

Q 1 paper

~~ETL~~

ETL Life Cycle real life

- ① Cycle initiation
- ② Build reference data
- ③ Extract (from sources)
- ④ validate
- ⑤ Transform (Clean, Apply business rules, Check for integrity, Create aggregates or disaggregates)
- ⑥ Stage (Load into staging tables, if used)
- ⑦ Audit reports (for example, on Compliance with business rules, also in case of failure, helps to diagnose/repair)
- ⑧ publish (to target tables)
- ⑨ Archive
- ⑩ clean up

① Cycle Initiation:

- i) Start the ETL process, often triggered by a schedule, an event, or specific business need.
- ii) ex: initiation the ETL process every night to update a data warehouse with the latest sales data.

② Build Reference Data:

- i) Create or update reference data needed for transformation and validation processes.
- ii) ex: Building a reference table or product categories that can be used to categorize and enrich sales data

③ Extract (from sources)

- i) Gather data from various sources, such as databases, files, or APIs

- ii) ex:
Extracting customer information from an online store's database and order data from a legacy system

④ Validate

- i) Check the extracted data for accuracy, completeness, and adherence to predefined rules

- ii) ex:
verifying that all order records have valid customer ID and that the total order amount is within an acceptable range.

at p. 10

⑤ Transform

- i) clean, apply business rules, and manipulate the data to fit the target format.
- ii) ex converting date formats, aggregating sales data by region and calculating average order value.

⑥ Stage (Load into Staging)

- i) Load the transformed data into staging tables for further processing.
- ii) ex storing cleaned and transformed customer and order data in temporary tables before loading them into the main database.

⑦ Audit Reports

- i) Generate reports to track compliance with rules & diagnose issues in case of failures.
- ii) ex creating a report that highlights any data anomalies, such as missing values or discrepancies between source and target data.

⑧ Publish (To target tables)

- i) move the validated and transformed data to the final destination or target tables.
- ii) ex loading the cleaned and validated sales data into the main data warehouse for business analytics.

⑨ Archive:

- i) Store the historical or older versions of data for future reference or regulatory compliance.
- ii) ex Archiving quarterly financial data to maintain a historical record for auditing purposes.

⑩ Clean Up:

- i) Remove temporary files, tables or any other artifacts created during the ETL process
- ii) ex Deleting staging tables and temporary files used during the ETL cycle to free up storage space

* Q2 paper

Data Source Selection Criteria:

- ① Relevance to project Objectives
 - * meaningful insights
- ② Data quality & Accuracy
 - * accurate info
- ③ Timeliness of data
 - * up-to-date information
- ④ Data Volume and Scope
 - * sufficient data
- ⑤ Data Compatibility and format
 - * integration
- ⑥ legal and ethical compliance
 - * copyrights
- ⑦ cost and resource constraint
 - * financially

Q2

Data source selection criteria

- ① credible (reputable & trustworthy)
- ② complete (all necessary information)
- ③ verifiable (accuracy & authenticity)
- ④ accurate (free from errors)
- ⑤ current (up-to date)
- ⑥ compliance (legal & regulatory)
- ⑦ accessible (easily accessible)
- ⑧ cost (reasonable cost)
- ⑨ legal (copyright, intellectual properties)
- ⑩ security (protection)
- ⑪ storage (efficient stored)
- ⑫ provenance

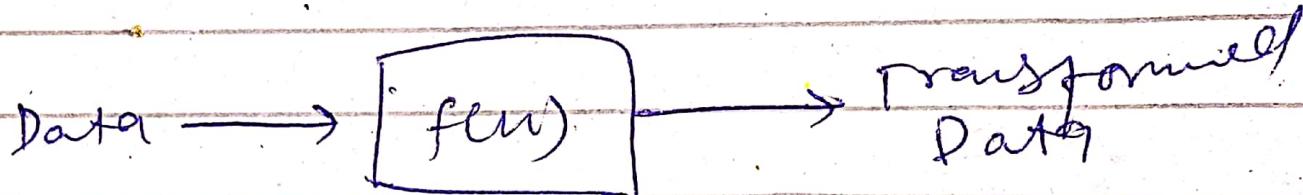
↳ Understanding the origin & history of the data, including how it was collected & processed.

