

Q1] Explain the property and F distribution using R software

? The Fisher's F distribution: The F distribution is a distribution of all possible value of F statistic. It is a continuous random variable. Consider two random variables χ^2_1 and χ^2_2 such that $\chi^2_1 \sim \text{Chi square}_n$ and $\chi^2_2 \sim \text{Chi square}_m$. Where n and m are degrees of freedom. The random variable F is defined by $F = \frac{\chi^2_1}{\chi^2_2}$.

Properties of F distribution

- 1) The F distribution is positively skewed and with increase in degrees of freedom n or m its skewness decreases.
- 2) The value of F distribution is always positive or zero since the variance are the square of the deviation and hence cannot be negative.
- 3) It is used to calculate the value of mean and variance.
- 4) The shape of F distribution depends on its parameter and n and m degrees of freedom.

Code

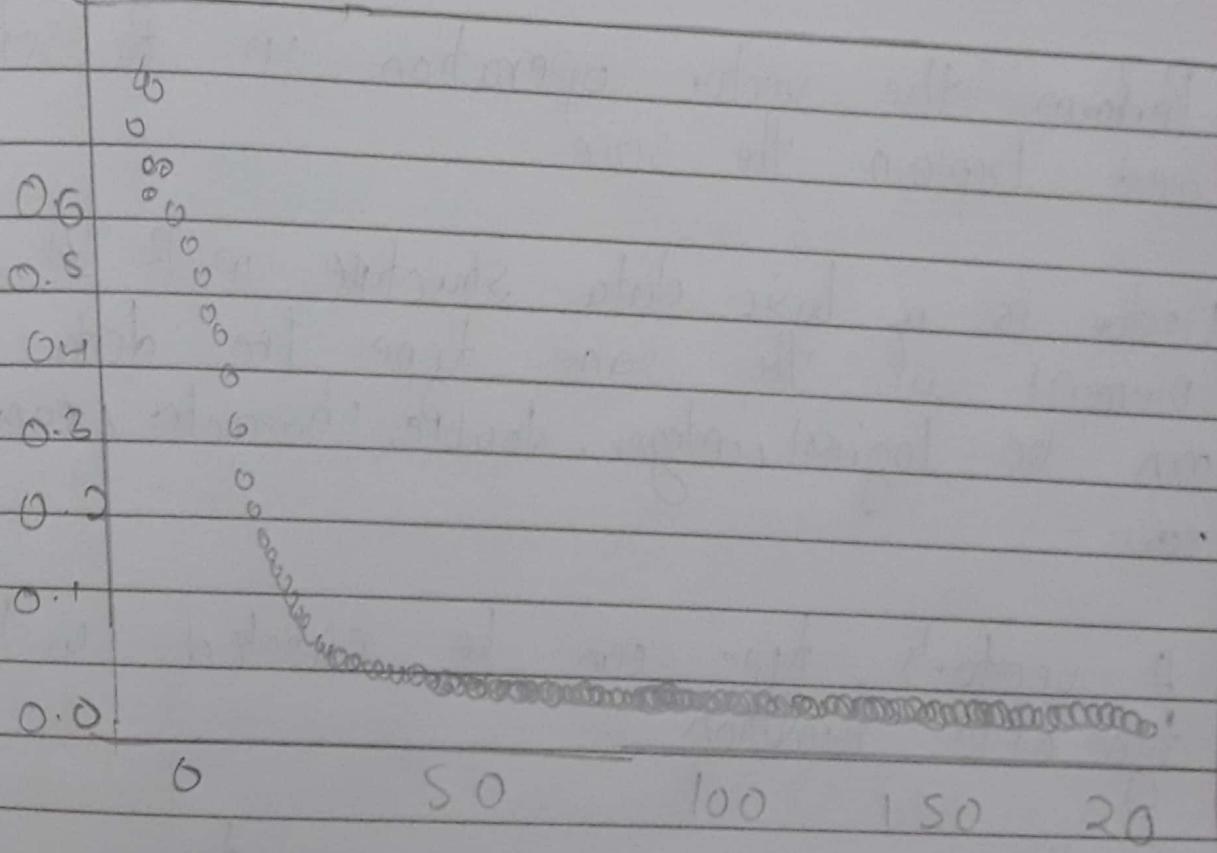
```
?x_df <- seq(0, 20, by = 0.1)
```

```
?y_df <- df(x = df, df1 = 3, df2 = 5)
```

```
?plot(y = df)
```

F-12a 41ppay

Page No. _____
Date. ALANRAO

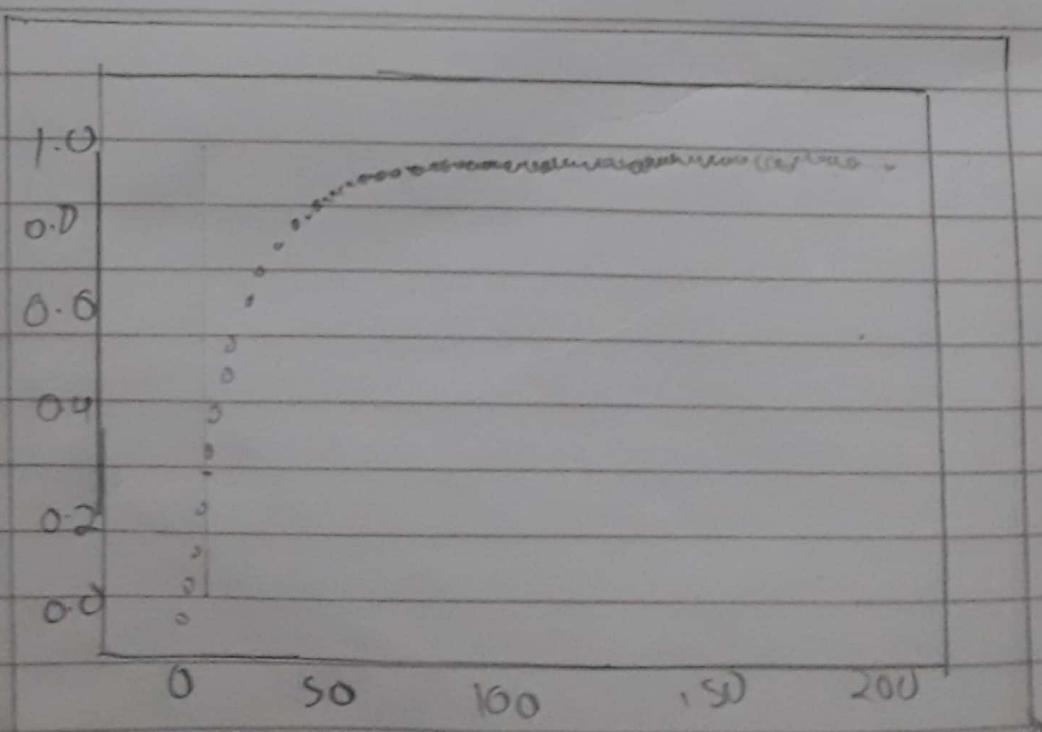


(Q42) Generate the plot for F distribution also mention the application of F distribution

→ A probability density function that is used especially in analysis of variance and is a function of the ratio of two independent random variables, each of which has a chi square distribution and is divided by its number of degrees of freedom.

