

Q.1) Attempt any eight of the following.

1) What is data science?

→ Data science is deep study of massive amount of data, which involves extracting meaning full insights from raw, structured & unstructured data.

2) Define Data source.

→ A data source is simply the source of the data. A data source in data science is the initial location where data that is being used come from.

3) What is missing value?

→ Some values in the data may not be filled up for various reason is called as missing value.

4) List the visualization libraries in python.

→ i) matplotlib ii) feather iii) Geoplotlib
iv) seaborn v) plotly vi) Gleam
vii) Ggplot viii) Bokeh
ix) missing viii) pygal

5) List application of data science.

→ i) image Recognition Speech recognition
ii) Gaming world
iii) Internet search
iv) Transport
v) Healthcare
vi) Recommendation system
vii) Risk detection

6) What is Data transformation?

→ Data transformation is the process of converting raw data into a format that would be more suitable for data analysis.

7) Define Data cleaning

→ Data cleaning is the process of correcting or removing incorrect incomplete or duplicate data from a dataset.

8) Define Standard Deviation?

→ Standard deviation is a measure of how dispersed the data is in relation to mean.

Q.4 Attempt any four of the following

1) List the tools for data scientist

→ The tools are following :-
 i) Pytorch ii) Numpy iii) Pandas iv) SQL v) Excel
 vi) Tableau vii) Tensor Flow viii) Apache spark

2) What is Data cube?

→ i) A Data cube is a multi-dimensional array of values. A data cube generally used to easily interpret data.

ii) Data cube is especially useful when representing data together with dimensions as certain measures of business requirements.

3) What are the purpose of Data visualization?

→ i) Data visualization is the presentation of data in graphical format.

ii) Data visualization is a generic term is used which describes any attempt to help understanding of data by providing visual representation.

iii) It is the very important part of the data analysis. We can use it to explore our data.

Q.4) Define statistical data analysis

- i) Statistical data analysis is the collection & interpretation of data in order to uncover pattern & trends.
- ii) It is component of data analytics.
- iii) It is the science of collecting exploring & presenting large amount of data to discover underlying patterns & trends.

Q.3.) Attempt any 2 of the following

i) What are the measures of central tendency?

→ Explain any 2 of them in brief.

→ i) One of the simplest & yet important measure of statistical analysis is to find one such value that describe the characteristic of the entire huge set of data.

ii) A measure of central tendency is a summary statistic that represent the point or typical value of dataset.

iii) There are mainly 3 ways to measure central tendency, mean, median, mode as describe below

i) Mean :-

i) It is the popular & widely used measure of representing the entire data by one value.

ii) Mean is the ratio of the sum of all observations in the data to the total number of observations.

2) Median :-

i) The median of a distribution with a discrete random variable depends on whether the number of terms in distribution is even or odd.

Q.4) Attempt any 2 OF the following

- i) Explain different data format in brief
→ i) In data science the data appears in different size of shapes it can be numerical data, text, audio, video or a few other types of data.
ii) some data formats in data science are explained

a) Integers :- Integers data types may be different sizes & may not be allowed to contain negative value.

b) Floats :- A floating point has decimal point even if that decimal point value is 0.

c) Text Data :- Text data type is known as string in python, or object in pandas, strings can contain numbers & characters.

d) Dense Numerical Arrays :- Most image files or sound files consist mostly of dense arrays of numbers, packed adjacent to each other in memory.

e) Compressed or Archived Data :- Many data files when stored in a particular format take up a lot more space compared to the file in question logically needs less space.

f) CSV format :- CSV stands for comma separated value CSV files are the commonly used data format for data science.

Q3) What is venn diagram? How to create it? Explain with example.

i) A venn diagram is also called primary diagram set diagrams or logic diagrams is a diagram that shows all possible logical relationships between a finite collection of different sets.

ii) Each set is represented by circle, the circle size sometimes represents the importance of the group but not always.

iii) A venn diagram that visually displays all the possible logical relationships between collection of sets.