↳ Documenting the source and processing steps for research data.

~~#Q3 paper~~

Data Transformations

- ① ⚡ Apply a function $f(x)$ to adjust scales of data
- ② Done usually when data is skewed, so that it becomes easier to perform modelling.
- ③ Done to convert non linear relationship into ~~linear~~ linear relationship.



(Natural) log Transformation

- ① To transform data that is positively skewed.
- ② Usually done when data is concentrated to zero
(most values are close to zero but few are much larger)

log Transformation

⇒ Makes the relationship b/w two variables more linear.

↳ many simple modeling methods assume a linear relationship b/w variable

⇒ linear regression

Other Transformation

- ① Square Root
- ② Square
- ③ Inverse

f

Q4

- ↳ Accuracy is not good evaluation measure when dealing with imbalanced datasets
- ↳ Unbalanced datasets occur when one class significantly outnumbers the other.

Example

Total transaction = 10000

non fraudulent = 99000

Fraudulent = 1000 (1% of total)

* Model prediction all transaction as non fraudulent to maximize accuracy

=> TP = 0 (model incorrectly predicts all fraudulent transaction as non fraudulent)

=> True N: 99000 (Model ^{correctly} predicts no fraudulent transaction)

=> FP: 0 (model incorrectly predict non as fraud)

=> False N = 1000 (model pairs to ^{predict} fraud. from)

$$\text{Accuracy} \Rightarrow \frac{TP}{(TP+FP)} = 0/0+0 = 0\%$$

$$\text{Accuracy} = \frac{(TP+TN)}{\text{Total}}$$

$$= 0 + 99000 / 100000 = \underline{\underline{99\%}}$$

* Q5 paper

Explain the role of data scientist

- ① Explorers^{DS}: Explore and analyse large sets of data to uncover meaningful patterns, trends and insights.
- ② problem solver: They use data to solve complex problems and make informed decisions, helping business and organizations.
- ③ Predictor: Build models to predict future outcomes.
- ④ Story teller: Turn data into actionable insights.
- ⑤ continuous learner: Adapting new technologies, methods and data scientists stay updated.

* Paper

Q6:

a) `df[["color"]] = df[["Department"]].map(
 {"IT": "red", "HR": "green", "Finance": "blue",
 "Marketing": "orange"})`

b) `plt.hist(df[["Department"]], bins=len(df["Department"].
unique))`

c) `Sns.lineplot(x="Age", y="Experience", hue="Department")`

d) `plt.hist(df[df["Department"] == "IT"]["Age"],
 df[df["Department"] == "HR"]["Age"],
 color=["red", "green"], label=[IT, HR])`