Code

$x = \text{pf}^{-1}(0.28, df1=3, df2=5)$
 $y = \text{pf}^{-1}(x, df1=3, df2=5)$
 $\text{plot}(y)$



(iii) Calculate the f software in Anova using R Software and explain process

Code

```
→ > deep.Data = read.csv(file.choose())
> summary(deep.Data)
> hist(deep.Data$yield)
Deep One-way L-aov (yield) ~ fertilizer, Data = Deep.Data
> summary(Deep.OneWay)
> deep.TwoWay L-aov (yield ~ fertilizer + density,
Data = Deep.Data)
> summary(Deep.TwoWay)
> deep.TwoWayInteraction L-aov (yield ~ fertilizer * density,
Data = Deep.Data)
> summary(Deep.TwoWayInteraction)
> deep.TwoWayBlocking f-aov (yield ~ fertilizer + density +
block, Data = Deep.Data)
> summary(Deep.TwoWay)
```

30

25

20

15

10

5

0

175 176

177

178

179

Q1) Generate the RUN chart in R software to understand the RUN in data set

- A run chart is a line graph of data plotted over time. By collecting and charting data over time, you can find trends or pattern in process. Because they do not use control limits, run charts cannot tell you if a process is stable. However, they can show you how the process is running.
- The run chart is can be valuable tool at the beginning of the project, as it reveals important information about a process before you have collected enough data.

Code:

```
library(gcharts)
```

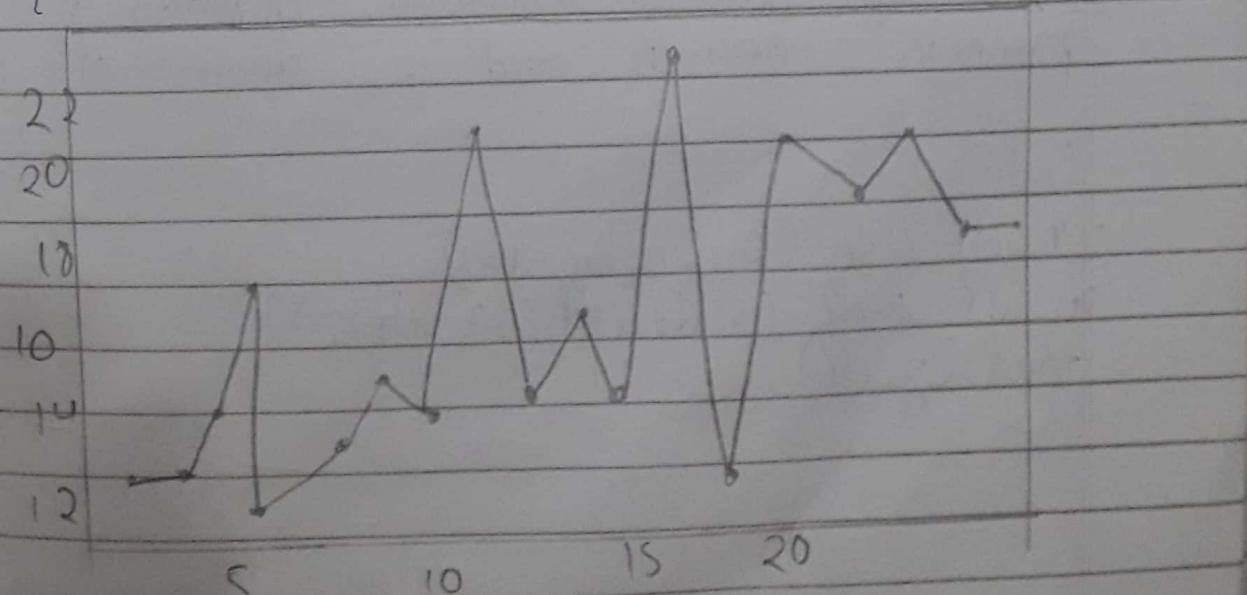
```
set.seed(9)
```

```
y2=rpois(24,16)
```

```
qic(y)
```

```
y[13:24] <- rpois(12,24)
```

```
qic(y)
```



(Q4) For a data set generate the stem and leaf plot and explain the same.

- 1) A stem and leaf plot is a special type where each data value is split into a 'stem' (the first digit) and a 'leaf' (usually the last digits or digits).
- 2) For example "32" is split into "3" (stem) and "2" (leaf). The stem values are listed on the left and the leaf values are listed next to them.
- 3) This way the stem groups the scores and each leaf indicates a score within that group.
- 4) Unlike histogram stem and leaf display retain the original frequency data at least two significant digits.

Code :

grades = [87, 53, 74, 82, 76, 93, 68, 47, 88, 25, 71, 79, 91, 80, 65, 93, 86, 87, 63, 74]
→ stem (grades)

The decimal point is 3 digits to the right of the unit.

5 | 3

6 | 3 5

7 | 1 1 4 6 9

8 | 0 2 5 6 7 7

9 | 1 3 5

Name - Neeraj Appu
f-129

Page No. _____
Date 17.11.2018

System (grads) scale = 3

The decimal point is 1 digit (.) to right of thy

417

513

51

613

6158

7144

7169

8102

(a) Explain the post hoc analysis in R software

- 1) Post-hoc means to analyse the results of your experimental data
- 2) Without all data in an experiment post hoc analysis is not possible.
- 3) Post-hoc analysis is performed to analyse the difference between data of groups
- 4) Post-hoc is the tests attempt to control the experiment wise error rate.

Code

Drug testing - Data frame (`Text`) = `data.frame`(
 `Drug 1` = c(5.40, 5.50, 5.55,
 5.85, 5.70, 5.75,
 5.20, 5.60, 5.50
 5.55, 5.50, 5.40
 5.90, 5.25, 5.70
 5.40, 5.40, 5.3)
 `Drug 2` = c(5.45, 5.50, 5.35,
 5.25, 5.20, 5.10
 5.65, 5.55, 5.25
 5.05, 5.00, 4.95
 5.50, 5.50, 5.40
 5.45, 5.55, 5.35
 6.30, 6.30, 6.35))

`Sci` = factor (`rep` (c("Drug 1", "Drug 2", "Drug 3"), 10),
 levels = factor (`rep` (1:20, 3)))
with (Drug testing, boxplot (Text ~ Sci)),
 manova (Drug testing ~ Sci),
 post.hoc (test `Sci`)

Output

Drug 2 - Drug 1 - 0.79300209

Drug 3 - Drug 1 - 0.01602425

Drug 3 - Drug 2 - 0.0880336

(Q7) Explain Pairwise mean difference test in R software?

