

Q.3) Explain KDD process & Different Data mining tech.
→ - Data mining is a synonym for another popularly used term 'KDD (Knowledge Discovery from Data)', while other view Data mining as merely as an essential step in process of Knowledge Discovery.

- KDD Process -

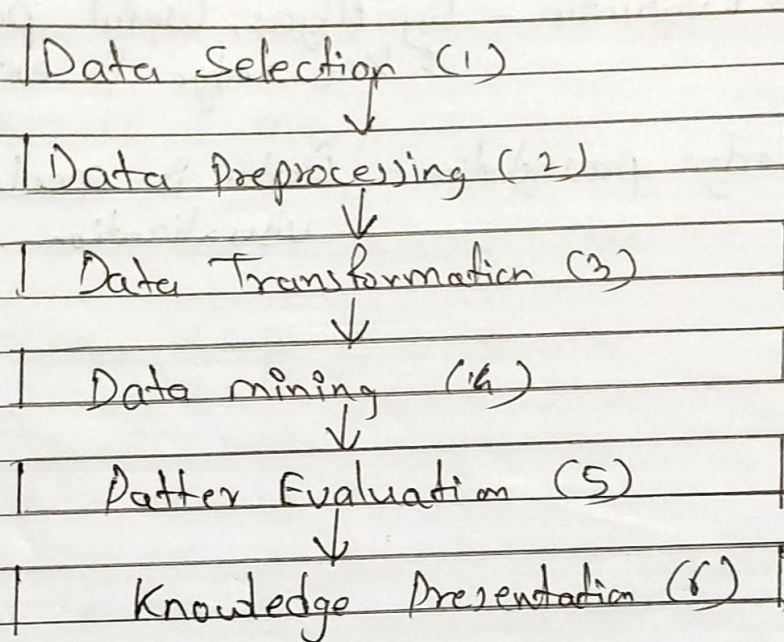


Fig. KDD Process

- Data Selection : Relevant Data is selected from various sources.
- Data Preprocessing : Data is made into a consistent state.
Unnecessary information is removed.

- Data Transformation - We convert the Data into Req. format.
- Data mining - We use various Data mining Algo for pattern finding.
- Pattern Evaluation - By Algos, Useful patterns & knowledge is extracted.
- Knowledge presentation - Data is made for visualization.

Q.4) Suppose that the Data for Analysis includes the attribute age. Age values for the Data tuples are: 13, 15, 16, 16, 19, 20, 23, 29, 35, 41, 49, 53, 62, 69, 72. Use min-max normalization to transform the value 45 for age onto the range 0 to 1.

$$V_i' = \frac{V_i - \text{min}_A}{\text{max}_A - \text{min}_A} (\text{new-max}_A - \text{new-min}_A) + \text{new-min}_A$$

$$\text{min}_A = 13 \text{ \& } \text{max}_A = 72$$

$$\text{new-min}_A = 0 \text{ \& } \text{new-max}_A = 1 \text{ - (Given)}$$

$$V_i = 45 \text{ ... (Given age to be Transformed)}$$

$$V_i' = \frac{(45 - 13)}{(72 - 13) \times (1 - 0) + 0}$$

$$V_i' = 0.542$$

\therefore Age 45 is transformed to 0.542, within the range 0 to 1.

Q.5) Diff betn Classification & Prediction.

→

Classification

i) Classification is a major type of prediction problem where classification is used to predict or nominal values.

ii) Classification is the use of prediction to predict class.

iii) Eg. Group patients based on their known medical data & treatment outcome then it's a classification.

iv) Eg. use cases are like Spam email Detection, Disease Diagnosis.

Prediction.

i) Prediction can be viewed as the construction & use of a model to assess the class of an unlabeled sample.

ii) It is used to access the values or value ranges of an attribute that a given sample is likely to have.

iii) Eg. If a classification model is used to predict the treatment for a new patient, then it would be a prediction.

iv) Eg. use cases are like Stock prediction, Sales forecasting.

Q.6) Explain Accuracy & error measure for classifiers.

→ - Accuracy of a classifier M , $\text{acc}(M)$ is the percentage of test set tuples that are correctly classified By model M .

A) Partition the Data randomly into 3 sets - Training set, validation set & test set.

- Training set is the subset of data used to train / build the model.

- Test set is a set of instances that have not been used in the training process. The model's performance is evaluated on unseen Data. Testing just estimates the probability of success on ~~unknown~~ unknown data.

- Validation data is used for parameter tuning but it can't be the test Data. Validation data can be the training data, or a subset of training data.

• General Error : Model error on the test Data.

B) Success :- Instance (record) class is predicted correctly.

C) Error :- Instance class is predicted incorrectly.

D) The confusion matrix :- It is a useful tool for analyzing how well your classifier can recognize tuples of different classes.

If we have only 2 way classification then only four classification outcomes are possible ~

	Actual class	Class Label	Predicted Class		Total
			C ₁	C ₂	
	C ₁	→	True Positive (TP)	False Negative (FN)	P
	C ₂		False Positive (FP)	True Negative (TN)	N
	Total		P'	N'	All

TP - Class member which are classified as class member.

TN - Class non-members which are classified as non-members

FP - Class non-members which are classified as class members

FN - class members which are classified as class non-members.

P - No. of positive tuples

N - No. of Negative tuples

P' - No. of tuples that were labelled as positive.

N' - No. of tuples that were labelled as negative.

All - Total no. of tuples.

Sensitivity : TP/P

Specificity : TN/N

Classifier Accuracy

$$= \frac{TP+TN}{P+N}$$

$$\text{Error Rate} = \frac{FP+FN}{P+N}$$

$$\text{Precision} = \frac{TP}{TP+FP}$$

$$\text{Recall} = \frac{TP}{TP+FN}$$

1

F₁/F-Score

$$= \frac{2 \times \text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$

$$= \frac{2 \times \text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$