

II - ASSIGNMENT

(Start Writing From Here)

- ① Discuss K-medoids clustering with an example.
- ② K-medoids is an unsupervised method with unlabelled data to be clustered. It is an improvised version of the K-Means algorithm mainly designed to deal with outlier data sensitivity. Compared to other partitioning algorithms, the algorithm is simple, fast and easy to implement. A medoid can be defined as a point in the cluster whose dissimilarities with all the other points in the cluster are minimum. The dissimilarity of the medoid (C_i) and object (P_i) is calculated by using $E = |P_i - C_i|$.

The cost in K-Medoids algorithm is given as

$$C = \sum_{C_i} \sum_{P_i \in C_i} |P_i - C_i|$$

EX:-

	X	Y
0	8	7
1	3	7
2	4	9
3	9	6
4	8	5
5	5	8
6	7	3
7	8	4
8	7	5
9	4	5

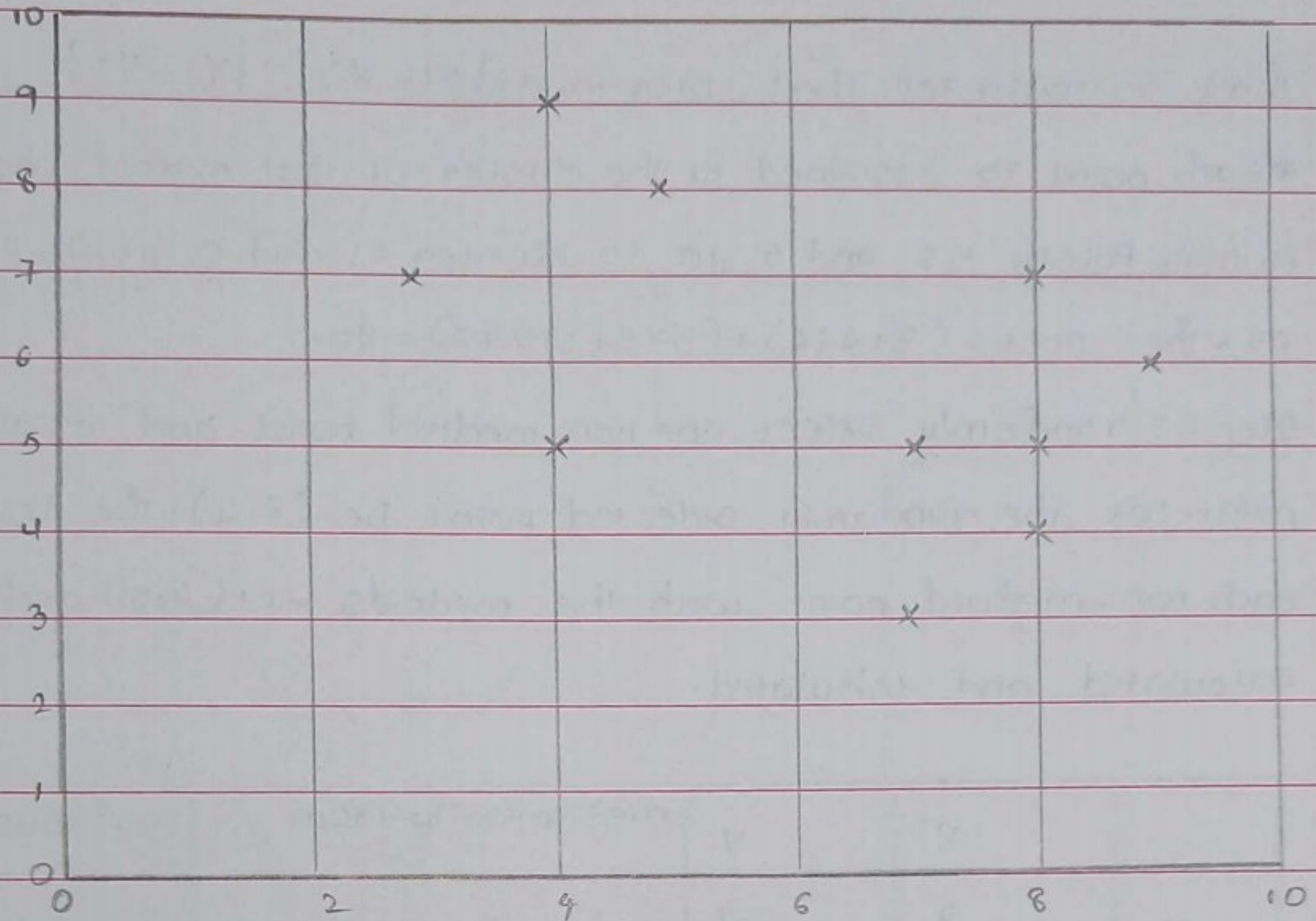
Let's consider the following example:

If a graph is drawn using the above data points, we obtain the following:

Step 1:- Let the randomly selected 2 medoids, so select $K=2$, and let $C_1 = (4, 5)$ and $C_2 = (8, 5)$ are ^{two} medoids.