- 1) What is Pairwise mean difference test?
- 2) Pairwise mean difference test is also called as the paired t-test
- 3) The paired t-test is used to compare the means between two related groups of samples
- 4) In this case, you have two values for samples
- 5) The purpose of test is - to determine whether there is statistical evidence that the mean difference between paired observations on a particular outcome is significantly different from zero. The Paired T-Test is a parametric test

Code

before <- c(200.1, 190.9, 192.7, 213, 241.4, 196.9, 172.3, 185.5, 203.2, 193.7)

after <- c(392.4, 393.2, 345.2, 345.1, 393, 434, 427.9, 422, 38.3, 292.3, 352.2)

res <- t-test(before, after, paired = TRUE)

res

Output

data : before and after

$t = -20.883$ df = 9, p-value = 6.2×10^{-9}

alternative hypothesis: true difference in mean is not equal to 0. 95 % interval

-2.15 - 55.81 - 173.4219

Mean of differences

-194.49

~~(Q18)~~ ~~(Q8)~~ Explain Tukey's method and relate to R package TukeyHSD to understand the output.

→ Tukey's range test (also known as Tukey's test or Tukey method)

- The Tukey method applies simultaneously to the set of all pairwise comparisons
- The confidence coefficient for the statement, all the sample sizes are equal, exactly $1 - \alpha$. for unequal sample size

$$M = g \max(x) - \min(x) / NMS/n$$

Code

library(multcompView)

set.seed(1) treatment <- rep(c("A", "B", "C", "D", "E"), each = 20)

value <- c(sample(2:5, 20, replace = T), sample(6:10, 20, replace = F), sample(1:7, 20, replace = T))

model <- lm(data \$ value ~ data \$ treatment)

anova <- anova(model)

Tukey <- Tukey HSD(anova, data \$ treatment, cont.lvl = 0.95)

plot(Tukey, los = 1, col = brown)

Name - Alenraj Appu

ET29

Page No.
Date

ALLENRAJ

Output

B-A

C-A

D-A

E-P

L-B

N-S

I-O

D-C

E-C

F-D

-S

O

S

10

(Q9) Write a code to understand the critical region in R software of Normal distribution

→ The function that represents the distribution of many random variables as a symmetrical bell-shaped graph is known as normal distribution
 Critical region

The set of all the values of the test statistic that would cause us to reject the null hypothesis is called Critical region or rejection region. It indicates that if the value of the test statistic lies in the region.

If the test statistic fall into the region of non-rejection, you do not reject the null hypothesis.

If the test statistic is

Ex

$$\bar{x}_{\text{bar}} = 9.0$$

$$\mu_0 = 10$$

$$s_{\text{sig}} = 1.1$$

$$n = 30$$

$$z = (\bar{x}_{\text{bar}} - \mu_0) / (s_{\text{sig}} / \sqrt{n})$$

?

$$(1) -0.4979296$$

$$\alpha/\text{beta} = 0.5$$

$$\gamma = \alpha/\text{beta} = 0.5$$

$$\gamma = 2 \cdot \alpha/\text{beta}$$

(Q3) Display the types of Error in Hypothesis testing using R software

Type I error

It occurs when rejecting null hypothesis when true

```
typeI.test <- function(mu0, sigma, n, alpha, iteration = 100000){
```

```
  pvals <- rep(NA, iteration)
```

```
  for (i in 1:iteration){
```

```
    temporary.sample <- rnorm(n = n, mean = mu0, sd = sigma)
```

```
    temporary.mean <- mean(temporary.sample)
```

```
    temporary.Sd <- sd(temporary.sample)
```

```
    pvals[i] <- 1 - pt(((temporary.mean - mu0) / temporary.sd
```

```
      lsq rt(n)), df = n - 1)
```

```
  }
```

```
  return (mean(pvals < alpha))
```

```
}
```

```
typeI.test(<math>\mu_0=5</math>, sigma = 0.25, n = 100, alpha = 0.03)
```

```
output - 0.00472
```

Type 2 error

It occurs when failing to reject null hypothesis when it is actually false

The probability

```
typeII.test <- function(mu0, TRUEmu, sigma, n, alpha,
```

```
  iteration = 10000){
```

```
  pvals <- rep(NA, iteration)
```

```
  temporary.sample <- rnorm(n = n, mean = TRUEmu,
```

```
  sd = sigma)
```

(Q1) Plot the properties of discrete random variable and explain the same.

- 1) A discrete random variable has a countable number of possible values. The probability of each value of a discrete random variable is between 0 and 1, and the sum of all probabilities is equal to 1. A continuous random variable takes on all the values in some interval of numbers.

Code
library(discreteRV)
x <- RV(outcomes = 1:6, probs = 1/6)
X prob(x)
plot(X)

Sum (prob(x)*outcomes(x))

The easy way is to use the Expectation function in the discreteRV package:
 $E(X)$

Similarly, you can calculate variance by calculating the weighted average of the squared deviations $\sum p(x)(x - E(x))^2 \sum p(x)(x - E(x))^2$

Friday Appovi

Page No.
Date
PINRA

Sum ($\text{prob}(x) \cdot (\text{outcomes}(x) - \bar{x}(x))^2$)

Or you can use the variance function V , and standard deviation function, SD ;
 $V(x)$ & $\text{SD}(x)$

More dice.

What if you roll dice and add them?
You could treat this as the sum of two independent and identically distributed (iid) random variables. In the discrete sum package the `SOLID` function for sum of iid random variables:

`SOLID(x, n=2)`

Two dice

`<- SOLID(x, n=2)`

What if you roll 20 dice and add them up?

Twenty-dice <- `SOLID(x, n=20)` plot(twenty-dice)

(Q2) Plot the properties of continuous variable and explain the same.

