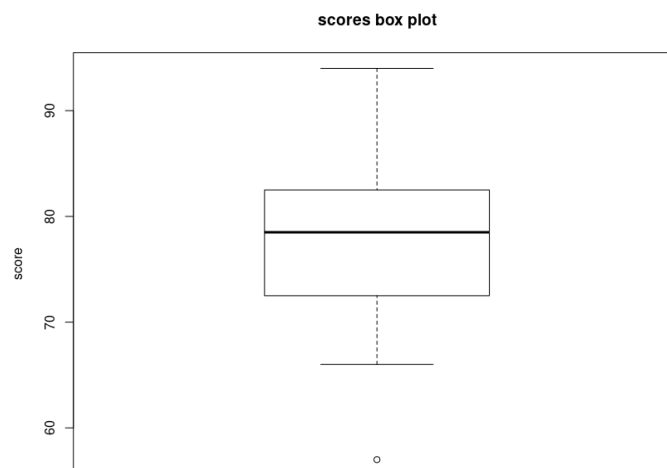
Histogram:

The histogram and normal curve have similar shapes in general but they are different before first quartile of data.



Box Plot:

Because the mean and median are close to the center of the box, the distribution is close to normal but the box plot is not symmetric.



Problem 8

a) P_1 (Nobody gets back his own raincoat and umbrella) =

$1 - P_2$ (there is at least one person who gets back both his raincoat and umbrella)

E_i : person i gets back his own raincoat and umbrella

$$\rightarrow 1 - P_2 = 1 - P\left(\bigcup_{i=1}^n E_i\right) = 1 - \left(\sum_{i=1}^n P(E_i) - \sum_{\substack{i,j=1 \\ j \neq i}}^n P(E_i \cap E_j) + \dots\right)$$

$$= 1 - \left(\binom{n}{1} \times \frac{1}{n^2} - \binom{n}{2} \times \frac{1}{n(n-1)^2} + \dots\right)$$

$$= 1 - \left(\frac{n!}{1!(n-1)!} \times \frac{1}{n^2} - \frac{n!}{2!(n-2)!} \times \frac{1}{n^2(n-1)^2} + \dots\right)$$

$$= 1 - \left(\frac{1}{n} - \frac{1}{2!n(n-1)} + \frac{1}{3!n(n-1)(n-2)} - \dots\right)$$

$$= 1 - \frac{1}{n} + \frac{1}{2!n(n-1)} - \frac{1}{3!n(n-1)(n-2)} + \dots + \frac{(-1)^n}{n!} = \sum_{i=0}^n \frac{(-1)^i (n-i)!}{i! n!}$$

b) C_i = person i gets ^{back} neither his raincoat nor his umbrella

$\rightarrow P(\text{everybody gets back at least his raincoat or umbrella}) = 1 - P(\text{there is at least a person who does not get back his umbrella and raincoat})$

$$c) (1 - 0.1)^{\frac{N}{2}} (1 - 0.2)^2 = (0.9 \times 0.8)^{\frac{N}{2}} = 0.72^{\frac{N}{2}}$$

Problem 9

a.

$$1) \quad V_1 = P(\text{the tried VPN is the better one})$$

$$V_2 = P(\text{the " " " " worse one})$$

$$P(V_1 | \text{stablisthed}) = \frac{P(\text{stablisthed} | V_1) P(V_1)}{P(\text{stablisthed} | V_1) P(V_1) + P(\text{stablisthed} | V_2) P(V_2)}$$

$$= \frac{0.3 \times 0.5}{0.3 \times 0.5 + 0.1 \times 0.5} = \frac{3}{4} = \underline{\underline{0.75}}$$

9.b.

Code:

```
1 connect_vpn <- function(p) {
2   did_stablish <- as.logical(rbinom(1, size=1, prob=p))
3   return(did_stablish)
4 }
5
6 select_vpn_randomly <- function() {
7   if(as.logical(rbinom(1, size=1, prob=0.5)))
8     return(1)
9   else
10    return(2)
11 }
12
13 experiment <- function() {
14   num_try = 0
15   while(probabilities[1] < 0.9) {
16     if(probabilities[1] < 1e-4)
17       return(0)
18     i <- select_vpn_randomly()
19     stablished_vpn <- connect_vpn(success_rates[i])
20     if(stablished_vpn) {
21       probabilities[i] <- (success_rates[i] * probabilities[i]) /
22         (success_rates[i] * probabilities[i] + success_rates[3-i] * probabilities[3-i])
23     }
24     else {
25       probabilities[i] <- ((1-success_rates[i]) * probabilities[i]) /
26         ((1-success_rates[i]) * probabilities[i] + (1-success_rates[3-i]) * probabilities[3-i])
27     }
28     probabilities[3-i] <- 1 - probabilities[i]
29
30     num_try <- num_try + 1
31   }
32   return(num_try)
33 }
34
35 success_rates <- c(0.3, 0.1)
36
37 tries = c()
38 for(j in (1:100)) {
39   print(j)
40   probabilities <- c(0.5, 0.5)
41   num_try = experiment()
42   tries <- c(tries, num_try)
43 }
44 cat("mean_trials to be 90% certain = ", mean(tries))
45
```

Result:

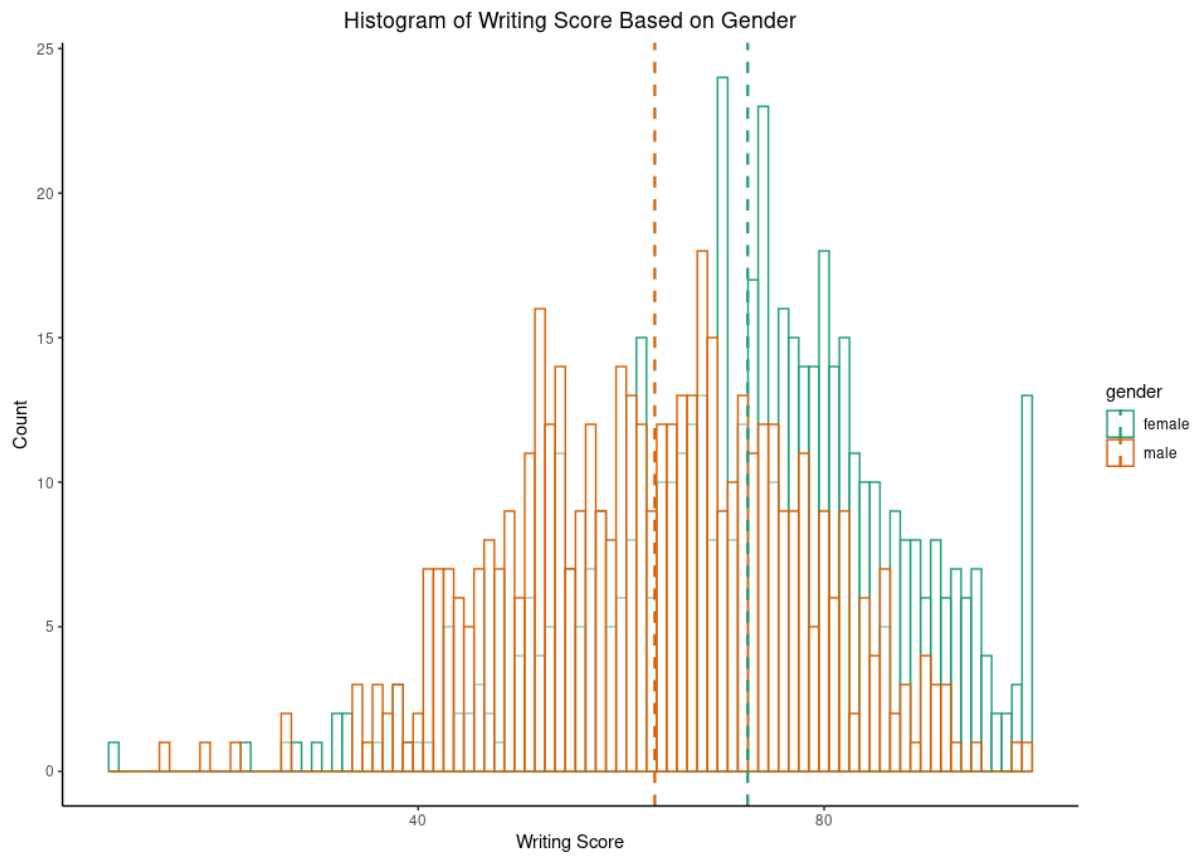
The result changes each time we run the experiments but it is about 50

Problem 10

a.

```
1 library(ggplot2)
2
3 data <- read.csv('StudentsPerformance.csv')
4
5 # a
6
7 males = data[data["gender"]=="male", ]
8 females = data[data["gender"]=="female", ]
9
10 means <- c(mean(males$writing_score), mean(females$writing_score));
11
12 p_a <- ggplot(data,
13               aes(x=writing_score,
14                   color=gender)) +
15   geom_histogram(binwidth=1,
16                 fill='white',
17                 alpha=0.5,
18                 position="identity") +
19
20   geom_vline(data=data.frame(gender=c('male', 'female') , means=means),
21             aes(xintercept=means, color=gender),
22             linetype="dashed",
23             size=0.8) +
24
25   scale_color_brewer(palette="Dark2") +
26   labs(title="Histogram of Writing Score Based on Gender", x="Writing Score", y = "Count")+
27   theme_classic() +
28   theme(plot.title = element_text(hjust = 0.5)) +
29   scale_x_continuous(breaks=seq(0, 100, 40))
30
31 show(p_a)
```

a. result:



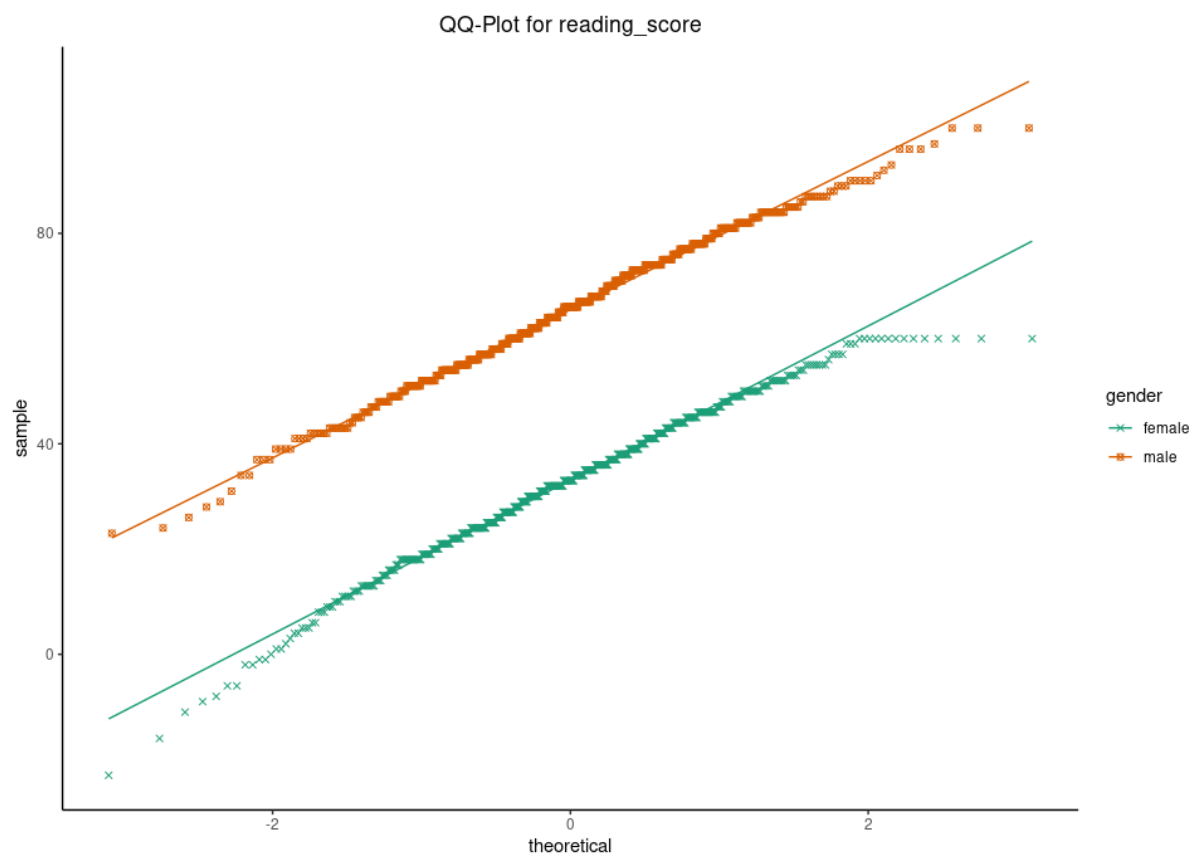
b.

