

Q.1) Attempt any eight of the following

a) Define volume characteristic of data in reference to data science.

→ Volume refers to sheer scale of data that is being considered for analysis.  
Characteristics → Large datasets → Technological impact.

b) Give example of semistructured data. (B)

→ ① Markup language XML  
② Open standard JSON (JavaScript Object Notation)  
③ No SQL-like query for blunder search

c) Define data discretization.

→ Data discretization is characterized as a method of translating attribute values of continuous data into a finite set of intervals with nominal information.

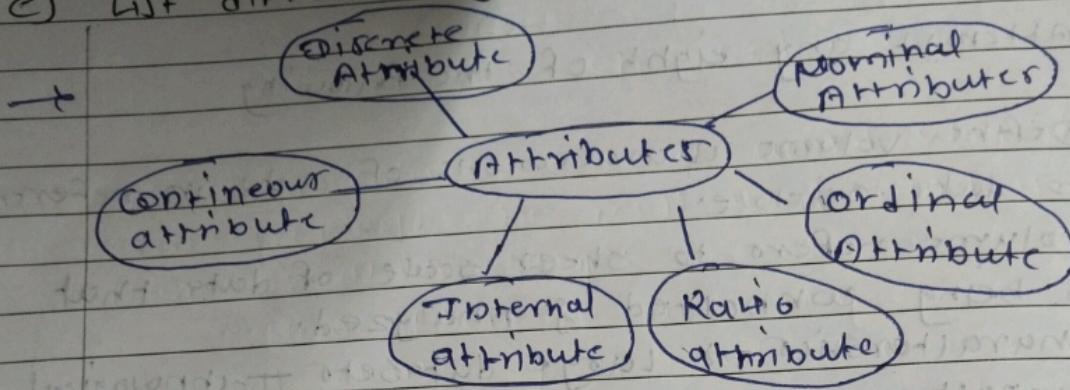
d) What is quartiles?

→ A quartile is a statistical term that refers to division of dataset into equal parts, each containing one-fourth of data points when data are arranged in ascending order.

e) Define data object.

→ A data object is a collection of attributes that together describe an entity or instance in dataset.

c) List different type of attributes.



g) What is data transformation?

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

b) Write the tools used for geospatial data.

- ① GIS software
- ② programming libraries & framework
- ③ spatial dataset and other two things
- ④ Remote sensing toolkit
- ⑤ visualization & mapping tools
- ⑥ Geospatial data formats and tools
- ⑦ web-based geospatial formatter

Q.4 Attempt any four of the following

a) Explain any three ways in which data is stored in files.

→ ① Flat files

- format - Data is stored in next format usually structured or semi-structured such as CSV, TSV, JSON

② Relational Databases -

Data is stored in a more structured format typically in relational database management.

b) Explain role of statistical in data science.

→ - some roles

① predictions of classification from basis of statistical help in prediction if classification of data whether it would be correct or client by either previous usage.

② Help to create probability distribution and estimation of crucial in understanding basis of machine learning & logistic regressions.

③ cross-validation &looou techniques

they are also in hencty statistical tools that have been brought into machine learning data analytic world.

d) Explain any tools in the data science.

→ ① Python programming.

- python refers various libraries designed explicitly for data science operations.

- It was found in 1990 by Guido van Rossum

- The rich set of libraries are core strengths.

② Tableau public

- tableau is data visualization software public which has its free version named as tableau

- It data visualization software packed with powerful graphics to make interactive visualization.

a)

e) Write a short note on word clouds.

→ A word 'cloud' is a word visualization that displays the most used words in a text from small to large.

- How often it appeared

- a word cloud tag or cloud that is a visual representation of textual data that present words as a list of word frequency in descending order.

- most often used for aesthetic purpose at depicting categorical data.