→ The probability density function $f(x)$ of a continuous random variable x with a supporting set B

- $f(x)$ is positive everywhere in the supporting set B , that is $f(x) > 0$ for all x in B
- The area under the curve $f(x)$ in the support B is 1, i.e. $\int_{-\infty}^{\infty} f(x) dx = 1$
- If $f(x)$ is the P.d.f of x then the probability that x belongs to A (where A is some interval within the range), is given by the integral of $f(x)$

Code
 $x \leftarrow \text{seq}(\text{from} = -3, \text{to} = +3, \text{length.out} = 100)$
 $\text{dnorm}(x)$
 $\text{plot}(x, \text{dnorm}(x))$

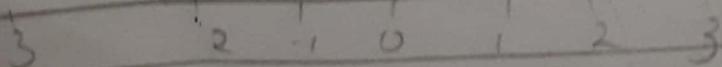
24

0.2

2.2

0.1

0.2



(Q3) A Plot a pie diagram using R software

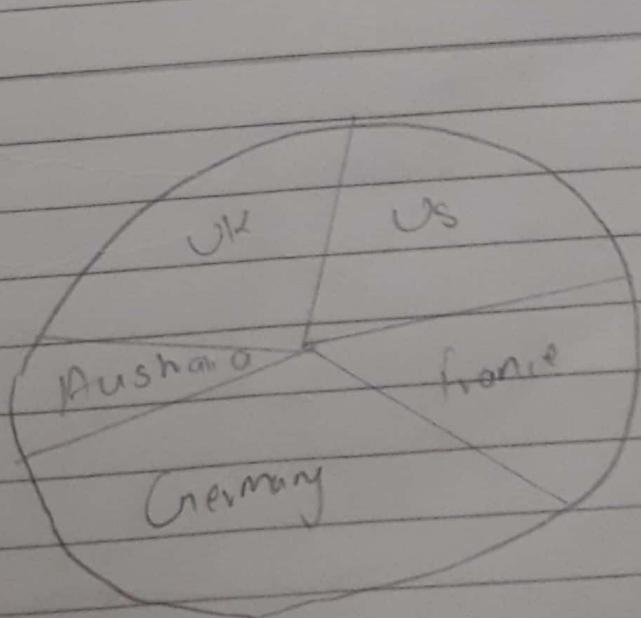
- 1) A pie chart is a representation of values as slices of a circle with different colors
- 2) In R pie chart is created using the `pie()` function which takes positive numbers as a vector input
- 3) Pie charts are not recommended in the R documentation, and their features are somewhat limited

Code:

Slices <- c(10, 12, 4, 16, 8)

lbls <- c("US", "UK", "Australia", "Germany", "France")
Pie (slices, labels = lbls, main = "Pie Chart of Countries")

Output



(Q8) Explain and plot the bar diagram using R software.

- A bar graph is a visual tool that uses bars to compare data among categories.
- 2) A bar graph may run horizontally or vertically. The important thing is,
- 3) Bar graph consists of two axes. On a vertical bar graph, the horizontal axis (X-axis) and on vertical y axis

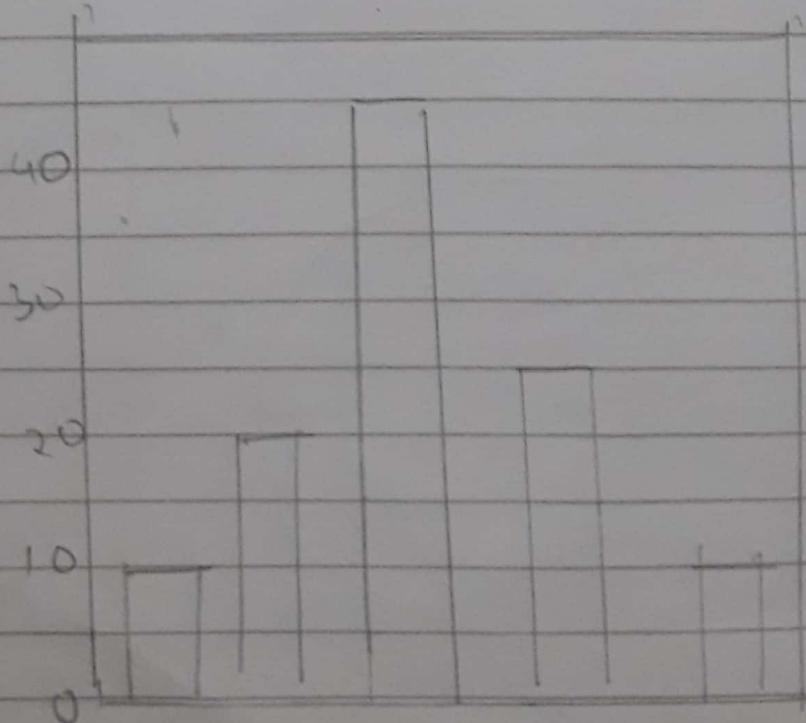
Code

H <- c(7, 12, 28, 3, 41)

png(file = "bar-chart.png")

bar.plot(H)

dev.off()



(ass) Calculate the mean and variance in R software for a data set and explain the application of both in day to day life.

Any numeric variable will have mean, median and mode.

Mean

It is calculated by taking the sum of the values and dividing with the number of values in data sets.

Variance (standard deviation) is the expectation of the squared deviation of the random variable from its mean.

2) Variance is much more used measure of variation.

Neeraj Appar F72a

Page No.

Date

ALFIRAO

Mean

Suppose there are frequencies of bus coming in a bus station

$$I.P = (13, 16, 20, 11, 18, 17, 15, 16)$$

$$\sum x C - c (13, 16, 20, 17, 18, 17, 15, 16)$$

\bar{x}

? mean C

16.5

Variance

The company produce random sample of 65 such high quality selected off of which 25% are defective. To detect 60 defective. Find variance and SD

$x_C \sim \text{Binomial} (\text{size} = 65, \text{prob} = 0.25)$

? $E(x)$

$$f(x) = 11.25$$

? $\text{Var}(x)$

$$\text{Var}(x) = 8.4375$$

? $s_d(x)$

$$s_d(x) = 2.904738$$

Plame - New 10
F. 12^a

Appri

Q56) Perform the vector operation in R software and Explain the same.

Vector is a basic data structure in R. It contains elements of the same type. i.e. data types can be logical, integer, double, character, complex or raw.

A A vector's type can be checked with the type() function

3) Creating a vector using : operator

```
7> x <- 1:7, x  
[1] 1 2 3 4 5 6
```

2) Creating a vector using seq() & function

```
7> seq(1, 2, by = 0.2)  
[1] 1.0 1.2 1.4 1.6 1.8 2.0
```

Using integer vector as index

Vector index in R starts from 1, unlike most programming language where it starts from 0.

```
7>  
[1] 0 2 4 6 8 0  
7> x[3]  
[1] 4
```

2) calculate the relationship between the coefficient "d" and Nernst-Neher's Apparatus

3) Using logical vector as index

```
? x <- c(TRUE, FALSE, FALSE, TRUE)  
? x[1:3]
```

```
? x[x < 0]  
[1] -3 -1
```

4) Using character vector as index
This type of indexing is useful when dealing with named vectors.

```
? x <- c(first = 3, second = 0, third = 9)  
? names(x)
```

```
? x["first"]  
[1] 3  
? x["second"]  
[1] 0
```

5) To modify a vector in R

We can modify a vector using assignment operator

```
? x  
[1] -3 -2 -1 0 1 2  
? x[1] <- 0;  
? x
```

Page No. 112
Date 10/12/2020

Name - Neeraj Appavu
I-12A

Page No.
Date
A-L-N-298

To Delete a vector

We can delete a vector by simply assigning a NULL to it

? 76
? [1] -3 -2 -1 0 12

? x <- NULL

? x
NULL

Q5] Plot a line of regression in R software and explain the same.