e) `df = df.drop("Salary", axis=1)`

Q7

(a) `data[“Title”].value_counts()`

output

Title	
SE	1
PM	1
HR M	1
QA	1
A	1
M A	1
HR Inter	1

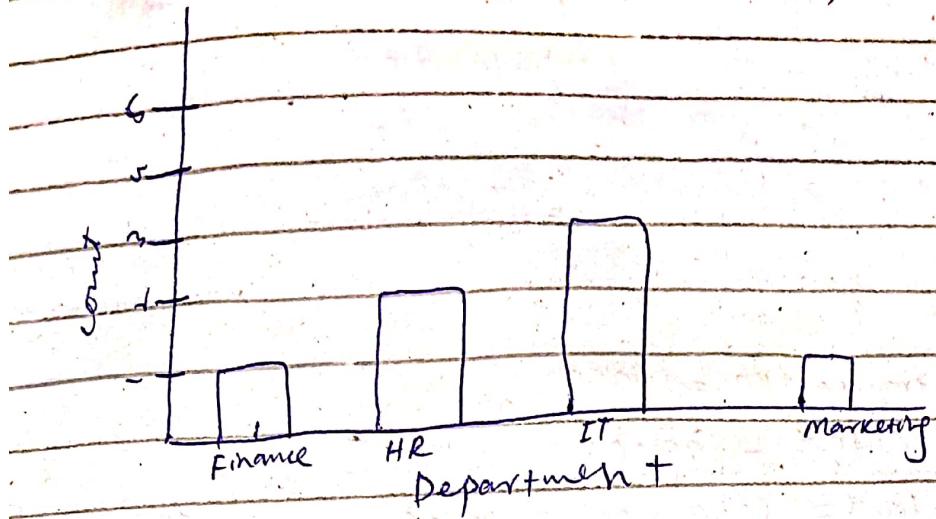
Name : count, dtype: int64

(b) `data[“Department”].describe()`

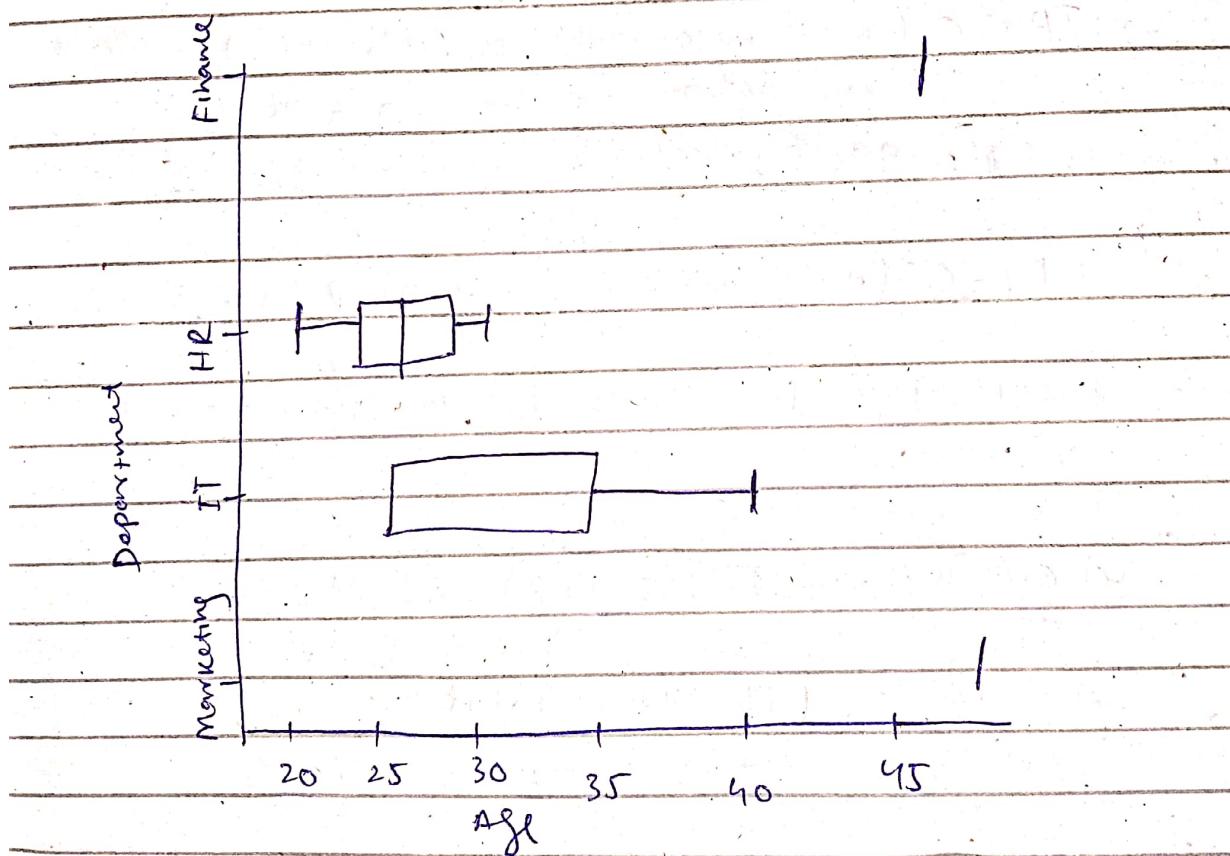
output

Count	7
Unique	4
top	IT
freq	3

c) sns.countplot(data[["Department"]])



d) sns.boxplot(data[["Age"]], data[["Department"]])



