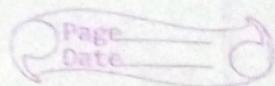


Assignment: 1



Q:1) What is Data science? What are its applications?

Data Science is a multidisciplinary field that involves using scientific methods, algorithms, processes, and systems to extract knowledge and insights from data.

- It combines elements of statistics, mathematics, computer science, and domain expertise to analyse and interpret data.

Applications of data science are vast and diverse, including:

1. Business Analytics:

Forecasting sales, customer segmentation, market analysis, and decision-making support.

2. HealthCare:

Predictive modeling for disease diagnosis, drug discovery, and personalized medicine.

3. Finance:

Fraud detection, risk assessment, algorithmic trading, and credit scoring.

4. Marketing:

Customer behaviour analysis, targeted advertising, and campaign optimization.

5. Social Media Analysis:

Sentiment analysis, network analysis, and recommendation Systems

6. Natural Language Processing:

Language translation, sentiment analysis and chatbots.

7. Internet of Things (IoT):

Analyzing sensor data for predictive maintenance and optimization.

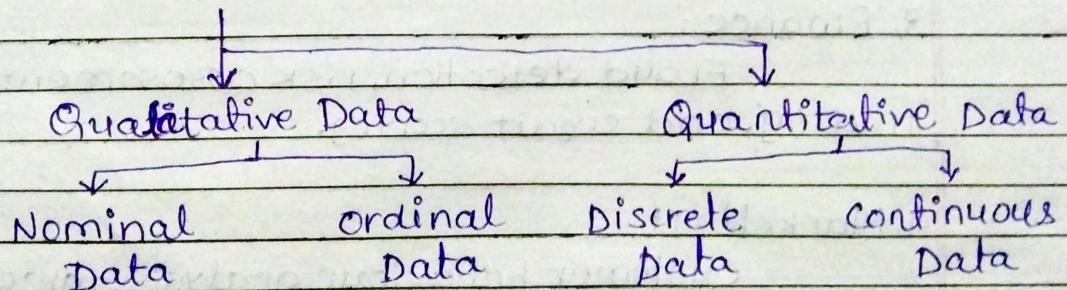
8. Environmental Science :

climate modeling, resource management, the environmental impact assessment.

Q:2) What are the different types of data in Data Science? Explain this?

→ Data is a collection of facts and figures. It is a set of characters used to collect & store information for specific purpose.

* Types of Data :-



(1) Qualitative Data :-

1. Nominal Data : (to label variables without any other)

→ Can't arrange in sequence (order)

Ex: Colour, gender, etc.

2. Ordinal Data :

- Having natural ordering (on the scale)
- Arrange in sequence (order)
- Feedback, Rank, Grade, etc.

(2) Quantitative Data :

1. Discrete Data :

- means separated
- It contains values under integer or whole number.
- It can't be broken into decimal or fraction.
Ex. NO. of students in class, etc.

2. Continuous Data :

- In form of fractional numbers.
- can be height of a person, height of object, etc.
- It represents information, that is divided into small parts.

Q:3)

calculate Descriptive statistics (mean, median, mode, standard deviation, variance, range) of the following list.

82, 93, 91, 69, 96, 61, 88, 58, 59, 100, 93, 71, 78, 98

Solⁿ

$$\text{Sum} = 1137, n = 14$$

$$\text{mean} = \frac{\text{sum}}{n} = \frac{1137}{14} = 81.2143 \approx 81$$

~~Data~~ Data in Ascending Order :

58, 59, 61, 69, 71, 78, 82, 88, 91, 93, 93, 96, 100

98,

$$\text{median} = \frac{\left(\frac{n}{2}\right)^{\text{th}} \text{ term} + \left(\frac{n}{2} + 1\right)^{\text{th}} \text{ term}}{2}$$

$$= \frac{82 + 88}{2}$$

$$= \frac{170}{2}$$

$$\rightarrow = 85$$

mode = 93

$$\text{Range} = 100 - 58$$

$$\rightarrow = 42$$

$$\text{Standard deviation (sample)} = \sqrt{\frac{\sum (x_i - \bar{x})^2}{n-1}}$$

$$\sum (x_i - \bar{x})^2 = (23)^2 + (22)^2 + (20)^2 + (12)^2 + (10)^2 + (3)^2 + (1)^2 + (7)^2 + (10)^2 + (12)^2 + (12)^2 + (15)^2 + (17)^2 + (19)^2$$

$$= 529 + 484 + 400 + 144 + 100 + 9 + 1 + 49 + 100 + 144 + 144 + 225 + 289$$

$$+ 361$$

$$\rightarrow = 2979$$

$$\text{standard deviation (sample)} = \sqrt{\frac{2979}{13}}$$

$$= \sqrt{229.154}$$

$$\rightarrow = 15.1378$$

$$V(x) = E(x^2) - (E(x))^2$$

~~$$E(x) = \frac{\sum x_i}{n} = 81 \quad (\sum x_i)$$~~

$$E(x^2) = \frac{\sum x_i^2}{n}$$

= $\frac{95319}{14}$
 = 6808.5

$$V(x) = (15.1378)^2$$

↳ = 229.154

$$V(x) = \sigma^2 E(x^2) - (E(x))^2$$

= $6808.5 - 6561$
 = 247.5

Q:4) Calculate sample covariance of given set of numbers.