→ 1) A linear regression is a statistical model that analyzes the relationship between a response variable (X) and one or more variables (Y) and often interaction.

2) It is used to predict the value of an outcome variable X based on one or more input predictor.

3) It is the line that best fits the data

$$\text{formula } \hat{Y} = \beta_0 + \beta_1 X + E$$

Code

Ex → weight $\leftarrow c(70, 115, 50, 102, 59, 57, 63, 49, 50, 50)$

→ height $\leftarrow c(164, 164, 133, 175, 148, 164, 173, 155, 148, 175)$

→ plot (weight, height)

→ lm (height ~ weight)

Call:

lm(formula = height ~ weight)

(Coefficients:

(Intercept)

117.7507

weight

0.6400

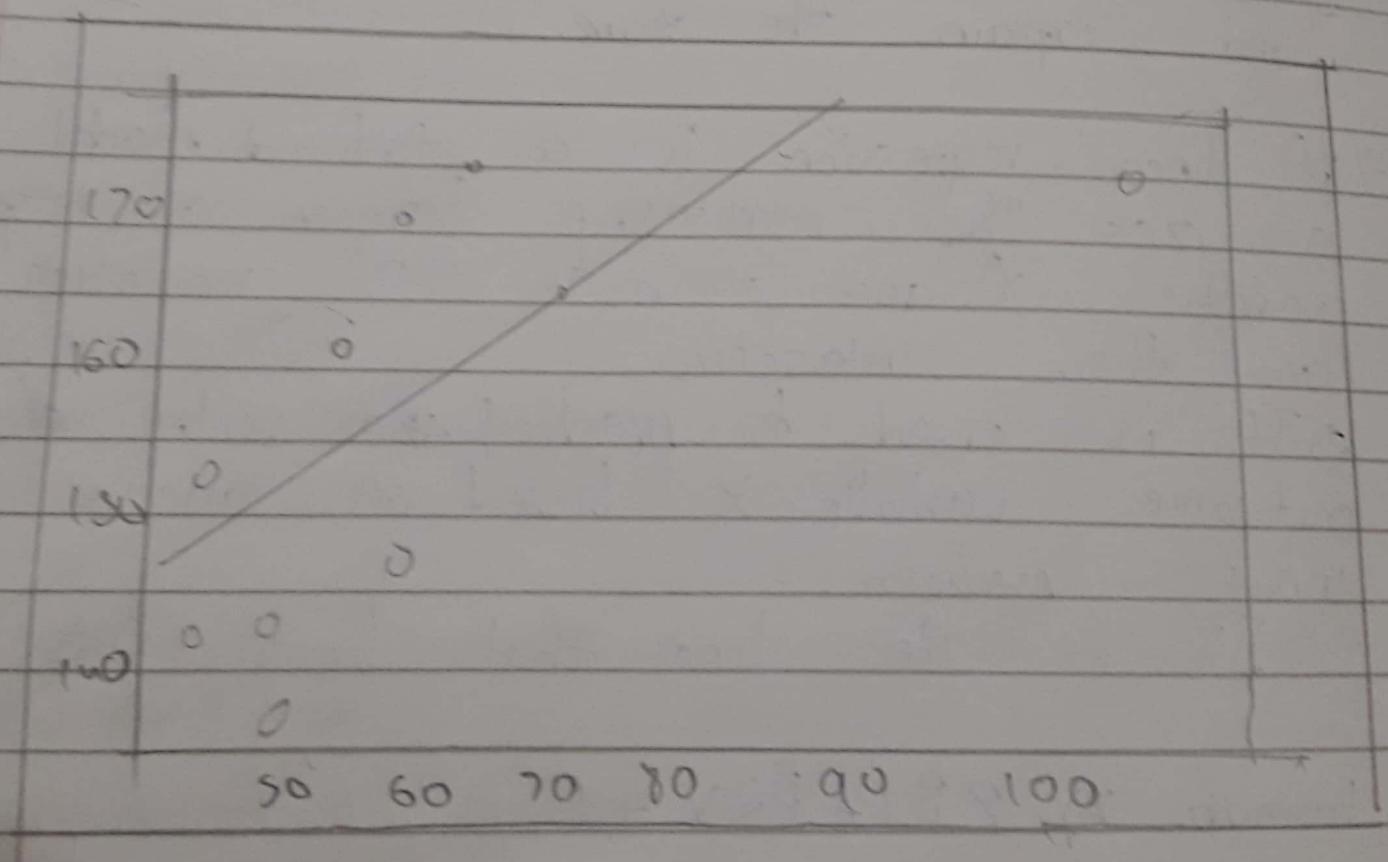
→ abline (117.7507, 0.6400)

F-129

Page No.
Date

A.L.NIRAO

Output.



- (Q8) Prepare the data set and plot the bayes theorem for statistical Analysis and explain Bayes theorem.
- 1) TD is mathematical formula for calculating conditional probability
 - 2) This theorem provides a way to revise or existing predictions or theories given new or additional evidence