Step 2:- Calculating cost. The dissimilarity of each non-medoid point with the medoids is calculated and tabulated.

Graph:



	x	y	dissimilarity from C ₁	dissimilarity from C ₂
0	8	7	6	2
1	3	7	3	7
2	4	9	4	8
3	9	6	6	2
4	8	5	—	—
5	5	8	4	6
6	7	3	5	3
7	8	4	5	1
8	7	5	3	1
9	4	5	—	—

Here we used Manhattan distance formula to calculate the distance matrices between medoid and non-medoid points.

That formula tell that, $\text{Distance} = |x_1 - x_2| + |y_1 - y_2|$.

* Each point is assigned to the cluster of that medoid whose dissimilarity is less. Points 1, 2 and 5 go to cluster c_1 and 0, 3, 6, 7, 8 go to cluster c_2 . The cost = $(3+4+4) + (3+1+1+2+2) = 20$.

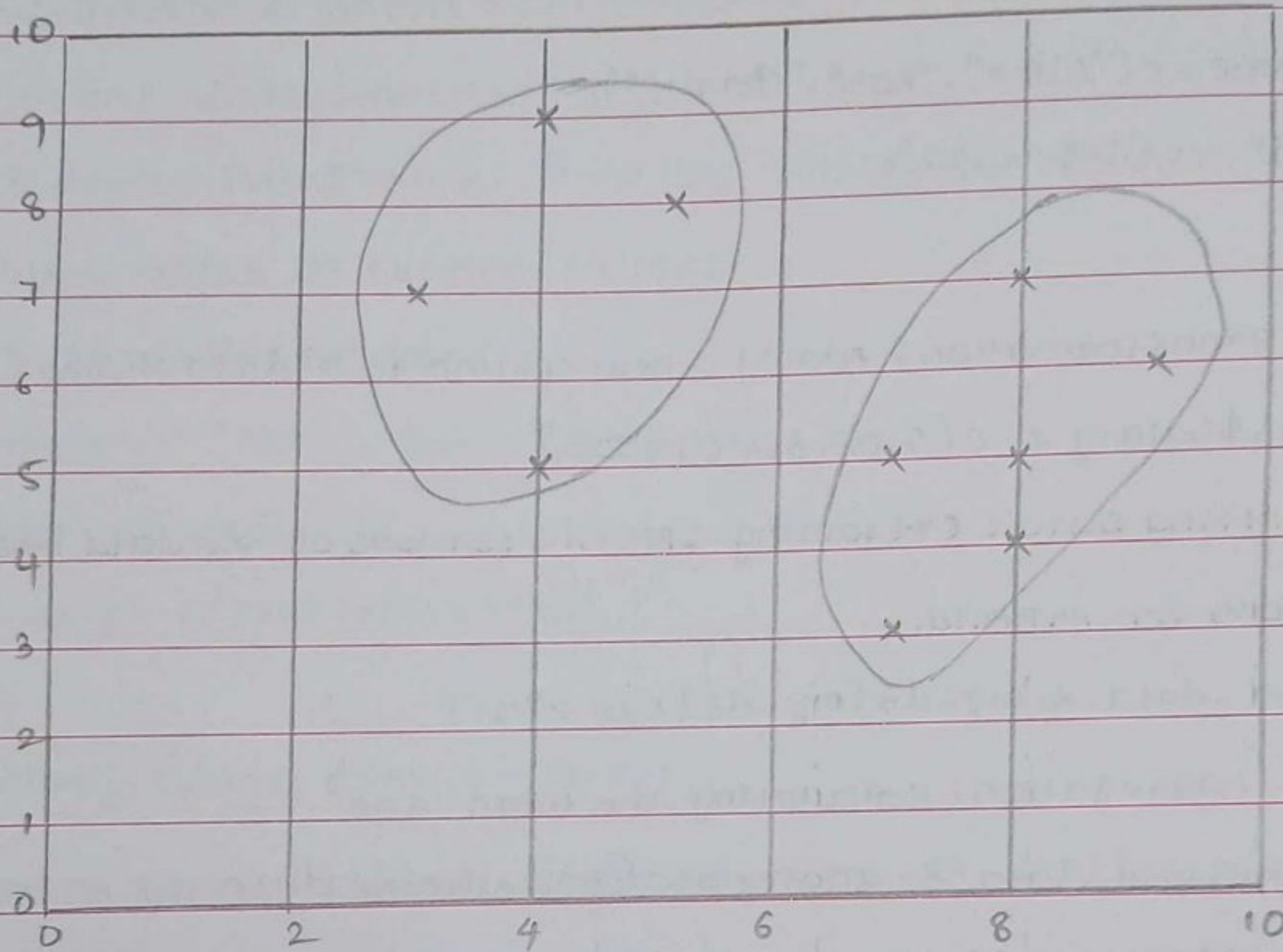
Step 3: randomly select one non-medoid point and recalculate the point. Let the randomly selected point be (8, 4). The dissimilarity of each non-medoid point with the medoids - $c_1(4, 5)$ and $c_2(8, 4)$ is calculated and tabulated.