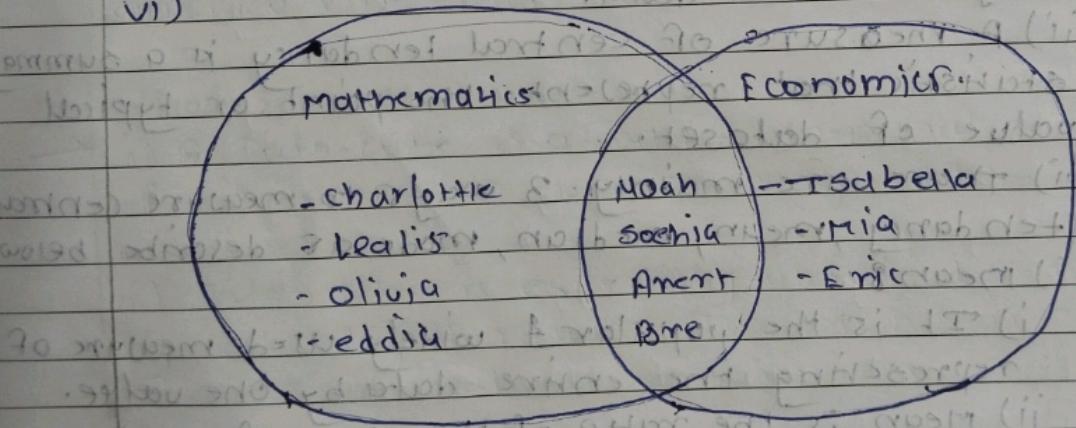
iv) A venn diagram is used a) visually organize

b) compare contrast.

c) find correlation & predict probability.

v) The simplest way to create a venn diagram is by drawing two overlapping circles.

vi)



- 2) Write details notes on Basic data visualization tools? Q.5) AT
→ i) Data visualization tools such as Histogram
Bar chart, scatter plot, line chart area.
plots, pie charts, donut chart if so ? → i) What
ii) Histogram :- A histogram is a graphical
display of data using base of different height
measured in cm.
iii) Bar charts :- Bar chart is one used for
comparing the quantities of different categories
or groups.
iv) Line plot :- The line plot is a two-dimensional
plotting of values connect to their order.
v) Scatter plot :- A scatter plot is a 2-D chart
showing the comparison of two variables
scattered across 2 axes.
vi) Pie chart :- A pie chart shows the proportion
or percentage of data elements in a circular
form.
vii) Area plot :- Area charts are used to plot
data trends over a while to show how a value
is changing.
viii) Donut chart :- A donut chart is an
extension of a pie chart.

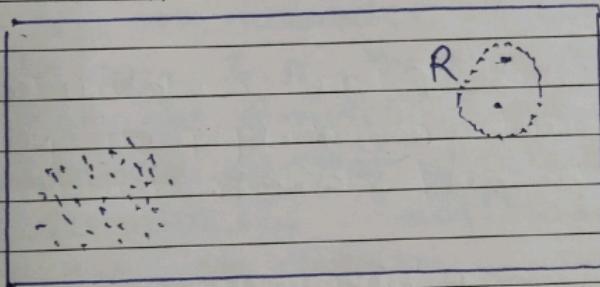
Ques. Q.5) Attempt any one of the following.

i) What is Outlier & state types of outliers?

→ i) Outliers are a very important aspect of data analysis.

ii) An outlier may indicate an experimental error or it may be due to variability in the measurement.

iii) Global Outlier: If an individual data point can be considered anomalous with respect to the rest of the data, then the donut is a point outlier.



ii) Contextual outlier: If an individual data instance is anomalous in a specific context or condition then it is termed as contextual outlier.

iii) Collective outlier: If a collective of data points is anomalous with respect to the entire data set it is termed as a collective outlier.