$$P(A|B) = \frac{P(B|A) \cdot P(A)}{P(B)}$$

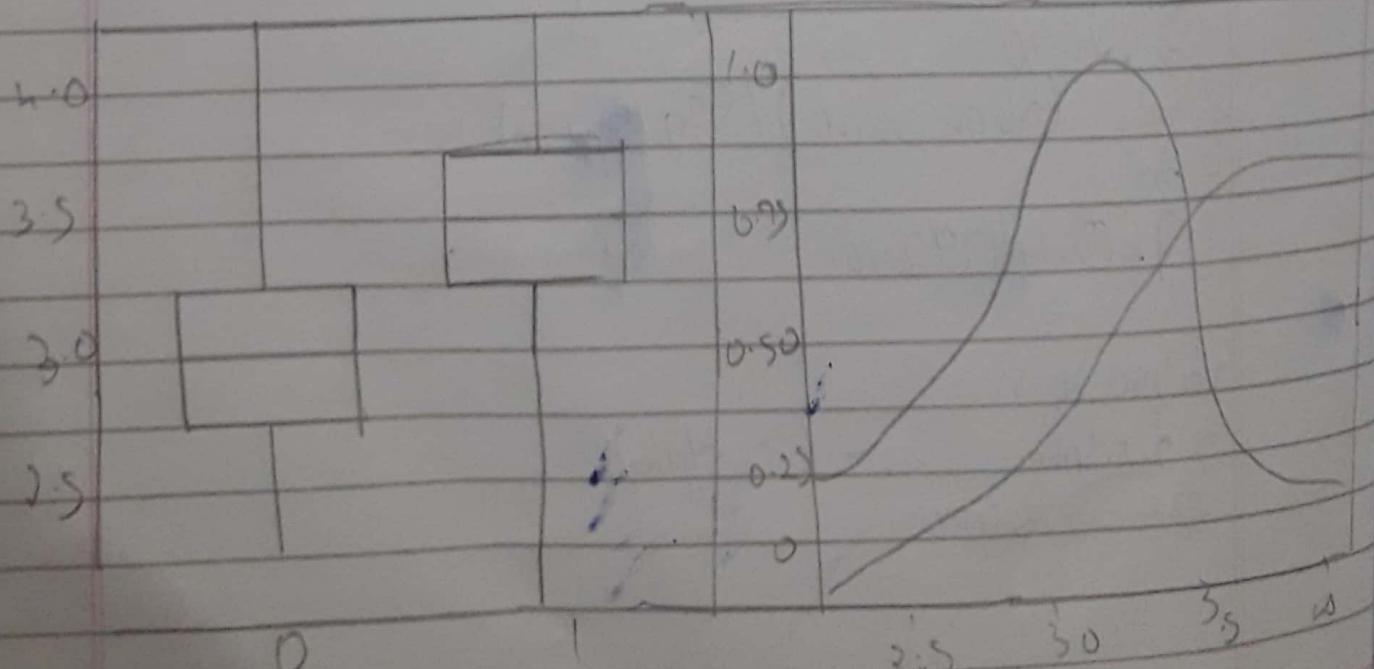
? pgm's panels (data (-1))

- data % > %

ggplot (aes (x=admit, y=gpa, fill=admit)) + geom_boxplot + ggtitle ("Box plot")

- data % > %

ggplot (aes (x=gpa, fill=admit)) + geom_density (color=black) + ggtitle ("Density")



(Q3) Plot the Decision Tree in R software only

```
install.packages("rpart")
install.packages("rpart.plot")
path <- "C:/Users/App/ESU"
fit <- read.csv(path)
tail(fit)
```

fit > tail fit

	Income	GymVisits	State	Hours	PayOrNot
q5	50	2	MP	5.2	No
q6	80	3	CA	6.2	No
q7	75	3	WA	7.4	Yes
q8	65	4	SD	6.0	No
q9	90	3	TY	7.6	Yes

fitness <- head(tail, n = 7)

library(rpart)

tree <- rpart(PayOrNot ~ Income + GymVisits + State, data = fitness)

tree

n = 100

node_id, split, n, loss, yval, yrob
(index), sum_n, node_id

1) node 00 28 No

2) Income < 95 71 9

3) Income >= 95 29 10

State: NY, 9, 1 NO

Roll no - 129

Page No.

Date

A.L.NRAO

library ("report.plot")

report.plot (here, extra=4)

(720.2)

~~Indonesia~~

Yes

34 66

State - NY

NO
87.13

NO
89.11

Yes
10.60

Q6. Karl Pearson's Coefficient of Correlation

→ The Karl Pearson's coefficient of correlation, a measure of strength of a linear association between two variables and used wherein the numerical expression is used to calculate the degree and direction of the relationship between linear related variables. The coefficient of correlations is denoted by "r"

formula

$$r = \frac{2(x - \bar{x})(y - \bar{y})}{\sqrt{\sum(x - \bar{x})^2} \cdot \sqrt{\sum(y - \bar{y})^2}}$$

where

\bar{x} = mean of x variable

\bar{y} = mean of y variable

Code

cor.test(trees \$Girth, trees \$Height)

→ cor.test(trees \$Girth, trees \$Height)

Person's product-moment correlation

Data trees \$Girth and trees \$Height

F = 3.32722, df = 29, p-value = 0.002751

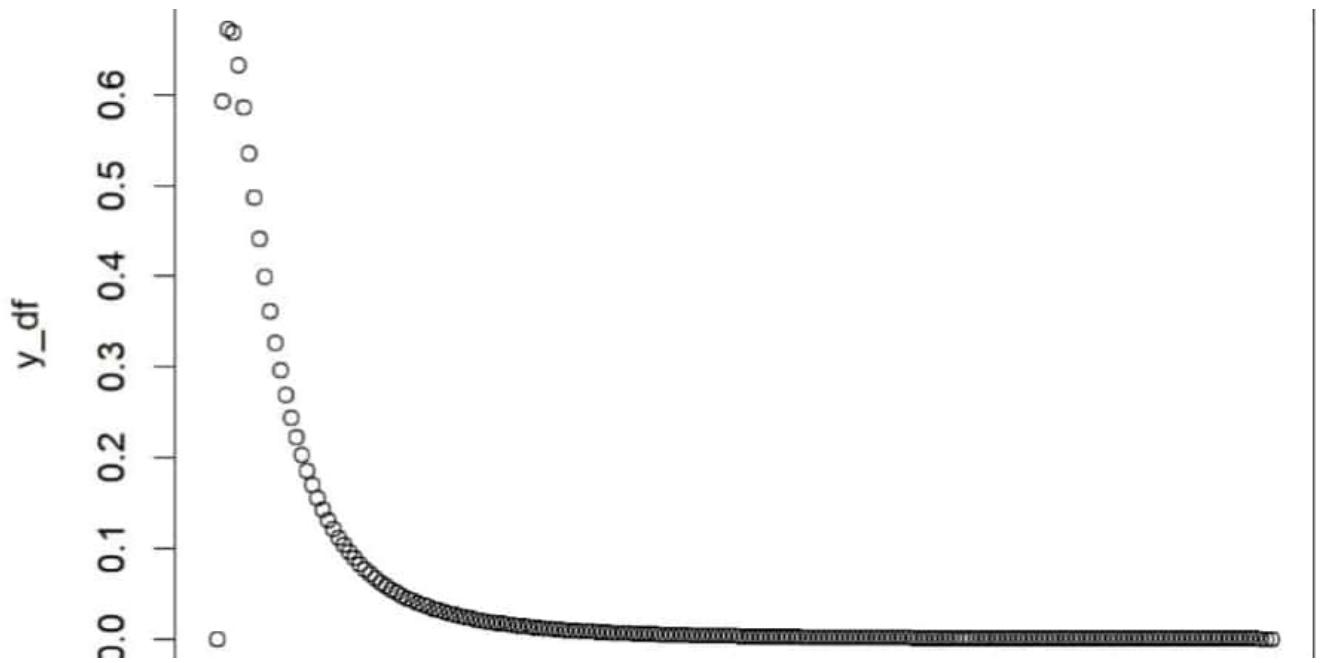
as percent (%)

alternative hypothesis: true correlation is not equal to 0.95 percent (confidence interval):