	X	Y	Dissimilarity from c_1	Dissimilarity from c_2
0	8	7	6	3
1	3	7	3	8
2	4	9	4	9
3	9	6	6	3
4	8	5	4	1
5	5	8	4	7
6	7	3	5	2
7	8	4	—	—
8	7	5	3	2
9	4	5	—	—

Each point is assigned to that cluster whose dissimilarity is less. So, points (1, 2) and 5 go to cluster c_1 and 0, 3, 6, 7, 8 go to cluster c_2 .

The new cost = $(3+4+4) + (2+2+1+3+3) = 22$ swap cost = New cost - previous cost = $22 - 20$ and $2 > 0$. As the swap cost is not less than zero, we undo the swap. Hence (4, 5) and (8, 5) are final medoids.

The clustering would be in the following way. The time complexity is $O(K \times (n-K)^2)$.



② Explain about data management and indexing in R with suitable examples.

① Data management in R involves the organization, storage, manipulation and analysis of data using the R programming language. Here are some aspects of data management in R with examples:

① Data Import/Export: 'read.table' or 'read.csv' functions to import data from a csv file into a data frame.

```
my-data <- read.csv("path/to/your/file.csv")
```

② Data cleaning: Removing missing values from a data frame.

```
clean-data <- na.omit(my-df)
```


③ Data structures: creating a data frame

```

my-df <- data.frame(
  ID = c(1, 2, 3),
  Name = c("Alice", "Bob", "Charlie"),
  Age = c(25, 30, 22)
)
  
```

④ Data Transformation: Adding a new column to a data frame.

```
my-df$Salary <- c(5000, 6000, 4500)
```

⑤ Subsetting data: Extracting specific portions of the data based on conditions or criteria.

```
subset-data <- my-df[my-df$Age > 25, ]
```

⑥ Data Aggregation: calculating the mean age

```
aggregate-data <- aggregate(Age ~ Name, data = my-df, FUN = mean)
```

⑦ Data Merging: Merging two data frames

```
merged-data <- merge(df1, df2, by = "ID")
```

⑧ Data Visualisation: Visualising the data using plots and charts

```
plot(my-df$Age, my-df$Salary, main = "Age vs salary", xlab = "Age",
```

⑨ Data storage: Saving the data for future use.

```
write.csv(my-df, file = "path/to/save/data.csv", row.name = FALSE)
```

Indexing: Indexing is the process of accessing or extracting specific elements, rows or columns from data structures like vectors, matrices and data frames.

① Indexing vectors: Vectors in R can be indexed using square brackets '['.

```
my-vector <- c(10, 20, 30, 40, 50)
```

```
second-element <- my-vector[2]
```


② Indexing Matrices: Matrices in R are two-dimensional and you can use row and column indices to access elements.

```
my-matrix <- matrix(1:9, nrow=3)
```

```
element <- my-matrix[2,3]
```

③ Indexing Dataframes: These are similar to matrices. You can use column names or numeric indices.

```
my-df <- data.frame(
```

```
  Name = c("Alice", "Bob", "Charlie"),
```

```
  Age = c(25, 30, 22),
```

```
  Salary = c(5000, 6000, 4500)
```

```
)
```

```
salary-column <- my-df$Salary.
```

③ ~~Discuss~~ Discuss about different types of operators in R programming with example.

① In R programming there are different types of operators and each operator performs a different task.

Arithmetic Operators: These are the symbols which are used to represent arithmetic math (operators) operations.

+ → This operator is used to add two vectors in R. $a \leftarrow c(2, 3, 3, 4)$

```
a <- c(2, 3, 3, 4) b <- c(11, 5, 3)
```

```
print(a+b)
```

- → used to subtract two vectors. $a \leftarrow c(2, 3, 3, 4)$

```
b <- c(11, 5, 3)
```

```
print(a-b)
```


$*$ \rightarrow used to multiply two vectors

`print(a * b)`

$/$ \rightarrow used to divide the vector from another one.