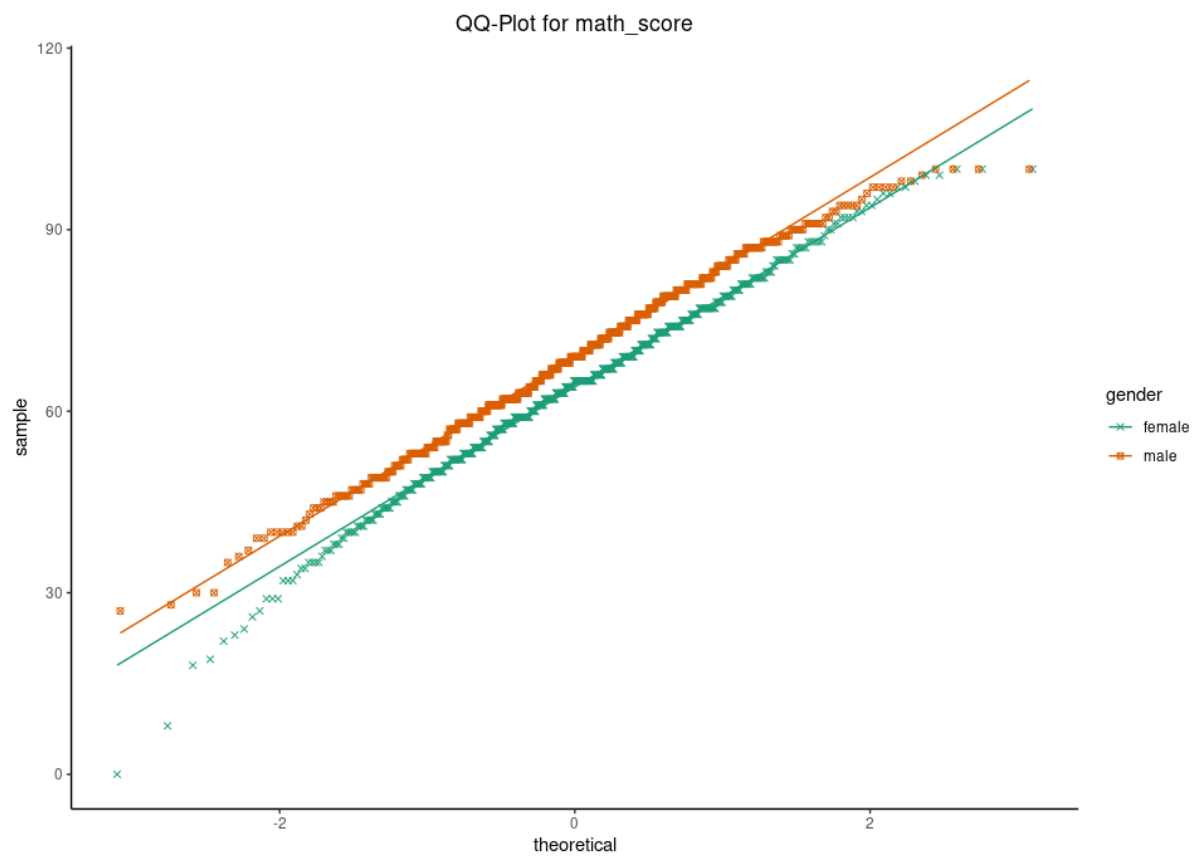
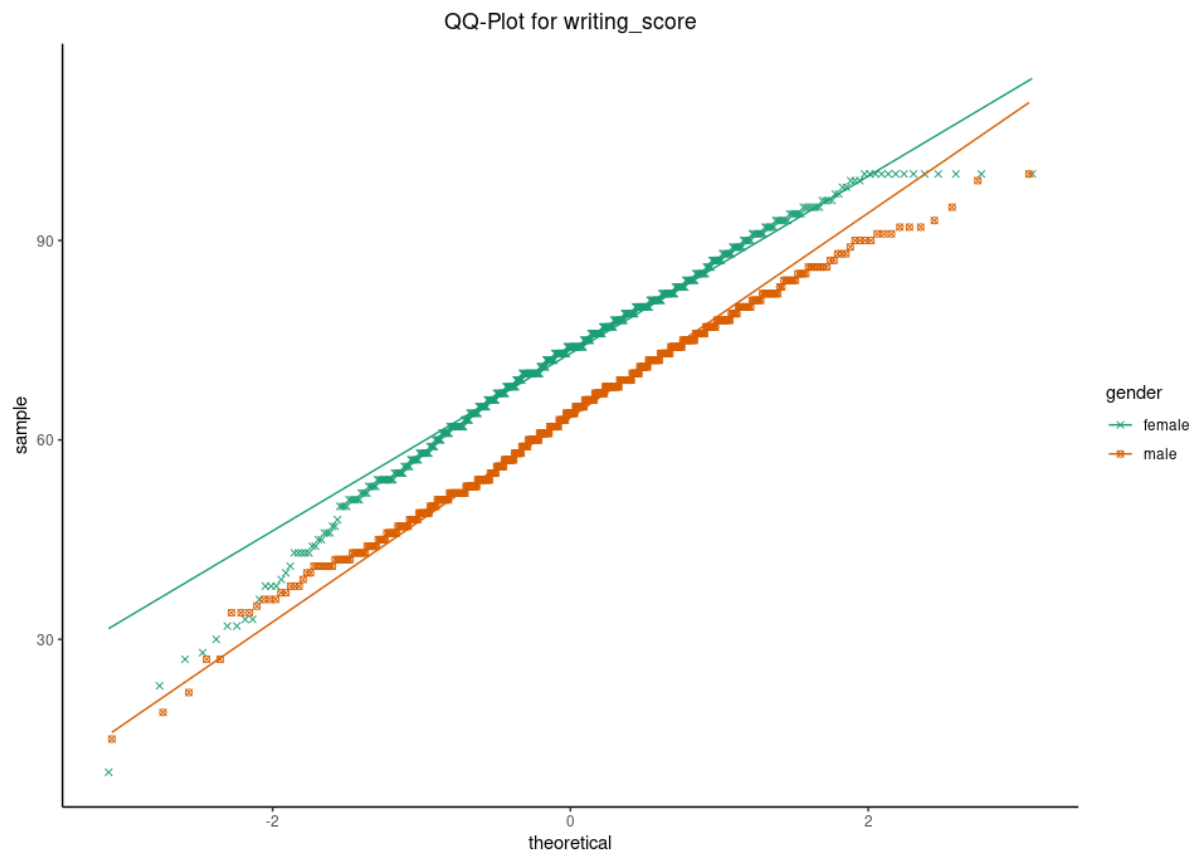
```

34 # b
35 for (subj in c('reading_score', 'writing_score', 'math_score')) {
36   p_b <- ggplot(data,
37     aes_string(sample=subj, shape='gender',
38       colour='gender')) +
39     stat_qq() +
40     stat_qq_line() +
41     scale_shape_manual(values=c(4,13)) +
42     scale_color_brewer(palette="Dark2") +
43     theme_classic() +
44     labs(title=paste("QQ-Plot for", subj)) +
45     theme(plot.title = element_text(hjust = 0.5))
46     show(p_b)
47 }

```

b.results:

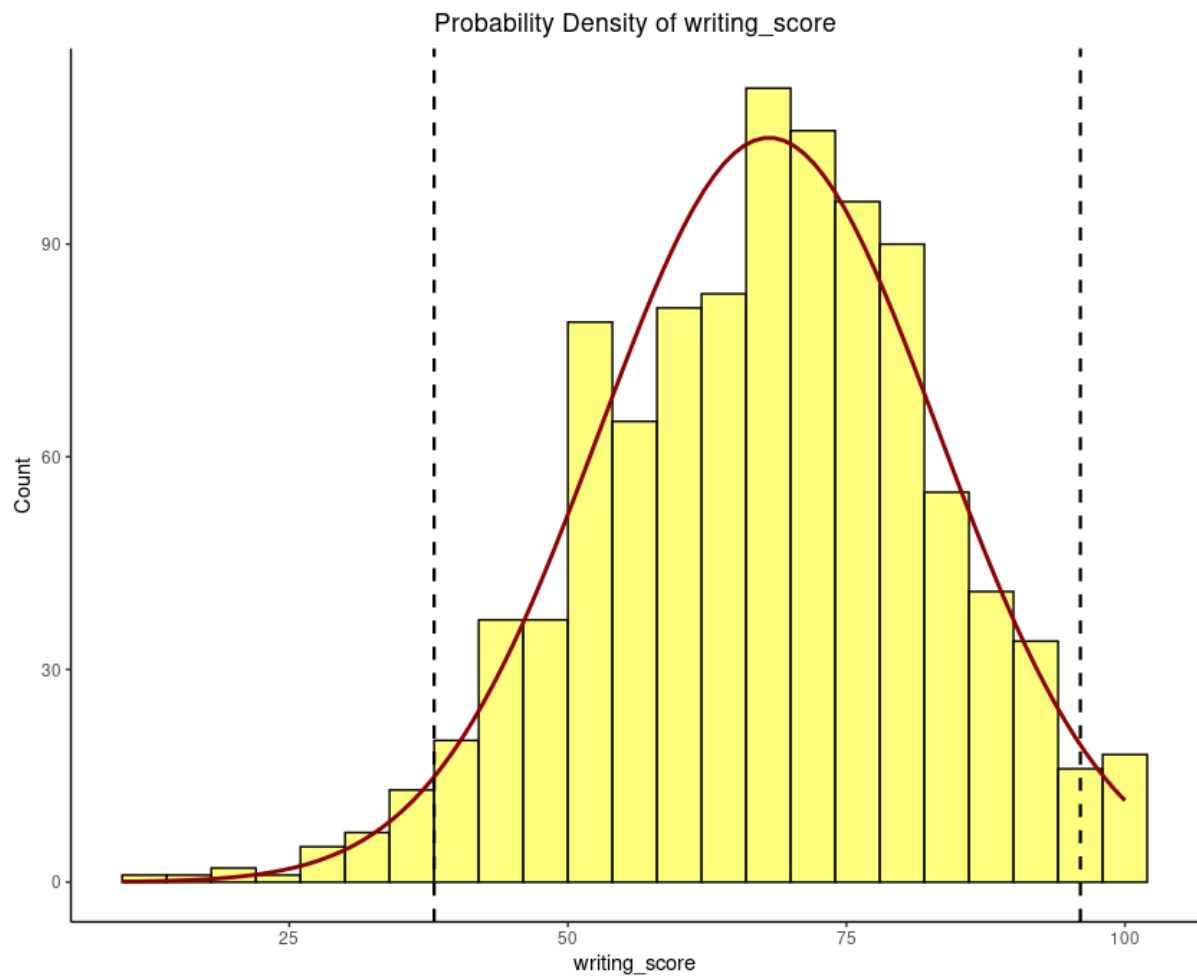


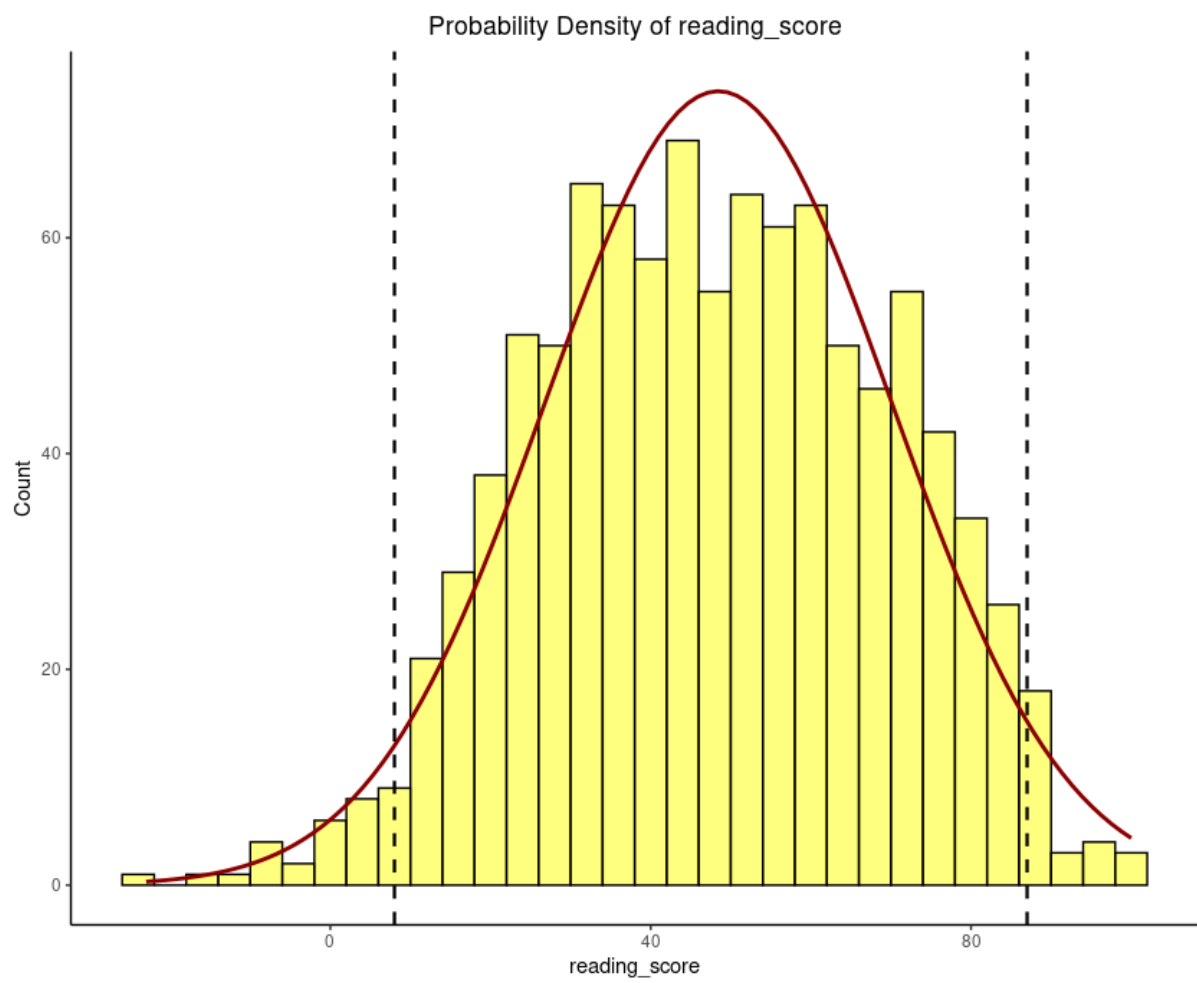


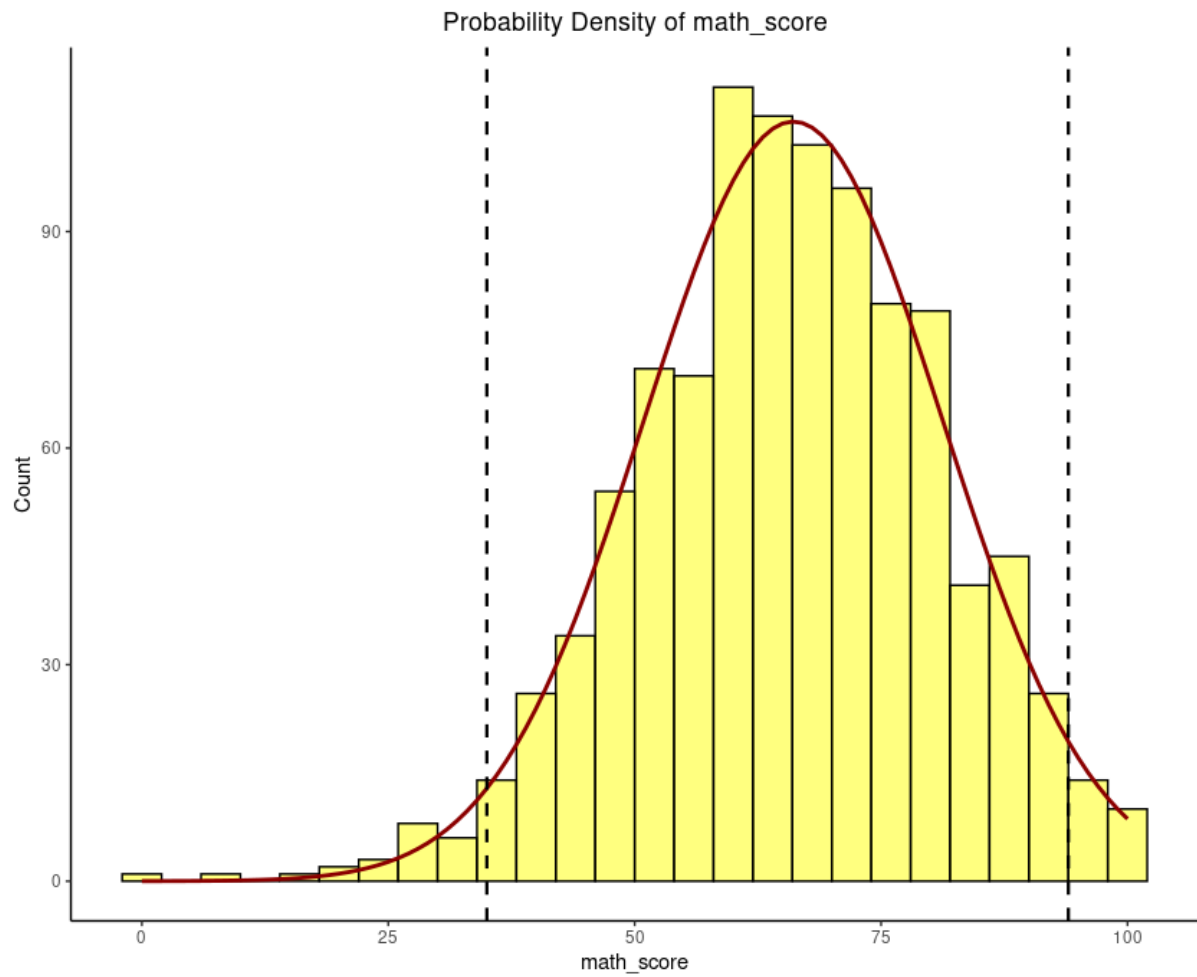
C.

```
49 # c
50
51 reading_score_mean <- mean(data$reading_score)
52 writing_score_mean <- mean(data$writing_score)
53 math_score_mean <- mean(data$math_score)
54
55 reading_score_sd <- sd(data$reading_score)
56 writing_score_sd <- sd(data$writing_score)
57 math_score_sd <- sd(data$math_score)
58
59 bin_width = 4;
60 num_samples = nrow(data);
61 |
62 ▾ scaled_dnorm = function(x, mean, sd, n) {
63   bin_width * num_samples * dnorm(x, mean, sd)
64 ^ }
65
66 ▾ for (subj in c('reading_score', 'writing_score', 'math_score')) {
67   p_c <- ggplot(data,
68     aes_string(x=subj)) +
69
70     geom_histogram(binwidth=bin_width,
71       color='black',
72       fill='yellow',
73       position="identity",
74       alpha=0.5) +
75     stat_function(fun=scaled_dnorm,
76       args=list(data[[subj]],
77         mean=mean(data[[subj]]),
78         sd = sd(data[[subj]]),
79         colour = "darkred", size = 1) +
80     geom_vline(xintercept=unname(quantile(data[[subj]], c(0.025, 0.975))),
81       linetype="dashed",
82       size=0.8) +
83     labs(title=paste("Probability Density of", subj), x=subj, y = "Count")+
84     theme_classic() +
85     theme(plot.title = element_text(hjust = 0.5))
86     show(p_c)
87 ^ }
```

c.Results:



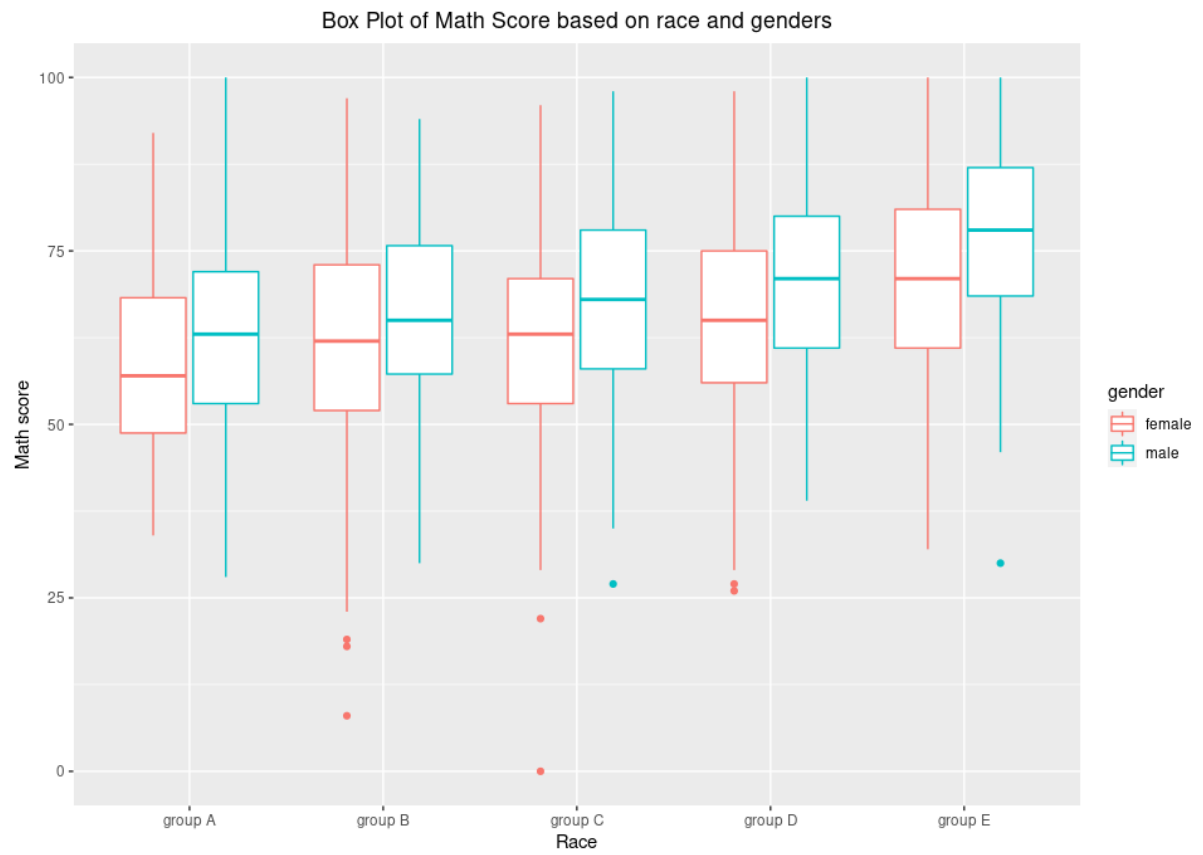




d.

