

Q.1) Attempt any eight of the following.

a) List any 2 applications of Data Science.

→ i) Gaming world

ii) Health care sector

b) What is outlier?

→ The outlier is an observation point that is distant from other outlier observation.

c) What is missing value?

→ Some values in the data may not be filled up for various reasons & hence are considered missing values.

d) Define variance.

→ Variance is the measure of dispersion that is related to the standard deviation. It is calculated by finding the square of the standard deviation of given data distribution.

e) What is data transformation?

→ Data transformation is the process of converting, structuring data into a usable format that can be analyzed support decision making process to the growth an organization.

f) What is Nominal attribute?

→ Nominal means relating to name's the value of nominal attribute are symbol or names of things.

g) What is one hot coding.

→ One hot coding is one method of converting data to prepare it for an algorithm & get a better prediction.

h) What is the use of Bubble plot?

→ A bubble plot is scatter plot where a third dimension is added. The value of an additional numeric variable is represented through the size of the dots.

i) Define Data visualization

→ Data visualization is the presentation of data in graphical format. Data visualization is generic term used which describes any attempt to help understanding of data by providing visual representation.

j) Define standard deviation.

→ A standard deviation (or  $\sigma$ ) is a measure of how dispersed the data is in relation to the mean.

Q. 2) Attempt any four the following

a) Differentiate structure & unstructured data

→ i) Structure Data :- i) Structure data has name suggest type data is well organized. ii) Structured data is data that depends on a data model & resides in a fixed field within a record.

ii) Unstructure Data :-

i) Unstructured data is data that is not organized in a pre-defined manner does not have pre-defined data model.

ii) Unstructured data has internal structure but not structured via pre-defined data models or schema.

b) What is inferential statistics.

→ In inferential statistics, we make an inference from a sample about the population. The main aim of inferential statistics is to draw some conclusions from the sample, generalize them for the population data. Statistical inference mainly deals with two different kinds of problems: hypothesis testing & estimation of parameter values.

c) What do you mean by data preprocessing?

→ It is the task of transforming raw data to be ready to be fed into an algorithm. It is a time-consuming yet important step that cannot be avoided for the accuracy of result in data analysis.

d) Define Data discretization?

→ i) Data discretization is the process of converting continuous data into an a set of discrete inference or categories. ii) This technique is used to for data reduction, data simplification, or to make the data suitable format for analysis & it is typically used for large datasets.

e) What is visual encoding.

→ i) The visual encoding is the way in which data is mapped into visual structure, upon which we build the images on the screen.

ii) Encoding in data visualization means transforming the data into a visual element on a chart or map through position, shape, size, symbol & colour.

Q.3) Attempt any two of the following (d)  
pt

a) Explain outlier detection methods in brief.

→ The outlier detection methods can be divided into supervised method, unsupervised method & semi-supervised method.

i) Supervised method:

a) Supervised method model data normality & abnormality.

b) Domain experts examining of table a sample of underlying data.

ii) Unsupervised data:

a) In same application scenarios labelled as "normal" or "outlier" are not available. Thus an unsupervised learning method has to be used.

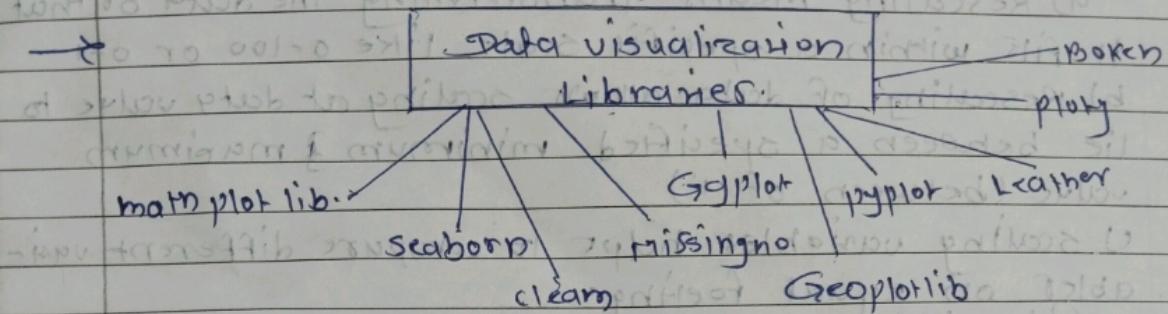
b) Unsupervised outlier detection methods like an implicit assumption, the normal objects are somewhat "clustered".

iii) Semi-supervised method:

a) Semi-supervised outlier detection method were developed to take such scenarios.

b) Building a model for outlier based on only a few labeled outlier. If outlier is unlikely to be effective.

b) Write different data visualization libraries in python.



Q.4) Attempt any two of the following:

a) Explain 3V's of Data science.

→ i) The 3V's are (volume, velocity, variety)

ii) Due to the expansion of data of the turn of the 21st century coined by the so-called 3V's of data science which are volume, velocity, variety.

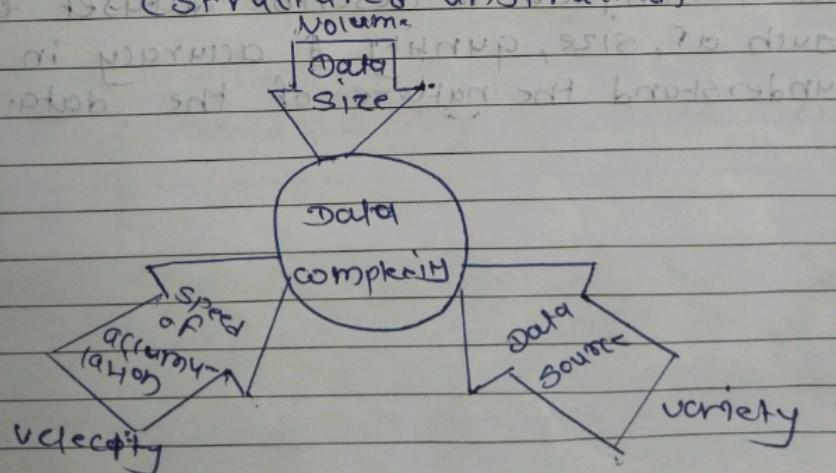
iii) Volume refers to the increasing size of data, velocity the speed at which data is acquired, variety the diverse types of data that are available.

iv) The 3V's are explained below:-

a) Velocity :- The speed at which data is accumulated.

b) volume :- The size of the scope of the data.

c) variety :- The massive array of data of types



b) Explain any two data transformation technique in details.

→ i) Rescaling :-

a) Rescaling means transforming the data so that it fits within a specific scale, like 0-100 or 0-1.

b) Rescaling of data allows scaling of data value to lie between a specified minimum & maximum value (between 0 & 1).

c) Scaling variables help to compare different variables on equal footing.

ii) Binarizing :-

a) It is the process of converting data to either 0 or 1 based on a threshold value.

b) All data values above the threshold value of one are marked where as all the data values equal to 0 or below the threshold value are marked as 0.

Q.5. Attempt any one of the following

a) Explain Exploratory data analysis (EDA) in details.

→ i) Exploratory Data Analysis is an analysis approach that identifies general patterns in the data.

ii) These patterns include outliers & feature of the data that might be unexpected.

iii) EDA is an important first step in any data science project.

iv) EDA is a crucial initial step in data science.

v) Data Analysis uses data visualization & statistical techniques to describe dataset characterization such as, size, quantity & accuracy in order to better understand the nature of the data.