

COMPUTER ENGINEERING DEPARTMENT

ASSIGNMENT NO-02

SUB: Data Warehousing & Mining

COURSE: T.E.

Year: 2020-2021

Semester: VI

DEPT: Computer Engineering

SUBJECT CODE: CSC603

SUBMISSION DATE: 04/05/2021

Name: Amey Thakur

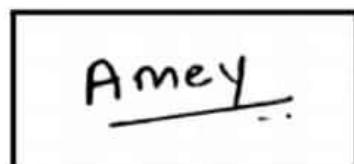
Roll No.: 50

Class: TE Comps B-50

ID: TU3F1819127

Assignment No 2

Sr. No.	Question	CO mapping
1	Differentiate between Classification, Prediction and Clustering.	C01
2	Explain Accuracy and Error measures for Multiple linear regression Model Evaluation & Selection.	C01
3	Explain the Apriori algorithm with an example.	C02
4	Explain Mining Multilevel Association Rules, Multidimensional Association Rules in detail.	C02
5	Differentiate between Spatial Vs. Classical Data mining.	C03
6	Write a short note on types of web mining.	C03


Amey

Q1. Differentiate between Classification, Prediction and Clustering.

Ans:

Classification	Prediction	Clustering
In classification, a training set containing data that have been previously categorised and based on this training set, the algorithm finds the category that the new data point belong to.	It is used to access the values or value ranges of an attribute that a given sample is likely to have.	The characteristics of similarity of data is not known in advance so using statistical concepts we split datasets into sub-datasets such that they have similar data called as clusters
Classification is supervised learning	Both Classification and Clustering are used to make predictions.	Clustering is unsupervised learning.
You're given an unseen tuple and you are suppose to set a label or a class to that tuple.	Construction and use of a model to access the class of an unlabelled sample.	You are given a set of transaction history that gives us details about which customer bought what item.
Example: - Group patients based on their known medical data and treatment outcome then its a classification	Example: - If a classification model is used to predict the treatment outcome for a new patient then it would be a prediction	Example: - By clustering techniques you can tell the segmentation of your customer.

Q.2. Explain accuracy and Error measures for multiple linear regression model evaluation and selection.

Ans:

Accuracy and Error Measures

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

Basic Concepts

Partition of data randomly into three sides.

Training Set:

- It is a set of instances that have not been used in training process. The models performance is evaluated on unseen data testing just estimate the probability of success of unknown data.

Validation Data:

- It is used for parameters tuning but it can not be the test data. Validation data can be training or a subset of training data.

Generalization error: Model's error on test data

- Success: instance class is predicted correctly.
- Error: instance class is predicted incorrectly.

The confusion matrix

- It is useful tool for analyzing how well your classifier can recognize types of different classes.

- TP: Class member which are classified as class member.
- TN: Class non-member which are classified as non-member.
- FP: Class non-member which are classified as member.
- P: Number of positive tuples.
- N: Number of negative tuples.
- N^+ : Number of tuples that were labelled negatively.
- P' : Number of tuples that were labelled positively.

Sensitivity.

- True positive recognition rate which is the proportion of positive tuples that are correctly identified sensitivity TP/P .

Specificity

- True negative recognition rate which is the proportion of negative tuples that are correctly classified.

$$\text{Specificity} = \frac{TN}{N}$$

Classifier accuracy on recognition rate.

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

$$\text{Accuracy} = \text{Sensitivity} \frac{P}{(P+N)} + \text{specificity} \frac{N}{(P+N)}$$

$$\text{Error rate} = 1 - \text{Accuracy}$$

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

AMEY

B

50

Amey -

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

F measure

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

$$F_p = \frac{(1 + \beta^2) \times \text{Precision} \times \text{Recall}}{\beta^2 \times \text{Precision} + \text{Recall}}$$

β = non-negative real number.

Q3. Explain the apriori algorithm with example.

Ans:

Apriori Algorithm

- The Apriori algorithm uses frequent itemsets to generate association rules, and it is designed to work on the dataset that contains transactions.
- With the help of association rule, it determines how strongly or how weakly two objects are connected.
- This algorithm uses a breadth-first search and Hash tree to calculate the itemset associations efficiently.
- It is the iterative process for finding the frequent itemsets from large dataset.

Steps for Apriori Algorithm

Step 1 - Determine the support of itemsets in the transactional database, and select the minimum support and confidence.

Step 2 - Take all supports in the transaction with higher support value than the minimum or selected support value.

Step 3 - Find all the rules of these subsets that have higher confidence value than the threshold or minimum confidence.

Step 4 - Sort the rules as the decreasing order of lift.