$$x = 2.1, 2.5, 4.0, 3.6 \quad y = 8, 12, 14, 10$$

x	y	$x - \bar{x}$	$y - \bar{y}$	$(x - \bar{x})(y - \bar{y})$
2.1	8	-0.95	-3	2.85
2.5	12	-0.55	1	-0.55
4.0	14	+0.95	3	2.85
3.6	10	+0.55	-1	-0.55
				4.6

$$\bar{x} = 2.2305$$

$$\bar{y} = 11$$

$$\text{cov}(x, y) = \frac{\sum (x - \bar{x})(y - \bar{y})}{n}$$

= $\frac{4.6}{4}$
 = 1.15

Q:5) What is Sampling? what are the different Sampling techniques?

Ans.

Sampling is the process of selecting a subset of a population to study. This is done in order to make inferences about the population as a whole.

There are lot of sampling techniques, which are grouped into two categories as:

- Probability Sampling
- Non-Probability Sampling

* Non Probability Sampling

- This sampling technique uses randomization to make sure that every element of the population gets an equal chance to be part of the selected sample. It's alternatively known as Random Sampling.

1) Simple Random Sampling:

A method where every individual or item in the population has an equal chance of being selected for the sample, without any specific criteria or grouping.

2) Stratified Sampling:

Divides the population into mutually exclusive subgroups (strata) and then randomly selects samples from each stratum, ensuring representation of various characteristics within the population.

3) Systematic Sampling:

Selects samples at fixed intervals from the population, often using a random starting point, providing an evenly spread representation.

4) Cluster Sampling:

Divides the population into clusters, randomly selects some clusters, and then includes all the individuals within the selected clusters in sample.

5) Multi-Stage Sampling:

Involves a combination of various sampling techniques, using multiple stages of selection, suitable for large and complex populations.

* Non-Probability Sampling:

- It does not rely on randomization. This technique is more reliant on the researcher's ability to select elements for a sample. This type of Sampling is also known as non-random sampling.

1) Convenience Sampling:

Involves selecting individuals who are readily available and easily accessible, leading to potential biases due to the non-random nature of selection.

2) Purposive Sampling:

Also known as judgemental or selective sampling, it involves deliberately choosing specific individuals who possess particular characteristics or expertise relevant to the study.

3) Quota Sampling :

Selects Participants to meet pre-determined proportions or quotas of for certain characteristics, often used when random sampling is difficult.

4) Referral/snowball Sampling :

Relyes on initial participants to refer or recruit others from their network, commonly used in studies involving hard-to-reach populations.

Q: 6) Explain Hypothesis Testing. Explain any 2 types of testing.

Hypothesis testing is a method in inferential statistics used to make assumptions about a full population based on a representative sample since observing the entire population is often impossible. It involves establishing hypotheses, such as a difference between two groups or a correlation between variable in the population. The Alternative Hypothesis (H_A) represents the thesis to be proven, while the Null Hypothesis (H_0) states that nothing new is happening in the population.

The purpose of hypothesis testing is to determine if the Null Hypothesis can be rejected or not. Rejecting the Null Hypothesis does not necessarily prove the Alternative Hypothesis true, but it may lead to acceptance of the Alternative Hypothesis in some cases.

→ Two Types of Hypothesis testing are:

1) T-test:

- The T-test is used when the population standard deviation is unknown, and sample size is relatively small.
- The test statistics for that t-test is calculated as the difference between the sample mean and the population mean, divided by the standard deviation ~~error~~ of the mean.

$$t = \frac{\bar{x} - \mu}{\sigma/\sqrt{n}}$$

σ : standard deviation

μ : Population mean

\bar{x} : Sample's mean

n : sample's size.

2) Z-Test

- The z-test is used when the population standard deviation is known, or the sample size is large.

- It is often employed when comparing the means of two groups or mean when comparing the mean of a sample to a known population mean.

$$z = \frac{\bar{x} - \mu_0}{\sigma/\sqrt{n}}$$

$H_0: \mu = \mu_0$ or $H_A: \mu \neq \mu_0$

μ_0 : Population mean

σ : standard deviation

n : Sample size

\bar{x} : sample mean

Q:7)

Difference between Classification and clustering.

Classification

clustering

- It is used for supervised Learning.
 - It is a process of classifying the input instances based on their corresponding class labels.
 - It has labels so there is need of training and testing dataset for verifying the model created.
 - more complex as compared to clustering
 - Logistic Regression, Naive Bayes classifier, SUPPORT Vector machines etc.
 - Produces class labels or category predictions for new, unseen data based on the learned model.
- It is used for unsupervised Learning.
 - grouping the instances based on their similarity without the help of class labels.
 - There is no need of training and testing dataset.
 - less complex as compared to classification.
 - k-means clustering algo, Fuzzy - Cmeans clustering algo, Gaussian (EM) clustering algo.
 - Outputs clusters or groups of data instances that are similar to each other based on their features.

