

U-2 (ML)

DOMS

Page No.

Date

/

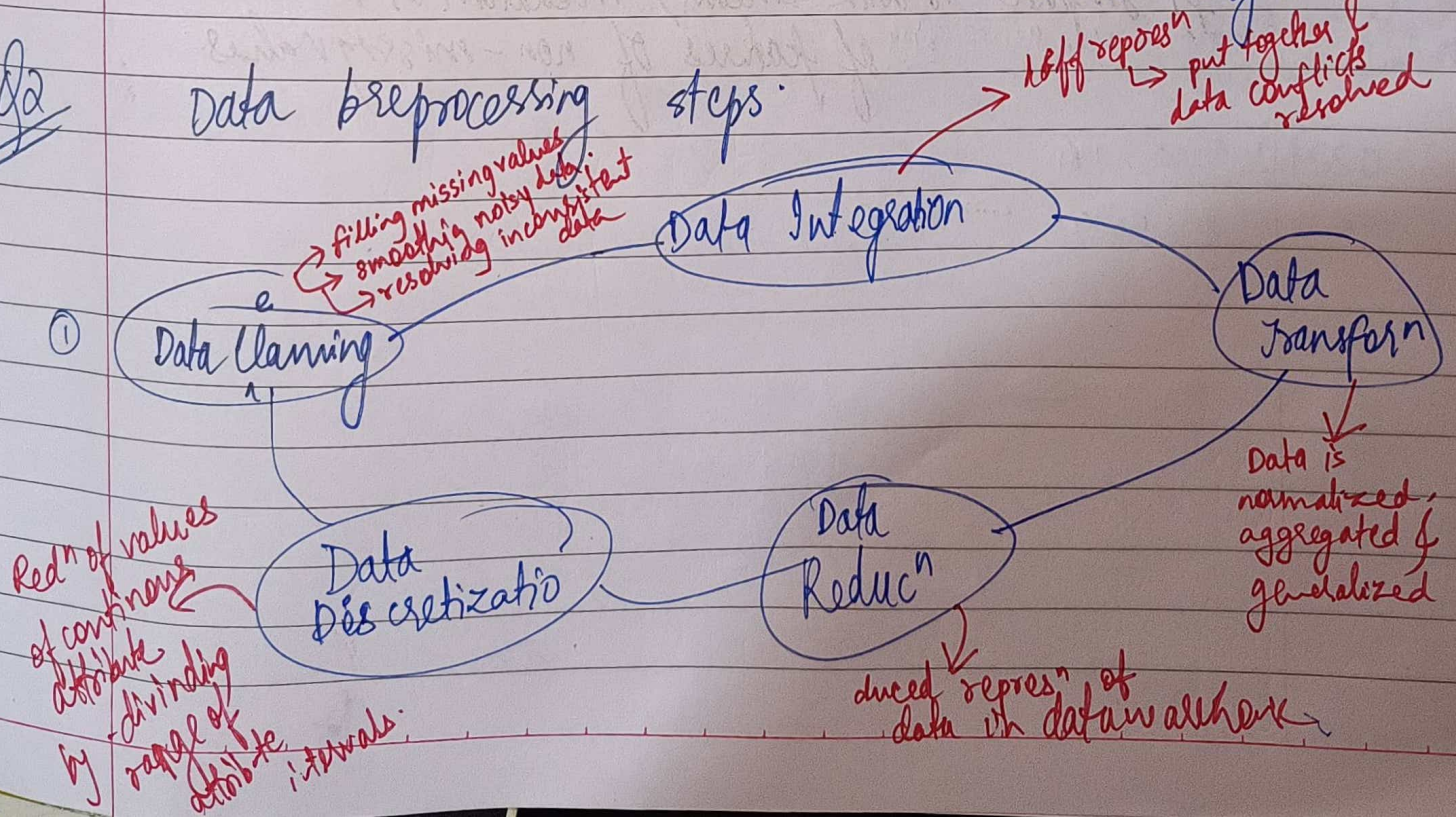
/

Q1 Feature Engineering -

Raw data \rightarrow set of meaningful insights & features.
 \rightarrow improve performance of ML model.

- 1) Feature Selection - subset of most relevant features from original set of variables
- 2) Feature Extraction - creates new features by applying mathematical or statistical transformations
- 3) Feature transformation - existing feature \rightarrow more suitable for model
• scaling, normalization, etc.
- 4) Feature Creation \rightarrow new feature \rightarrow domain knowledge or prior experience
• specific patterns or relⁿ not evident in original data

Q2 Data preprocessing steps.



Steps → Data Preprocessing -

- ① Get dataset
- ② Importing Libraries
- ③ "—" Dataset
- ④ Handle missing Data
- ⑤ encoding categorical data
- ⑥ split data → Train Test
- ⑦ Feature Scaling

Q3

Handle missing data in dataset.

- 1) Ignore the tuple
- 2) Fill in missing value manually
- 3) Use global constant to fill missing value.
- 4) Use attribute mean to fill missing value
- 5) Use attribute mean for all samples belonging to same class as given tuple.
- 6) Use Most probable value to fill in missing value.
- 7) Impute it with mean, median or mode of features of non-missing values.

feature selection Techniques

↳ Method of reducing input variable to use model by using only relevant data & rid of noise in data.