Example :

TID	ITEMS
T1	L ₁ , L ₂ , L ₅
T2	L ₃ , L ₄
T3	L ₂ , L ₃
T4	L ₁ , L ₂ , L ₄
T5	L ₁ , L ₃
T6	L ₂ , L ₃
T7	L ₁ , L ₃
T8	L ₁ , L ₂ , L ₃ , L ₅
T9	L ₁ , L ₂ , L ₃

minimum support Count = 2

minimum confidence = 60%

Step 1: K = 1

Create a table containing support count of each item present in dataset called (candidate step)

Itemset	Support Count
L ₁	6
L ₂	7
L ₃	6
L ₄	2
L ₅	2

Compare candidate set items support count with minimum support count.

This gives itemset L₁.

Step 2 $\leftarrow k=2$

Generate candidate set C_2 using L_1

- Check all subsets of an itemset are frequent or not and if not frequent remove that itemset.
- Now find support count of three itemsets by searching in dataset.

itemset	Support Count
L_1, L_2	4
L_1, L_3	4
L_1, L_4	1
L_1, L_5	2
L_2, L_3	2
L_2, L_4	2
L_2, L_5	2
L_3, L_4	0
L_3, L_5	1
L_4, L_5	0

Compare candidate (C_2) support count with minimum support count that gives itemset L_2

itemset	Support Count
L_1, L_2	4
L_1, L_3	4
L_1, L_5	2
L_2, L_3	2
L_2, L_4	2
L_2, L_5	2
L_4, L_5	2

Step 3:

Similarly

itemset	Support Count
L_1, L_2, L_3	2
L_1, L_2, L_4	2

Compare C₃ this gives L₃

itemset	Support Count
L_1, L_2, L_3	2
L_1, L_2, L_5	2

We stop here because no frequent itemsets are found further.

Confidence:

$$\text{Confidence } (A \rightarrow B) = \frac{\text{Support count } (A \cup B)}{\text{Support count } (A)}$$

Step 4: Association rules

$$[L_1 \wedge L_2] \Rightarrow [L_3] \quad || \quad 2/4 * 100 = 50\%$$

$$[L_1 \wedge L_3] \Rightarrow [L_2] \quad || \quad 2/4 * 100 = 50\%$$

$$[L_2 \wedge L_3] \Rightarrow [L_1] \quad || \quad 2/4 * 100 = 50\%$$

$$[L_1] \rightarrow [L_2 \wedge L_3] \quad 2/6 * 100 = 33\%$$

$$[L_2] \rightarrow [L_1 \wedge L_3] \quad 2/7 * 100 = 28\%$$

$$[L_3] \rightarrow [L_1 \wedge L_2] \quad 2/6 * 100 = 33\%$$

So if minimum confidence = 50%

Then first 3 rules can be considered as strong association rules

Q4. Explain mining of multilevel association rule, multidimensional Association Rule.

Ans:

Multilevel Mining Association rules

- Items often form hierarchy.
- Items of the lower has lower support.
- Rules which contain associations with hierarchy of concepts are called multilevel association rules.

Department

Food stuff

Sector

Frozen

Refrigerator

Fresh

Bakery

Etc

Family

Vegetable

Fruit

Dairy

Etc

Product

Bananas

Apple

Orange

Etc

(i) Using uniform support level for all levels

- The same minimum support for all levels
- There is only one minimum support threshold no need to examine itemsets
- If support threshold is too low \rightarrow generate too many high level association.

