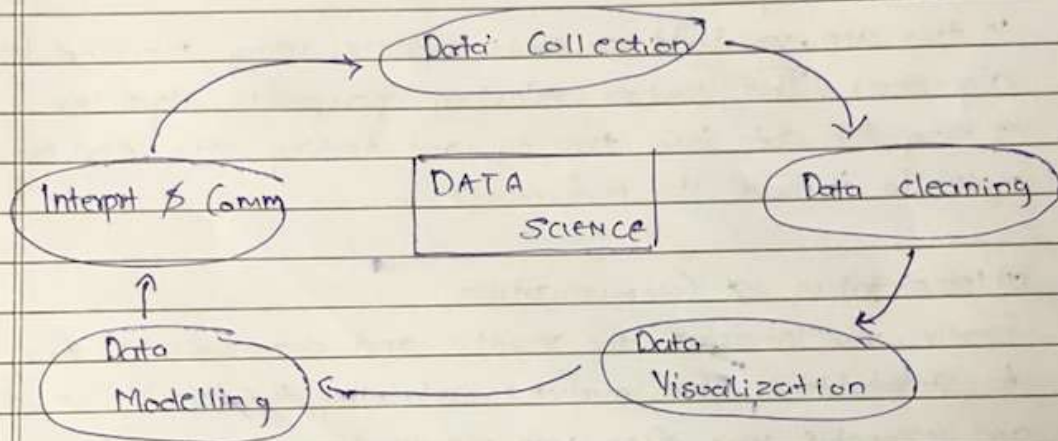


DATA SCIENCE - UNIT 3

Q. State and Explain 5-fundamental Steps of Data Science Process

Ans Data science is the art of science of extracting meaningful insights from data. It blends various fields, including statistics, CS and domain expertise to analyze large volume of data and uncover patterns.



① Data Collection

This is where you gather raw data from various sources. This can be includes databases, API, web scraping, surveys and more. It's crucial to ensure that data is collected is relevant and of high quality.

② Data Cleaning.

Raw data is often messy and filled with errors, missing values or duplicates. In this step, you clean and preprocess the data to make it usable. This might involve filling in missing values, correcting errors and normalizing the data.

⑩ Data Exploration & Visualization

In this the exploration of data is done to understand its structure, patterns, and relationships. Visualization tools are like chart and graph helps you spot trends and insight that might not be obvious from data alone.

⑪ Data Modelling

In this step, you build or train models using machine learning algorithms. This involves selecting appropriate algorithm, splitting the data into training and testing sets, and tuning the model to improve its performance.

⑫ Interpretation & Communication

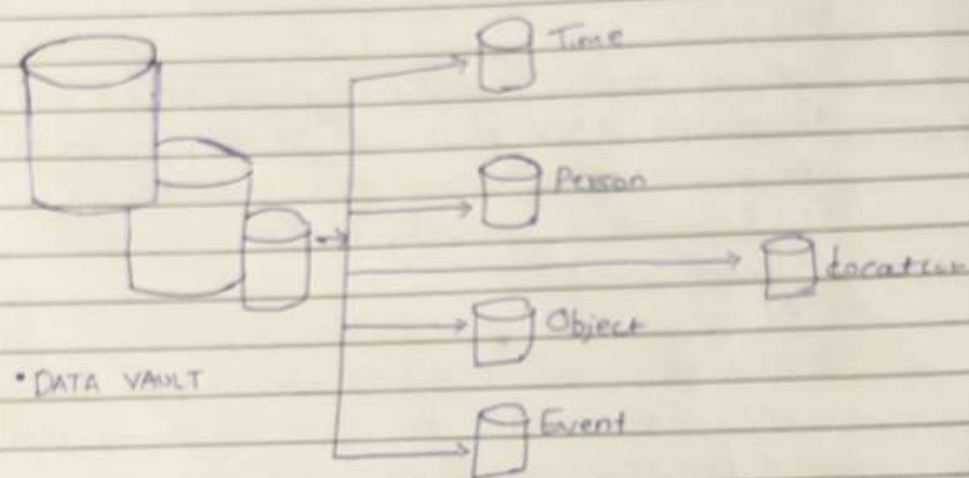
Finally, you interpret the results and communicate your findings to stakeholders. This involves explaining the insights in a clear and actionable way, often with visualization, reports & presentation.

Ex: Imagine you run a small e-commerce store and want to understand your customer's purchasing behaviour.

