

Q. 1.) Attempt any eight of following.

(a) What do you mean by primary data?  
→ Primary data is data that never collected before & can be generated in variety of ways such as conducting interview schedules etc.

(b) What do you mean by data quality?  
→ Data quality is meant by 'the ability of a given data set to serve an intended purpose'.

(c) Define outlier.  
→ An outlier is an observation point that is distant from other observations.

(d) Define interquartile range.  
→ It is a measure of statistical dispersion which represents range within which central 50% of data value lie.

(e) What do you mean by missing value?  
→ Values in data may not be filled up for various reasons & hence are considered missing.

(f) What are uses of zip files?  
→ Commonly used for mailing & sharing files over the internet.

Q) What do you mean by XML file data format?  
 → An XML is extensible markup language file designed to be human & machine readable & be used to store & transport data.

i) What is tag cloud?

→ A tag cloud is word visualization that displays most used words in terms from small to large according to how often each appears.

Q.2) Attempt any four of the following

a) Explain different application of data science.

→ i) Finance → Fraud detection → Data science is used to detect pattern in transactions.

ii) sales & marketing → Analyse customer data to segment the market for targeted marketing.

iii) E-commerce → Dynamic pricing models adjust price in real time based on demand & competition.

iv) predictive maintenance → Analysis equipment data to predict failure & schedule maintenance.

b) Explain null & alternate hypothesis.

→ Null Hypothesis → Null hypothesis state hypothesis that sample of statistic is equal to population statistic.

e.g. → There has been no diff found in average exam grade of students after coaching duration.

Alternate Hypothesis → Alternate hypothesis states hypothesis there is variation in sample statistic of population.

eg - there has been significant grade improvement found in exam grade of students after conducting coaching classes.

c) what do you mean by noisy data? Explain any two causes of noisy data.

→ - Noisy data contains error or outliers.  
ex - value of age attribute are within range 22-45 years whereas one record reflects the age attribute value as 80.

causes → ① facility data collection instruction errors / human errors  
② data entry problems.

d) what do you mean by data visualization? Give examples of any two data visualization libraries.

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

- libraries

i) matplotlib libraries

- It is common standard python library used for plotting 2D data visualization  
- mainly used for creating plots that zoomed behaviour of plot by using toolbar.

ii) seaborn library

- Seaborn library to create artistic charts with very few lines of code.  
- user creative styles, rich color palettes which user violin plots, Heatmaps, time oriented plots.

Q3 Attempt any two of the following: 3-12

(a) Explain data cube aggregation method in context of data reduction.

- - A data cube is multi-dimensional array of value. A data cube is generally used to easily interpret data.
- useful when representing data together with dimensions and certain measures of business requirement.
- Data cube aggregation is multi-dimensional aggregation which easier multidimensional analysis.

- Data cube present pre-computed summarized data which eases the data mining into fast accessible form.

b) What is mean, mode, & median & range of the following list of values.

24, 29, 29, 25, 27, 25, 32, 24,

→ Mean →  $\frac{24+29+29+25+27+25+32+24}{8}$

Range = 32 - 24 = 8

Mode = 29 (as it appears more than once)

Median = The median is the middle value

if there are even number of values  
the median is the middle value

since there are 9 numbers (odd number)  
the median is the 5th number

Median = 25

Mode  $\rightarrow$  value that appears most frequently  
in data

24 appears five times

$$\text{mode} = \underline{\underline{24}}$$

bimodal bimodal

### ① Histogram

- A histogram is graphical display of data using bars of diff heights.

- To create histogram we need divide the range of value into a series of intervals & second count how many values fall into interval.

### ② Bar chart

- Bar charts are used for comparing quantities of diff categories or groups.

- Values of category are represented by bars with vertical or horizontal bars

- Bar charts has gaps betw. bars compared to histograms

### ③ Line plot

- It is 2-dimensional plotting of value  
- value are displayed in scattered manner f connected.

line plot displays information as series of data points called "markers" connected by straight lines.

(Q.4) Attempt any two of the following

Q) Difference b/w structure & unstructured

structured	unstructured
- This data is typically well organised.	- This data is typically non-organised.
- organised by mean of relational database	- organised by simple character & binary data.
- management of concurrency of data present	- no transaction management of no concurrency present.
- less flexible as well as less scalable	- more flexible as well as more scalable.

⑥ Q) What do you mean by data attributes? Explain types of attribute with example.

→ - An attribute is property or characteristic of an object. A data attribute is a single value descriptor for data object.  
Ex. color of person, name of student

Type ① Nominal means "relating to name"  
Nominal attribute are symbol manner.

e.g. → Branch of relation of ATTONIC company  
has attributes like BID & NAME which are  
nominal type.

② Binary attribute → A binary attribute is  
a nominal attribute with only two categories:  
0, or 1 where 0 is absent, 1 is present

③ Ordinal attribute → An ordinal attribute is  
attribute with possible value that have meaning  
full order in successive value is not known.

e.g. → ranking (e.g. taste of potato chips on  
scale from 1-10).

grader, height in (tall, medium, short).

④ Numeric attributes → A numeric attribute is  
quantitative i.e. it is a measurable quantity  
represented in integer or real values.  
→ ① Interval Scaled

② Ratio Scaled

Interval ex → temperature in celsius or  
fahrenheit etc.

Ratio ex → temperature in kelvin, length  
time, counts, etc.

Q.S.) Attempt any one of the following

(a) What do you mean by data transformation  
complete strategies of data transformation

→ - Data transformation is process of converting raw data into a format or structure that would be more suitable.

(i) Rescaling

- Rescaling means transforming the data so that it falls within specific scale like 0 to 100

- Rescaling data allows scaling all data values.

- Rescaling data allows scaling all data values between specified min and max value.

Original data vs Scaled data

Col1	Col2	Col3	Scaled Data
[0.0, 1.0, 0.78]			

0.0 to 2.0	1.0 to 2.0	0.78 to 1.0	[0.33, 0.99, 1.0]
------------	------------	-------------	-------------------

4.0 to 9.0	6.0 to 9.0	3.6 to 7.0	[1.0, 0.55, 0.77]
------------	------------	------------	-------------------

L = 0.800, 1.9, 1.13  
S = 0.780, 1.0, 0.55

better representation

(ii) Normalizing

- Measurement unit used can affect data analysis.

Normalizing scaled attribute data so as to fall within smaller range 0.0 to 1.0

or -1.0 to 1.0

- Normalizing data attempt to give all attribute an equal weight.

⑧ Binarizing → It is process of converting data to either 0 or 1 based on threshold value.

④ Standardizing →

- It is also called as normalization - it is process of transforming attributes having gaussian distribution with differing mean & standard deviation values into a standard Gaussian distribution with mean of standard deviation 1

⑤ Label encoding → It is process used to convert textual labels into numeric form in order to prepare it to be used in readable form  
e.g. Gender → male → 0 female → 1

⑥ One Hot Coding →

- used to transform categorical variables into a format that can be binary.

- one hot coding refers to splitting the column which contains numerical categorical data into many columns.

depending on no. of categories present in column