Level 1

Milk

Support = 10%

Minimum support = 5%

2% Milk

Slim milk

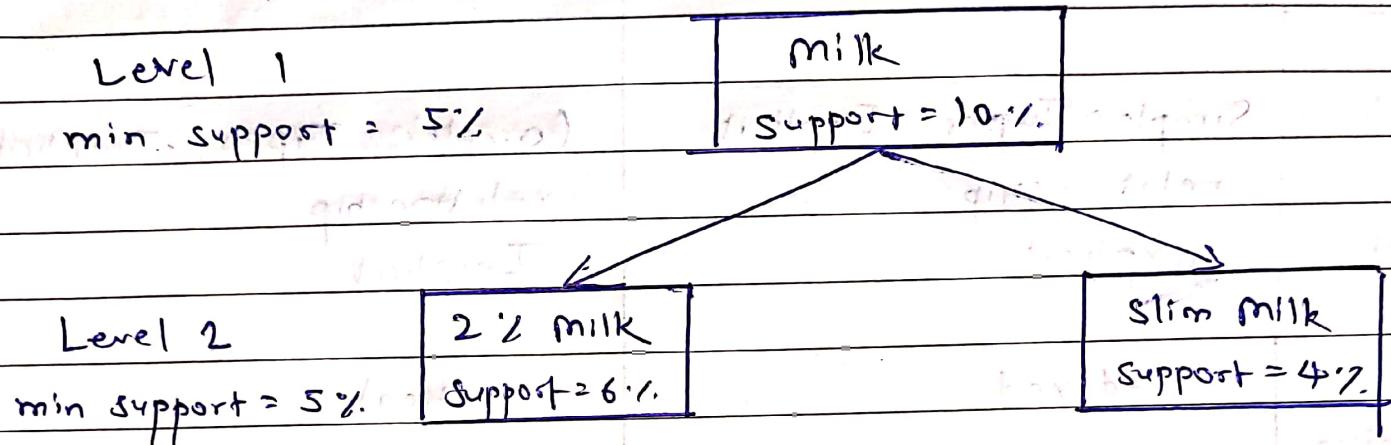
Support = 6%

Support = 4

min Support = 5%

② Using reduced minimum support level for all levels

- At every level of abstraction, there is its own minimum support threshold.
- So minimum support at lower level reduced.



Multidimensional Association Rule.

① In multidimensional association,

- Attribute can be categorical or quantitative
- Quantitative attributes are numeric and incorporates hierarchy
- Numeric attributes must be discretized.
- It consists of more than one dimension.

② Three approaches in mining multi-dimensional association rule.

- ① Using static discretization of quantitative attributes
- ② Using dynamic discretization of quantitative attributes
- ③ Using distance based discretization with clustering

① Using static discretization of quantitative attributes

- Discretization is static and occurs prior to mining.
- Discretized attributes are treated as categorical.
- Use apriori algorithm to find all k-frequent predicate sets. (This requires 'k' or 'k+1' table scans.)
- Every subset of frequent predicate set must be frequent.
- Eg: If in a data cube the 3D cuboid (age, income, buys) is frequent implies (age, income), (age, buys), (income, buys) are also frequent.
- Data cubes are well suited for mining since they make mining faster.
- The cells of an n dimensional data cuboid correspond to the predicate cells.

② Using dynamic discretization of quantitative attribute.

- Known as mining quantitative association rules.
- Numeric attributes are dynamically discretized.
- Eg: age($x, "20.25"$) & income($x, "30k.41.k"$) buys($x, "Laptop Computer"$)

	Age = 20	Age = 21	Age = 22	Age = 23
income, 38 to 41				
income, 34 to 37				
income, 30 to 33				

Grid for Tuples.

③ Using distance based discretization with clustering

- It involves a two step mining process
 - Perform clustering to find the intervals of attributes involved.
 - Object association rules by searching for groups of clusters that occur together.
- The resultant rules may satisfy
 - Clusters in the antecedent are strongly associated with clusters of rules in the consequent.
 - Clusters in the antecedent occurs together.
 - Clusters in the consequent occurs together.



Q5. Differentiate between spatial vs. Classical Data mining

Ans

Parameters	Classical Data Mining	Spatial Data Mining
Data Definition	Simple	Complex
Relationships	Explicit	Implicit
Data organization	Indexed	Vertical
Statistical Foundation	Independence of samples	Spatial auto correlation
Output	Set based	Spatial based
Example	Classification Accuracy	Spatial Accuracy
Computation Process	Combination and optimization numerical algorithm	Computational efficiency opportunity. Plane sweeping Spatial auto correlation

Q.6 Write a short note on types of web mining.

Ans:

Web mining

- Web mining is the process of data mining technique to automatically discover and extract information from web documents and services.
- Web mining can be broadly divided into 3 types.

① Web content mining

- It is the application of extracting useful information from the content of web documents.
- Web content consists of several types of data: text, images, audio, video, etc.
- Content data is the group of facts that a web page is designed.
- It can provide effective and interesting patterns about user needs.
- Text documents are related to text mining, machine learning and Natural Language Processing

② Web structure mining

- It is the application of discovering structure information from the web.
- The structure of the web graph consists of web pages as nodes and hyperlinks as edges connecting related pages.
- Structure mining basically shows the structured summary of a particular website.
- It identifies relationship between web pages linked by information or direct link connection.
- To determine the connection between two commercial websites, web structure mining can be very useful.

AMEY

B 50

Amey

Page No.:	100
Date:	youva

③ Web usage Mining

- It is the application of identifying or discovering interesting usage patterns from large datasets and these pattern enables us to understand the user behaviour or something like that.
- In usage mining, user access data is form so usage mining this is also called mining