`print(a / b)`

$\%$ \rightarrow used to find remainder of the first vector with second vector.

`print(a % b)`

$\% /$ \rightarrow used to find the division of first vector with second.

`print(a % / b)`

Relational operators: It is a symbol which defines some kind of relation between two entities.

$>$ \rightarrow will return TRUE when every element in the first vector is greater than the corresponding element of second vector.

$<$ \rightarrow will return TRUE when every element in the first vector is less than the corresponding element of second vector.

\leq \rightarrow will return TRUE when first vector element is less than or equal to second element.

\geq \rightarrow will return TRUE when first vector element is greater than or equal to the second element.

$==$ \rightarrow return TRUE when two elements are equal.

$!=$ \rightarrow return TRUE when two elements are not equal.

Logical operators: $\&$ \rightarrow takes the first element of both the vector and return TRUE if both the elements are TRUE.

$|$ \rightarrow takes the first element of both the vector and returns TRUE if one of them is TRUE.

Assignment operators:

$\angle - \text{or} = \text{or} \angle \angle - \Rightarrow$ These operators are known as left assignment operators.

$\rightarrow \text{or} \rightarrow \Rightarrow$ Known as right assignment operators.

Miscellaneous Operators:

$:$ \rightarrow used to create the series of numbers in sequence for a vector.

$v \angle - 1:8$

`print(v)`

o/p: [1] 1 2 3 4 5 6 7 8

$\%in\%$ \rightarrow used when we want to identify if an element belongs to a vector.

$\%*\%$ \rightarrow used to multiply a matrix with its transpose.

④ How to estimate the parameters of a model using maximum likelihood method?

④ Maximum likelihood estimation (MLE) is a statistical method used to estimate the parameters of a probability distribution that best describe a given dataset. The fundamental idea behind MLE is to find the values of the parameters that maximize the likelihood of the observed data, assuming that the data are generated by the specified distribution. A parameter is a numerical characteristic of a distribution. Mean (μ), variance (σ^2) as parameters, number of trials (n) & probability of success (p) as parameters. Gamma distributions have shape (k) and scale (θ) as parameters. Exponential distributions have the inverse mean (λ) as the parameter.

* Things aren't always that simple. Sometimes, you may encounter problems involving estimating parameters that do not have a simple one-to-one correspondance with common numerical characteristics. For instance,

If I give you the following distribution:

$$f_{\theta}(x) = \theta x^{-\theta-1}$$

The above equation shows the probability density function of a Pareto distribution with scale=1. It's not easy to estimate parameter θ of the distribution using simple estimators based because the numerical characteristics of the distribution vary as a function of the range of the parameter. For instance, the mean of the above distribution is expressed as follows:

$$\text{Mean} = \begin{cases} \infty & \text{if } \theta \leq 1 \\ \frac{\theta}{1-\theta} & \text{if } \theta > 1 \end{cases}$$

⑤ Explain how the F-value can be estimated in statistical learning.

① F-test is any test that utilises the F-distribution table to fulfill its purpose. It compares the ratio of the variances of two populations and determines if they are statistically similar or not. We can use this test when:

* The population is normally distributed.

* The samples are taken at random and are independent samples.

Formula:

$$F_{\text{calc}} = \frac{\sigma_1^2}{\sigma_2^2}$$

Steps involved:

① Use standard deviation (σ) and find variance (σ^2) of the data.

② Determine the null and alternative hypothesis.

$H_0 \rightarrow$ no difference in variances, $H_a \rightarrow$ difference in variances.

⑤ Find F_{calc} using eq-1.

⑥ Find degrees of freedom of the two samples.

⑦ Find F_{table} value using d_1 and d_2 obtained in step-4 from the F-distribution table. Take $\alpha = 0.05$

⑧ Interpret the results using F_{calc} and F_{table} .

$F_{calc} \leq F_{table} \rightarrow$ Can't reject null hypothesis.

$F_{calc} > F_{table} \rightarrow$ reject null hypothesis.

⑨ Illustrate the levels of measurement of data.

There are 4 levels of measurement.

① Nominal level: You can categorize your data by labelling them in mutually exclusive groups, but there is no order b/w the categories.

② Ordinal level: You can categorize and rank your data in an order, but you cannot say anything about the intervals between the rankings.

③ Interval level: You can categorize, rank and infer equal intervals between neighbouring data points, but there is no true zero point.

④ Ratio level: You can categorize, rank, and infer equal intervals between neighbouring data points, and there is a true zero point.