- ① Gathering data on customer purchases from your store's database.
- ② Remove any incomplete or duplicate entries from your dataset
- ③ Visualize the data to see which products are most popular
- ④ Use a ML algo to predict which products a customer is likely to buy next based on their past purchases.
- ⑤ Discover that customers who buy product A are very likely to also buy Product B. Use this insight to create targeted marketing campaigns & recommendation.

Q: Time - Person - Object - Event - Location Data Vault

Ans: TPOLE Data Vault, It is a powerful framework for organizing and analyzing data, especially ~~over~~ when dealing with a complex and interconnected information.



TABLE

① Time

Captures the temporal aspects of events. This includes dates, times and periods. It helps in tracking when event occur, enabling trend analysis over time.

Ex: A sales transaction on March 15, 2023, at 3:00 PM

② Person

It involves the individuals, groups or entities associated with the event. Essential for understanding who is involved in various events or transaction.

Ex: Abasufyan, a customer making a purchase.

(iii) Object

It Represents the items or entries central to the event. Provides context about what is involved in the event.

Ex: An Vintage Car is been purchased

(iv) Location

Captures spatial Aspects, such as where had the event or the transaction took place. It is crucial for geographic analysis and understanding where the event happen.

Ex: A transaction occurring at a store in Central Park, Mumbai

(v) Event

Describe occurrence or actions, serving as the core of the data model. Encapsulates the activities being recorded, providing a comprehensive view of what happened.

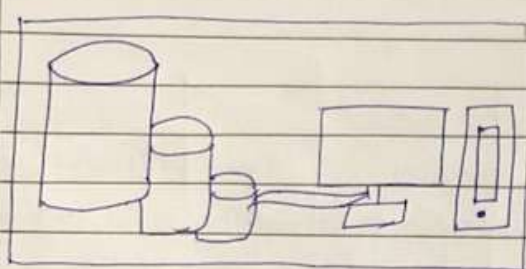
Ex: The actual sales transaction ~~the~~ linking the respective Time, Person, Object and Local Entries.

The TROLE can be written or displayed as:

- Sales Transaction on March 15, 2023 at 3:00 pm by customer AbuSufiyan, purchasing an Vintage Car at a Central Park Store in Mumbai, documenting the purchase event with all valid Papers.

Q. Building an Data Warehouse

Ans. Building an Data Warehouse is like creating a central hub for all of your organization's data.



• To build an Data warehouse all these process and steps are taken in consideration

- ① Req. Analysis
- ② Data Modelling
- ③ ETL Process
- ④ Data Integration
- ⑤ Performance Testing
- ⑥ Report and analysis

① Required Analysis

First you need to understand, what data you need and what you want to achieve with it. This involves talking to stakeholders, identifying data sources and defining key metrics and goals. It sets the foundation for your project.

② Data Modelling

You need to design a data model, which act as a blueprint for how the data will be organized. This involves creating diagrams that represents the structure of data, including facts and dimension tables. This ensures data is organized logically.

⑬ ETL Process

Extract, Transform, Load, after data model designing you start the ETL Process. You extract data from various sources, clean and transform it to ensure it's accurate and consistent, and load it into the data warehouse.

⑭ Data Integration

Once the data is loaded, you integrate it from different sources to provide a unified view. This involves mapping and linking related data, ensuring consistency. This step makes it easier to analyze the data comprehensively.

⑮ Performance Tuning

Optimizing the data warehouse for performance is crucial. This involves indexing and query optimization to ensure that you can retrieve data quickly and efficiently.

⑯ Reporting & Analysis

Finally, you create reports and dashboards to visualize the data and gain insights. This involves developing tools and visualizations that help stakeholders to understand and interpret the data.

Ex: Imagine you run an online bookstore and want to understand customer buying patterns.

① Identify data from sales transactions, cost info and inventory.

② Design schema with tables for sales, Cost - Books - Time

③ ETL: E = Extract sales data, cost info and inventory details

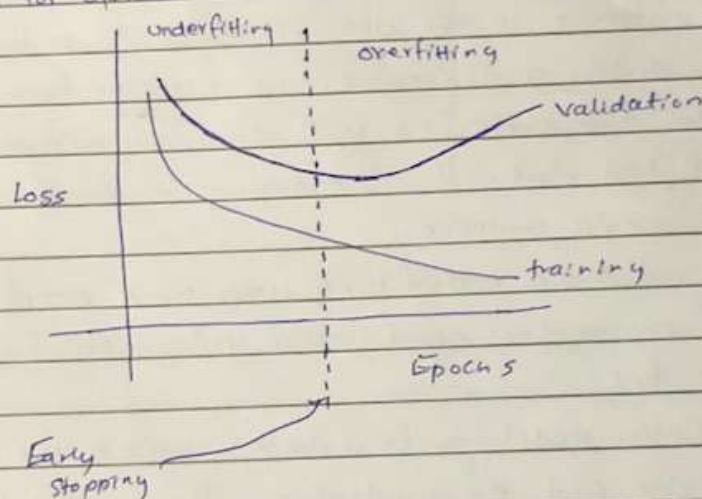
T = Transform data by cleaning and normalizing it

L = Load data into warehouse.