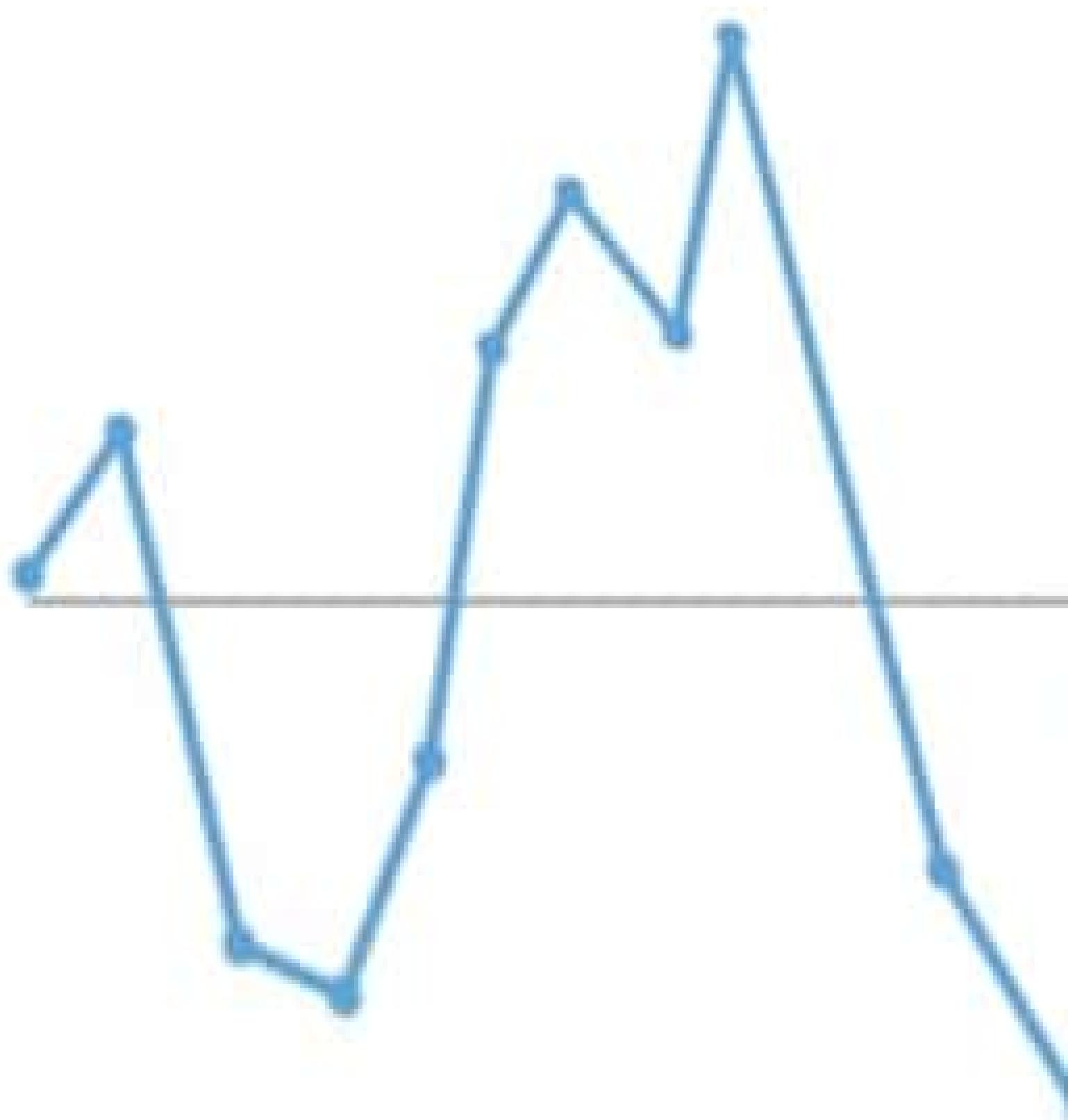
0.2021327 0.737853

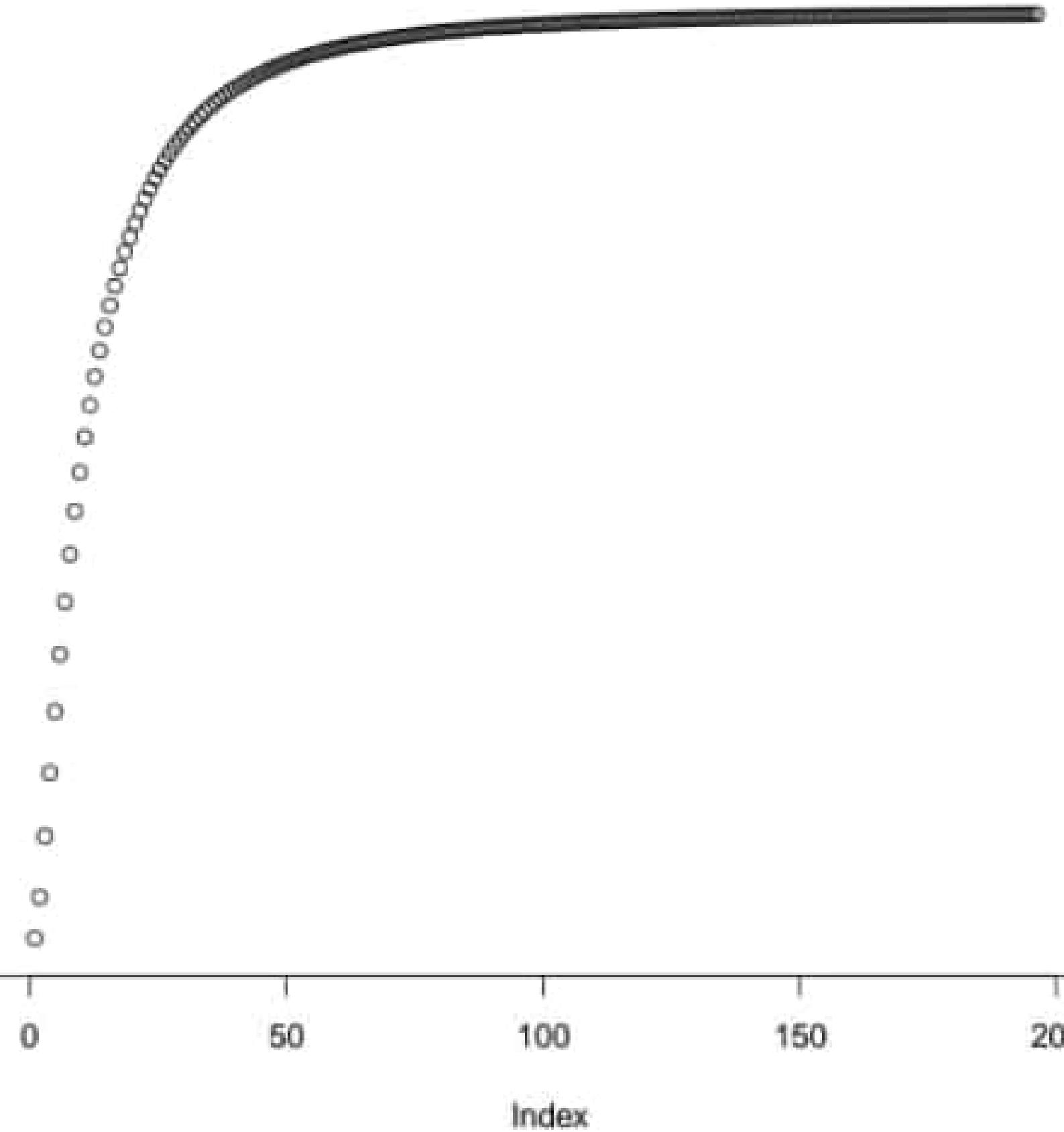
Sample estimate

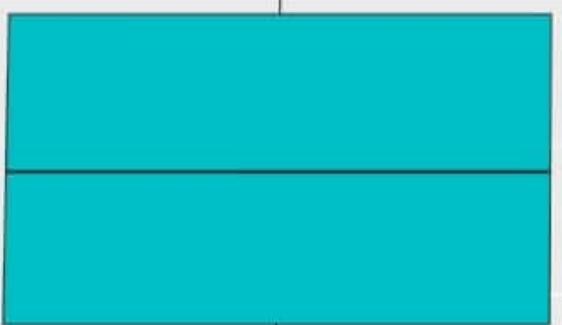