**Media**

**Social**

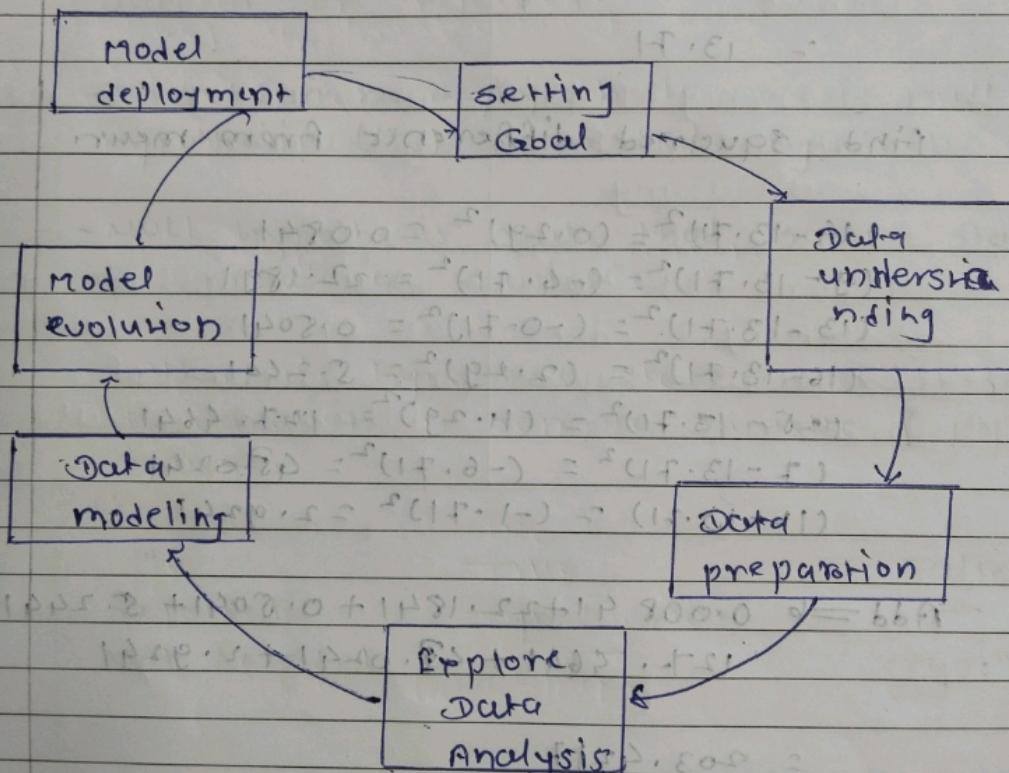
↳

↳

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

a) Explain data science life cycle with suitable diagram.

- - Life <sup>cycle</sup> of data science outliers phase from start to finish
- provides framework for best performance of each phase from creation to irr. completion



- Setting Goal: Entire circle revolves around the business or research goal
- Data understanding: Involves collection of available data.
- Data prep → includes selecting of relevant data.
- Exploratory data analysis: Involves getting idea.

about solution & factor affecting building  
model

→ Data modeling or In heart of data analysis.

Q.4. A

a) W

→ IT

c) calculated variance & standard deviation.  
for the following data.

$x: 14, 9, 13, 16, 25, 7, 12$

→ mean =  $\frac{14+9+13+16+25+7+12}{7}$

$$= 13.71$$

find squared difference from mean

$$(14 - 13.71)^2 = (0.29)^2 = 0.0841$$

$$(9 - 13.71)^2 = (-4.71)^2 = 22.1841$$

$$(13 - 13.71)^2 = (-0.71)^2 = 0.5041$$

$$(16 - 13.71)^2 = (2.29)^2 = 5.2441$$

$$(25 - 13.71)^2 = (11.29)^2 = 127.4641$$

$$(7 - 13.71)^2 = (-6.71)^2 = 45.0241$$

$$(12 - 13.71)^2 = (-1.71)^2 = 2.9241$$

$$\text{Add} = 0.0841 + 22.1841 + 0.5041 + 5.2441 +$$

$$127.4641 + 45.0241 + 2.9241$$

$$= 203.4287$$

$$\text{Variance} = \frac{203.4287}{7} = 29.0612$$

$$\text{Variance} = 29.0612$$

$$\text{Standard deviation} = \sqrt{29.0612}$$

$$= 5.39$$

Q.4. Attempt any two of the following. (d)

Q) Write a short note on hypothesis testing.

→ The process of determine whether stated hypothesis is accepted or rejected from sample data is called hypothesis testing.

→ It is most inferential statistical technique used to check whether a hypothesis is accepted or rejected.

- There can be 2 hypothesis namely null hypothesis ( $H_0$ ) & alternative hypothesis ( $H_a$ )

- Null hypothesis states that sample statistic fits population statistic.

- Alternative hypothesis states that there is variation in sample statistic of population statistic.

$H_0$	True	False
Rejected	Type I error	I
NOT Rejected	✓	Type II error

b) Difference between structure & unstructured data.

Structured	Unstructured
- Data is well organized	- Data is not well organized
- Get organised by mean of relational database	- Based on simple character or binary data.
- Concurrency of data is present & preferred in multitasking	- Concurrency is not present.
- support RDB so versioning is done over rows & tuples	- data is possible only on whole data as no support of database at all.

NYT  
NYT  
NYT  
NYT

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

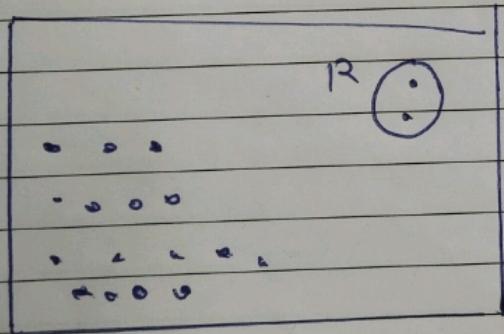
i) Write any two applications of data science.

- i) Image recognition & speech recognition  
ii) Gaming world.  
iii) Internet search.  
iv) Transport.

ii) Explain any type of outlier in detail.

→ i) Global outlier

→ If an individual data point can be considered as anomalous with respect to the rest of data, then data is termed as point outlier.



→ To detect global outlier critical issue is to find an appropriate measurement of deviation with respect to application in question.