Apriori algorithm

Q8 paper

Generate rules using Apriori Algorithm

Support = 50% confidence = 75%

Transaction ID

1

2

3

4

5

Items purchased

Bread, cheese, Egg, juice

Bread, cheese, juice

Bread, Milk, yogurt

Bread, juice, milk

Cheese, juice, Milk

① Find Frequent Item set and their support

Item	Freq	support in (%)
Bread	4	4/5 = 80%
Cheese	3	3/5 = 60%
Egg	1	1/5 = 20% X
Juice	4	4/5 = 80%
Milk	3	3/5 = 60%
Yogurt	1	1/5 = 20% X

Support(item) = Freq of item / no. of transactions

② Remove all the items whose support is below given min support

New Table

Item	Freq	Support
Bread	4	$4/5 = 80\%$
Cheese	3	$3/5 = 60\%$
Juice	4	$4/5 = 80\%$
Milk	3	$3/5 = 60\%$

③ Now form two items candidate set & write their freq

Items pair	Freq	Support (in %)
Bread, cheese	2	$2/5 = 40\%$
Bread, juice	3	$3/5 = 60\%$
Bread, milk	2	$2/5 = 40\%$
Cheese, juice	3	$3/5 = 60\%$
Cheese, milk	1	$1/5 = 20\%$
Juice, milk	2	$2/5 = 40\%$

Ans

- ④ Remove all items whose support is below the given min support

item pairs	freq	support %
Bread, juice	3	3/5 = 60%
Cheese, juice	3	3/5 = 60%

- ⑤ Generate rules

For rules Consider

① (Bread, juice)

Bread \rightarrow juice & juice \rightarrow Bread

② (Cheese, juice)

Cheese \rightarrow juice and juice \rightarrow cheese

$$\text{confidence } (A \rightarrow B) = \text{support}(A \cup B) / \text{support}(A)$$

$$\begin{aligned} \text{① confidence } (\text{Bread} \rightarrow \text{juice}) &= \frac{\text{support}(\text{Bread} \cup \text{juice})}{\text{support}(\text{Bread})} \\ &\Rightarrow \frac{3/5 + 5/4}{5/4} = \frac{3/4}{5/4} = 75\% \end{aligned}$$

$$\text{② confidence } (\text{juice} \cup \text{Bread}) / \text{support}(\text{juice})$$

$$\Rightarrow \frac{3/5 + 5/4}{5/4} = \frac{3/4}{5/4} = 75\%$$

③ Confidence (cheese \rightarrow juice) = $\frac{\text{support}(\text{cheese} \wedge \text{juice})}{\text{support}(\text{cheese})}$

$$\Rightarrow \frac{3}{5} * \frac{5}{3} = 1 = 100\%$$

④ Confidence (juice \rightarrow cheese) = $\frac{\text{support}(\text{juice} \wedge \text{cheese})}{\text{support}(\text{juice})}$

$$\Rightarrow \frac{3}{5} * \frac{3}{4} = \frac{3}{4} = 75\%$$

Q9

⑨

cons

- ① computational complexity: if no. of itemset grows exponentially.
- ② memory usage: requires multiple passes through data. leads to high memory usage.
- ③ fixed threshold
 - ↳ Setting min support threshold.
 - ↳ choice can impact the results.

pros

- ① simple concept: easy
- ② scalability: efficient
- ③ widely used:

Definition: ~~classical~~ algo in data mining.
for association rule learning
↳ used to discover frequent item sets in a dataset
↳ generate association rules based on these itemsets.

(a) $Q_0 = \text{min value} = 20$

$Q_{41} = \text{max value} = 47$

(c)

~~Q 10~~

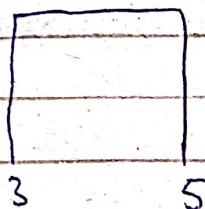
Required

FAC (Complete), Datalogram

Given

	P1	P2	P3	P4	P5
P1	0				
P2	9	0			
P3	3	7	0		
P4	6	5	9	0	
P5	11	10	2	8	0

①



② The ^{distance} matrix will update.