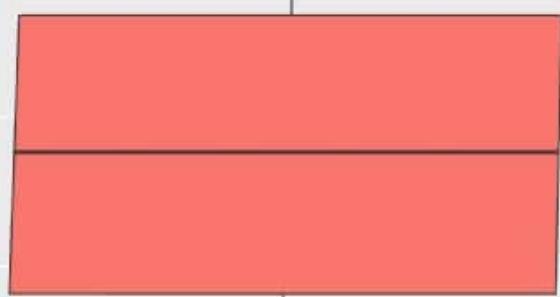
0.6192801

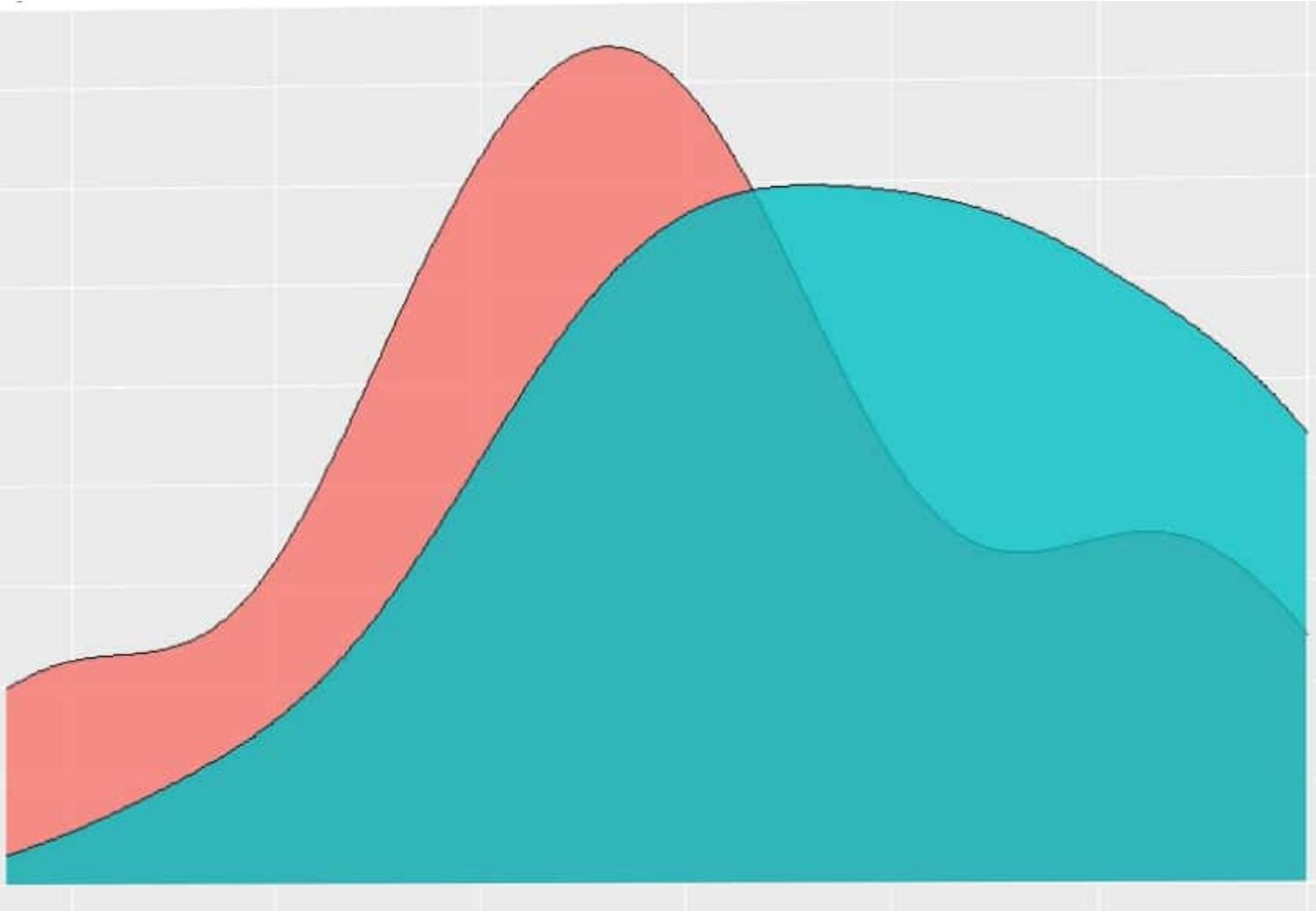


RUN Chart of y / n









Pie Chart of Countries