- (iv) Combines sales, customer and book data to provide a unified view of customer buying patterns
- (v) Optimize for queries, such as indexing frequently queried column
- (vi) Create Db to visualize top selling-books, cust demographic and sales trend.

Q. Overfitting and Underfitting

Ans. Overfitting and underfitting are key concepts in machine learning, representing two extremes of model performance that should be avoided for optimal truth.



• OVERFITTING:

Overfitting occurs when a model learns not just the underlying patterns in the training data but also the noise and random fluctuation. A model that overfits will show high accuracy on training data but low accuracy on test data. Overfitting happens when the model is too complex, such as having too many parameters or layers, which allows it to fit too closely to training data.

To prevent and avoid overfitting, use simpler models, prune complex models, employ regulation techniques like L1 and L2 regularization and validate the model using techniques like cross-validation.

Ex: If you train a decision tree to classify whether an email is spam, and the tree is too deep, it might learn specific, irrelevant details unique to the training emails rather than general spam indicators. As a result, it performs poorly on new emails.

• UNDERFITTING:

Underfitting happens when a model is too simple to capture the underlying patterns in the data. It fails to perform well even on the training data, and consequently also performs poorly on new data. A model that underfits will show low accuracy on both training and test data. It's like trying to summarize a complex book with a single sentence.

To prevent and avoid underfitting, use more complex models, increase the number of parameters, and ensure adequate feature engineering and training time.

Ex: Using a linear model to fit a dataset where the relationship between variable and ^{is} quadratic. The model fails to capture the curvature and thus performs poorly.

Balancing Two: The goal in machine learning is to find a balance between overfitting and underfitting.

Q. Cross Validation Test.

Ans. Cross validation test is a technique used in ML and statistics to evaluate the performance of a model. It helps to ensure that the model generalizes well to unseen data, rather than just performing well on the data it was trained on.

Working:

(i) Data Splitting: The original dataset is divided into smaller subsets. One common method is K-fold cross validation, where the data is split into K-equal sized folds.



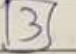
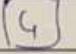
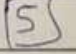
(ii) Training & Testing:

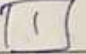

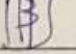
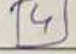
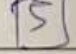
The model is ~~tested~~ trained on K-1 folds and tested on remaining fold. This process is repeated K-times, each time with a different fold as the test set.


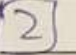


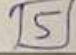
(iii) Performance Averaging:

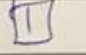
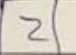
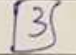

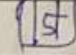
The performance of model is evaluated in each iteration. The results are averaged to produce a single performance task.

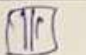
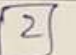
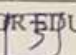
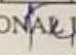

Ex: For a 5-fold cross validation.




i=1     

i=2     

i=3     

i=4     

i=5     

-  - Test set
-  - Training set
-  - Validation set

- (i) Split the data into 5 parts
- (ii) Train the model on 4 parts and test on 5th
- (iii) Repeat this process 5 times, each time with a different part as the test set
- (iv) Average the results to get final performance matrix

Q Hypothesis.

Ans: A Hypothesis is an statement or an assumption about population parameter that you want to test. It's often an educated guess based on prior evidence

Hypothesis

NULL H₀

- This hypothesis suggests there's no difference no effect. Think of it as an "default" assumption.

ALT H₁

- This hypothesis suggests there is an effect or a difference. It's what you want to prove

• Test in Hypothesis.

CHI SQUARE TEST.

A chi square test is a statistical test used to compare observed results with expected results. The purpose of this test is to determine if a difference between observed data and expected data is due to chance or if it is due to relationship between variables.