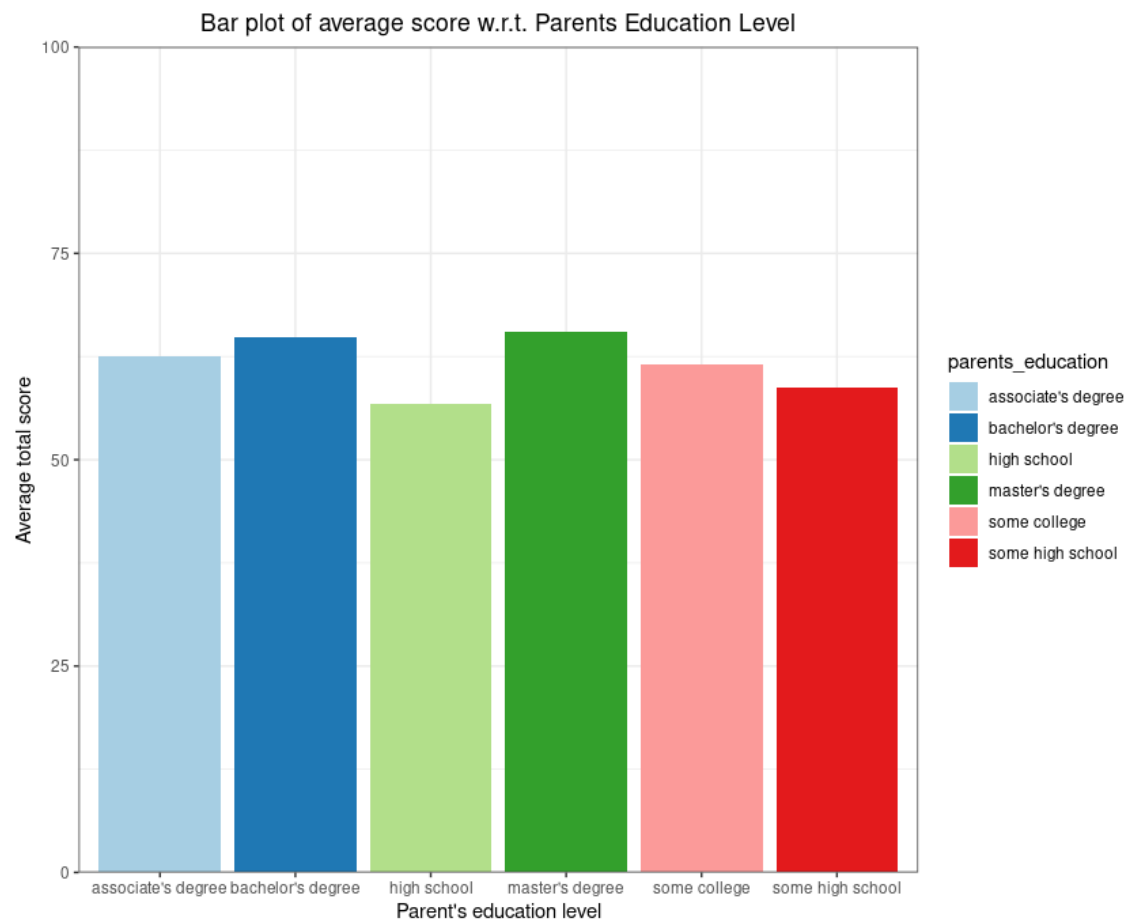
```
89 #d
90 p_d <- ggplot(data, aes(race, math_score)) +
91   geom_boxplot(aes(colour=gender)) +
92   labs(title="Box Plot of Math Score based on race and genders", x="Race", y = "Math score") +
93   theme(plot.title = element_text(hjust = 0.5))
94 show(p_d)
```


d.Results:

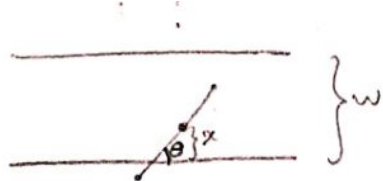


e.

```
96 # e
97 data['avg_score'] <- round(rowMeans(data[, c('math_score', 'reading_score', 'writing_score')]), 2);
98
99 p_e <- ggplot(data, aes(x=parents_education, y=avg_score, fill=parents_education)) +
100   geom_bar(stat="summary", fun="mean") +
101   scale_fill_brewer(palette="Paired") +
102   theme(axis.text.x=element_blank(),
103         axis.ticks.x=element_blank()) +
104   scale_y_continuous(expand = c(0, 0), limits = c(0, 100)) +
105   theme(axis.line.x = element_line(size=1.5, colour="grey")) +
106   labs(title="Bar plot of average score w.r.t. Parents Education Level", x="Parent's education level",
107        y = "Average total score") +
108   theme_bw() +
109   theme(plot.title = element_text(hjust = 0.5))
110 show(p_e)
111
```



Problem 11



x : distance between center of knife-trace and the closest line
 θ : acute angle between knife trace and closest line

always:

$$0 \leq x \leq \frac{w}{2} \rightarrow f(x) = \begin{cases} \frac{1}{\frac{w}{2}} = \frac{2}{w} & 0 \leq x \leq \frac{w}{2} \\ 0 & \text{o.w.} \end{cases}$$

$$0 \leq \theta \leq \frac{\pi}{2} \rightarrow f(\theta) = \begin{cases} \frac{1}{\frac{\pi}{2}} = \frac{2}{\pi} & 0 \leq \theta \leq \frac{\pi}{2} \\ 0 & \text{o.w.} \end{cases}$$

\rightarrow we know that x and θ are independent

$$\rightarrow f(x, \theta) = f(x)f(\theta) = \begin{cases} \frac{4}{\pi w} & 0 \leq x \leq \frac{w}{2}, 0 \leq \theta \leq \frac{\pi}{2} \\ 0 & \text{o.w.} \end{cases}$$



\rightarrow doesn't cross a line if $\frac{l}{2} \sin \theta \leq x \leq \frac{w}{2}$

$$\begin{aligned} \rightarrow P(\text{fall into one strip}) &= \int_{\theta=0}^{\frac{\pi}{2}} \int_{x=\frac{l}{2} \sin \theta}^{\frac{w}{2}} \frac{4}{\pi w} dx d\theta = \int_0^{\frac{\pi}{2}} \frac{4}{\pi w} \left(\frac{w}{2} - \frac{l}{2} \sin \theta \right) d\theta \\ &= \frac{4}{\pi w} \left(\frac{\pi}{2} \times \frac{w}{2} + \frac{l}{2} \cos \theta \right) \Big|_0^{\frac{\pi}{2}} \\ &= \frac{4}{\pi w} \left(\frac{w\pi}{4} - \frac{l}{2} \right) = 1 - \frac{2l}{\pi w} \end{aligned}$$

$$\rightarrow P(3 \text{ successive round}) = \left(1 - \frac{2l}{\pi w} \right)^3$$

because rounds are independent