	P1	P2	P3, P5	P4
P1	0			
P2	9	0		
P3, P5	11	10	0	
P4	6	5	9	0

① $\max [\text{dis}(P_3, P_5), P_1]$

$\Rightarrow \max [\text{dis}(P_3, P_1), (P_5, P_1)]$

$\Rightarrow \max [3, 11] \Rightarrow 11$

② $\max [\text{dis}(P_3, P_5), P_2] \Rightarrow \max [\text{dis}(P_3, P_2), (P_5, P_2)]$

$\Rightarrow \max (7, 10) \Rightarrow 10$

③ $\max [\text{dis}(P_3, P_5), P_4] \Rightarrow \max [\text{dis}(P_3, P_4), (P_5, P_4)]$

$\Rightarrow \max (9, 8) \Rightarrow 9$

P.T.O

2 4

③ The distance matrix will update

P_1	P_2, P_4	P_1	P_2, P_4	P_3, P_5
P_2, P_4	9		0	
P_3, P_5		11	10	0

$$\text{①} \Rightarrow \text{Max} \{ \text{dis} \{ (P_2, P_4), (P_1) \} \}$$

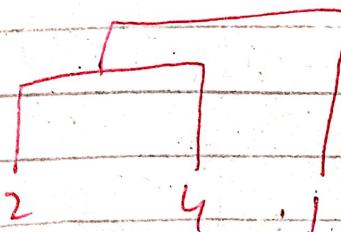
$$\Rightarrow \text{Max} \{ \text{dis} \{ (P_2, P_4), (P_4, P_1) \} \}$$

$$\Rightarrow \text{Max} \{ 9, 6 \} = 9$$

$$\text{②} \Rightarrow \text{Max} \{ \text{dis} \{ (P_2, P_4), (P_3, P_5) \} \}$$

$$\Rightarrow \text{Max} \{ \text{dis} \{ (P_2, (P_3, P_5)), (P_4, (P_3, P_5)) \} \}$$

$$\Rightarrow \text{Max} \{ \text{dis} \{ 10, 9 \} \} = 10$$



(4) The distance matrix will update

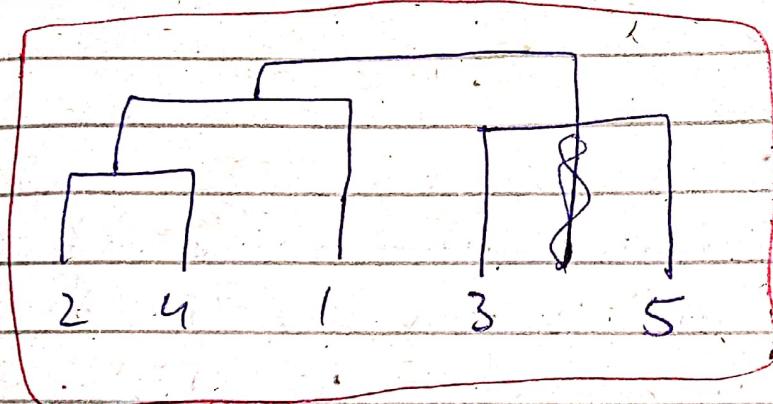
P_1, P_2, P_3	P_1, P_2, P_4	P_3, P_5
P_1, P_2, P_4	0	
P_3, P_5	11	0

$$1) \max [\text{dis}(P_3, P_5), \text{dis}(P_1, (P_2, P_4))]$$

$$2) \max [\text{dis}(P_1, (P_3, P_5)), \text{dis}(P_2, P_4), \text{dis}(P_3, P_5)]$$

$$\Rightarrow \max [11, 10]$$

$$\Rightarrow 11$$



Q8 Q9 paper Cons (b)

- ① Computational Complexity: can be computationally expensive, with a large number of data points.
- ② Sensitivity to noise: Susceptible to noise & outliers, which can influence the merging process.
- ③ Memory usage: Requires storing a potentially large distance matrix, leading to high memory usage.
- ④ Fixed Hierarchy: once clusters are merged, the hierarchy is fixed, and changes require recomputing from scratch.