Ex:

Gender	Action	Comedy	Total
Male	40	10	50
Female	20	30	50
Total	60	40	100

V_1 (value 1) V_2 (value 2)

Null = No Associ
HYP Between
 G / M
Alt : ASSO
 Between
 G & M

Soln:

$$\text{Form} = \frac{\text{Total 1} \times \text{Total 2}}{\text{Grand Total 1}}$$

$$= \frac{50 \times 60}{100} = 30$$

$$\frac{50 \times 40}{100} = 20$$

$$= \frac{50 \times 60}{100} = 30$$

$$\frac{50 \times 40}{100} = 20$$

Expected Value	Male	Female
male	30	20
Female	30	20

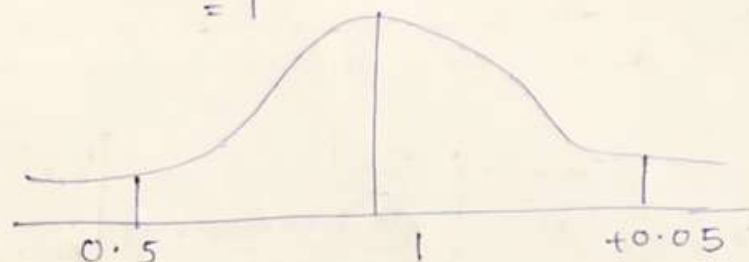
Chi square Test Statistics

$$\chi^2 = \frac{\sum (O - E)^2}{E}$$

Observed v	Expected E	O - E	(O - E) ²	χ^2 - chi
40	30	10	100	0.3
20	30	-10	100	0.3
10	20	-10	100	0.2
30	20	10	100	0.2
				1

(Critical Value of 1 = 2.706
from log table)

$$\begin{aligned}\text{Degree of Freedom} &= \text{Col} - 1 \times \text{row} - 1 \\ &= 2 - 1 \times 2 - 1 \\ &= 1\end{aligned}$$



$\therefore \chi^2$ calculated value (1) $> \chi^2$ tabular value critical (2.706), ~~Null Hypo rejected~~

\therefore There is no Association btwn G & M

If calc value is less than tab value we will accept Null hypo

Q T-test

Ans: It is an Statistical tool used to compare the means of two or more group. It's ratio that measure significance of difference btwn means of groups while taking their variance etc. -

Ex: Group 1 [Old Method] : [55, 60, 65, 70, 62]
Group 2 [New Method] : [68, 75, 80, 85, 78]

S1: Null: No relation btwn G1 & G2
Alt: Relation btwn G1 & G2

S2: G1: $\bar{x} = 62.4$
G2: $\bar{x} = 77.2$

S3: Variance:

$$\begin{aligned}& \text{G1} \\ &= \frac{(\text{Value} - \bar{x})^2 + (\text{Value} - \bar{x})^2 + \dots + (\text{Value} - \bar{x}_n)^2}{n-1} \\ &= \frac{(55 - 62.4)^2 + (60 - 62.4)^2 + (65 - 62.4)^2 + (70 - 62.4)^2 + (62 - 62.4)^2}{5-1} \\ &= \frac{(-7.4)^2 + (-2.4)^2 + (2.6)^2 + (7.6)^2 + (-0.4)^2}{4} \\ &= \frac{54.76 + 5.76 + 6.76 + 57.76 + 0.16}{4}\end{aligned}$$

$$= \frac{125.2}{4}$$

$$= 31.34$$

Variance σ^2 :

$$= \frac{(value - \bar{x})^2 + (value - \bar{x})^2 + \dots + (value - \bar{x})^n}{n-1}$$

$$= \frac{(68 - 77.2)^2 + (75 - 77.2)^2 + (80 - 77.2)^2 + (85 - 77.2)^2 + (78 - 77.2)^2}{5-1}$$

$$= \frac{(-9.2)^2 + (-2.2)^2 + (2.8)^2 + (7.8)^2 + (0.8)^2}{4}$$

$$= \frac{84.64 + 4.84 + 7.84 + 60.84 + 0.64}{4}$$

$$= \frac{158.8}{4}$$

$$= 39.7$$

$$t = \frac{\bar{x}_1 - \bar{x}_2}{\sqrt{\left(\frac{s_1}{n_1}\right)^2 + \left(\frac{s_2}{n_2}\right)^2}}$$

$$= \frac{62.4 - 77.2}{\sqrt{\left(\frac{31.3}{5}\right)^2 + \left(\frac{39.7}{5}\right)^2}}$$

$$= \frac{62.4 - 77.2}{\sqrt{\frac{31.3 + 39.7}{5}}}$$

$$= \frac{-14.8}{14.2} = -3.92$$

$$\text{Degree of Freedom} = n_1 - 1 + n_2 - 1$$

$$= 5 - 1 + 5 - 1$$

$$= 8$$

$$\text{By Def Consi } \alpha = 0.05 (\text{in } t_{\alpha/2})$$

$$= 2.306$$

Cal > Tabular

∴ Rejected