Q:8) Briefly discuss K-means clustering algorithm with its pros and cons.

Ans: K-Means is a popular clustering algorithm used to partition data into K clusters based on their similarity.

Algorithm:

1. choose the number of clusters K.
2. Randomly initialize K cluster centroids.
3. Assign each data point to the nearest centroid, forming K clusters.
4. calculate the mean of each cluster to obtain new centroid.
5. Repeat steps 3 and 4 until convergence (centroid do not change significantly) or a maximum number of iterations is reached.

Pros :

- Simple and easy to implement.
- computationally efficient, making it suitable for large databases-sets.
- Works well when clusters are well-defined and compact.
- converges to local optimum, which often provides reasonable results.

Cons :

- Requires the number of clusters K, to be specified beforehand, which may not always be known in advance.
- May produce suboptimal results when dealing with non-linearly separable or overlapping clusters.
- Does not work well with clusters of varying sizes or irregular shapes.

Q9)

Illustrate the steps of ANN algorithm with proper example.

→

To illustrate the steps of the Average Nearest Neighbor (ANN) algorithm we consider Business Analyst data. Analyst use ANN algorithm to perform classification analysis of datasets in his organization database. He compares employees based on the following five features:

- 1) Age
- 2) Number of children
- 3) Annual income
- 4) Seniority
- 5) Eligibility to Retire

Step 1: Input the spatial data set.

The database contains the 5 dimensional tuples representing the employees features.

Name	Age	childrens	Annual income	Seniority	Eligible to Retire
Mike	34	1	\$120,000	9	0
Liz	42	0	\$90000	5	0
JIN	22	0	\$60000	2	0
Mary	53	3	\$1,80,000	30	1

Step 2: calculate Difference Between Each 2 Employees.

1. Mike & Liz

$$\begin{array}{r}
 34 \quad 1 \quad 120000 \quad 9 \\
 - 42 \quad 0 \quad 90000 \quad 5 \\
 \hline
 -8 \quad 1 \quad 30000 \quad 4
 \end{array}$$

∴ calculate Difference
no need of Negative number

Same for all pair.

2. Mike & Jin:

$$(12, 1, 60000, 7)$$

3. Mike & Mary:

$$(19, 2, 60000, 21)$$

4. Liz & Jin:

$$(20, 0, 30000, 3)$$

5. Liz & Mary:

$$(11, 3, 90000, 25)$$

6. Jin & Mary:

$$(31, 3, 120000, 28)$$

Step 3: Finding Average of each pair of features values.

$$\text{Mike \& Liz} = 7503.25$$

$$\text{Mike \& Jin} = 15005$$

$$\text{Mike \& Mary} = 15010.5$$

$$\text{Liz \& Jin} = 7505.75$$

$$\text{Liz \& Mary} = 22509.75$$

$$\text{Jin \& Mary} = 30015.5$$

Step 4: Take Pairs that have minimum Average distance. compare that pair with other that is near to it.

- Mike & Liz and Liz & Jin

$$\left. \begin{array}{l} \{ 7503.25 \} \\ \{ 7505.75 \} \end{array} \right\}$$

so, we can group of Mike, Liz & Jin which will class-0.

- Now, Mary is

- So, Mary is now in class-1.

Q:10) What is confusion Matrix? Explain in detail all its measures.

A confusion matrix is a tabular representation of Prediction outcomes of any binary classifier, which is used to describe the performance of the classification model on set of data test data when true values are known.

The confusion matrix is simple to implement, but the terminologies used in this matrix might be confusing for beginners.

		Actual values	
		Positive	Negative
Predicted values	Positive	TP	FP
	Negative	FN	TN

- True Positive: when the actual value is positive and predicted is also Positive.
- True Negative: when the actual value is Negative and prediction is also Negative.
- False Positive: when the actual is negative but Prediction is Positive. Also Known as the Type 1 Error.
- False Negative: when the actual is positive but the Prediction is Negative. Also Known as the Type 2 error.

B

1. Accuracy :

The accuracy metric is one of the simplest classification metrics to implement, and it can be determined as the number of correct Predictions to the total number of Predictions.

$$\text{Accuracy} = \frac{\text{Number of Correct Predictions}}{\text{Total number of Predictions}}$$

2. Precision : ~~TP~~

It can be defined as the number of correct outputs provided by the model or out of all positive classes that have predicted correctly by the model, how many of them were actually true. It can be calculated using the below formula.

$$\text{precision} = \frac{\text{TP}}{\text{TP} + \text{FP}}$$

3. Recall or Sensitivity:

It is defined as the out of total positive classes, how our model predicted correctly. The recall must be as high as possible.

$$\text{Recall} = \frac{\text{TP}}{\text{TP} + \text{FN}}$$

4. F - Scores

F1 score can be calculated as the harmonic mean of both precision and Recall, assigning equal weight to each of them.

$$\text{F1 score} = 2 * \frac{\text{precision} * \text{recall}}{\text{precision} + \text{recall}}$$