- reduce dimensionality of feature space
- speed up a learning algorithm
- improve predictive accuracy
- improve comprehensibility of learning results

Filter Method

Wrapper methods

Embedded Methods

- | | | |
|---|---|---|
| <ul style="list-style-type: none"> • generic method
↓
don't incorporate specific ml ml algo. • faster than wrapper method • less prone to overfitting • ex → chi-square test | <ul style="list-style-type: none"> • evaluates on specific ml algo • high computation time • high chance of overfitting • Forward, Backward, Stepwise selection | <ul style="list-style-type: none"> • features → model building process
feature select → each iteration of model. • In b/w wrapper & filter time taken. • reduce overfitting by penalizing the coefficient of model. • ex → LASSO. |
|---|---|---|

Q5

Statistical Measures \rightarrow Feature Engineering

selecting & transforming variables/features in dataset for creating predictive model.

Count based feature selection \rightarrow 1) count of individual values within column

2) idea of distribution & range of values

1) Mean \rightarrow Mean of dataset is avg of all data value.

$$\text{Sample Mean} = \bar{X} = \frac{\text{Sum of values of the observation}}{\text{Total no. of obsv}^n \text{ in sample}} = \frac{\sum x_i}{n}$$

$$\text{Popl}^n \text{ Mean} = \mu = \frac{\text{Sum of value of } N \text{ obsv}^n}{\text{No. of obsv}^n \text{ in popl}^n} = \frac{\sum x_i}{n}$$

2) Median \rightarrow data set's value in middle when data items are arranged in ascending order.

3) Mode \rightarrow value that occurs is greatest frequency.

Q6

PCA -

- Principal Component Analysis
- unsupervised ML algo
- used for dimensionality reduction.
- converts observation of correlated function into set of linearly uncorrelated features.
- new transformed features \rightarrow principal component.
- popular tool \rightarrow EDA & predictive modelling.
- technique to draw strong patterns from reducing variance.
- Appln \rightarrow image preprocessing, movie recommendation system, optimizing power allocation.
- contains imp variables & deletes/drops least important variable.

Q7

Multi Dimensional Scaling (MDS)

- non-linear technique \rightarrow embedding data in lower-dimensional space.
- map points residing in MDS to L₂DS.
- dimensionality reduction \rightarrow input data \rightarrow not linearly arranged.
- visual representation of distance or dissimilarities b/w of high dimensional data.
- iterative & minimize diff b/w pairs of points in original data.
- preprocessing step \rightarrow dimensionality red'n \rightarrow classif'n & regression problem.

MDS \rightarrow Metric MDS / classical MDS \rightarrow preserve pairwise distance / dissimilarity measure as much as possible.

Non Metric MDS

\hookrightarrow ranks of dissimilarity metric is known.

Q8

Matrix Factorization -

- decomposing matrix in product of other matrix

appln \rightarrow collaborative filtering

- matrix \rightarrow smaller matrix \rightarrow uncover hidden patterns, simplify computations, etc.

- used in recommendation system

\rightarrow decomposes user item interaction data into user & item matrices to generate personalized recommendations based on latent features

Q9

feature scaling -

- every feature in same footing without any upfront importance.

- algorithm performance improvement
- preventing numeric instability
- each characteristics \rightarrow same consideration.

- 1) Min-Max Scaling \rightarrow find min & max value of column
 \rightarrow subtract min from entry & divide result by diff b/w max & min.

values are shifted & rescaled so their range can vary from 0 to 1.

$$X_n = \frac{X - X_{\min}}{X_{\max} - X_{\min}}$$

- 2) Simple Feature scaling \rightarrow divided each value by max. value for feature.

89 Normalization

- data preparation technique
- transforming columns in a dataset to same scale
- features more consistent with each other → model → accurate output
- make data homogenous over all records & fields.
- rescaling real-valued numeric attributes into 0 to 1 range.
- algo → KNN, SVM, Neural Networks & principle networks

Types → 1) Z-score or standard score.

• values are normalized based on mean & standard deviation of data.

$$X_{\text{new}} = \frac{X_{\text{old}} - \mu_A}{\sigma_A}$$

μ_A - standard mean σ_A - standard deviation

Q10 -

PCA helps in dimensionality reduction by -

- reduce no. of features or variable in dataset.

- 1) Data Compression -
- 2) Variance Maximization
- 3) Dimension Ranking
- 4) Noise Reduction
- 5) Interpretability
- 6) Improved Model Performance
- 7) Visualization
- 8) Feature Engineering
 - 1) Trade off → amt of variance retained & no. of dimensions reduced.

Q11 - feature extraction & types → Kernel PCA.
Local Binary Pattern.

- dimensionality reduction technique
- transform high dimensional data into lower dimensional data while preserving essential info & patterns.

* Kernel PCA -

- extension of standard PCA.
 - Kernel method captures non-linear relationship in data.
 - by mapping into high dimensional data.
 - result is nonlinear feature → analysis or modeling.
- Applⁿ → PCA insufficient → Kernel PCA helps.
in OpenCV, bioinformatics & etc.

* Local Binary Pattern (LBP).

- Text Descriptor → used in img analysis & computer vis.
- quantifies local texture patterns in image by comparing neighbouring pixels.
- texture analysis, object recognition, face recognition, etc.
- can capture intricate texture details.

Q12

Backward & Forward Selection Process.

- Used for feature selection in ML.
- choose subset of relevant features from original -
 - improve model performance, reduce overfitting, etc.

✱

Forward selection →

• Iterative Process.

• starts with empty set of features & gradually adds one feature at a time to build final feature subset

Steps → 1) empty set 2) evaluate feature 3) add selected feature
4) Iterate 5) Final Model → has selected features.

✱

Backward selection -

starts with all available features & progressively removes one feature at a time to determine final feature subset.

steps → 1) all feature - empty list / track of eliminated features
2) evaluate all feature
3) Remove Least significant
4) Iterate
5) Model → has required subset of features.