

# Assignment No. 1

## \* Problem Statement:

for an organization, choose a set of business processes. design Star/Snowflake Schema from analyzing these processes. Create a fact constellation Schema by combining them. Extract data from diff. data sources, apply suitable transactions and load into designation table using ETL tool.

## \* Objective :

- i.) understanding the basics of Star/ snowflake/ fact constellation Schema.
- ii.) Learn how to install & use Pentaho tool.

## \* Outcome:

- i.) Pentaho was successfully installed
- ii.) Process Schema Set was decided & Schema
- iii.) ETL was performed on dataset in Pentaho

## \* Theory:

- ETL Stands for Extract, Transform & Load.  
It extracts data from diff. RDBMS source system, transform data by applying calculation & concatenation then load data into data warehouse system. It is loaded in the form of dimension & fact table.

### - Extraction

A Staging area is req. during ETL load for various reasons.

- ① The Source System are only available for a specific period to extract data. This period of time is less than data load time. This Staging area allows to extract data from Source system & keep it in Staging area before the time slot.
- ② Staging area is also req. when you want to get data from multiple source together.

### - Transformation

In data transformation, you apply a set of functions on extracted data to load it into the target system.

Data which does not req. any transform is known as direct move pass thru data.

Diff transformations on extracted data from the Source System can apply, customised to the business decided.

### - Load

In this phase, data is loaded into the end target system, either a flat file or a Data warehouse.

### ① Star Schema ⇒

- Every dimension represented with only one dimension table.
- dimension table should contain a set of attributes
- dimension table joined to fact table using foreign
- dimension table are not joined to each other
- fact table contains key & measure
- easy to understand, widely supported by BI tools
- optimal disc usage, denormalized structure, faster queries
- dimension tables are not normalized
- cube processing is faster

### ② Snowflake Schema is an extension of Star Schema, & it adds additional dimensions

- dimension table normalized, which splits data into additional tables
- uses smaller disc space
- reduced query performance due to additional table
- complex database design
- low level data redundancy.

### ③ Pentaho data integration provides the ETL capabilities that facilitates the process of consuming, cleaning & storing data using a uniform & consistent format that is accessible & relevant to end users.

\* This tool enables users to -

- 1) load huge datasets into databases
- 2) clean data with steps ranging from very simple to very powerful complex transformations
- 3) data warehouse populations , and more

\* The PDI client 'Spoon' is a visual, drag & drop, no-code interface that is used to perform ETL.

\* Conclusion:-

- 1) Pentaho Data integration tool was successfully installed & used on a Sales dataset to successfully Extract, transform & load in MySQL.
- 2) A business was picked & a set of business processes developed & represented using Snowflake Schema.

## Assignment No. 2

### \* Problem Statement:

Consider a Suitable dataset, for clustering of data instances in diff. groups, apply diff. clustering techniques (min. 2) · Visualize the clusters using suitable tools.

### \* Objective :

- i) understand various clustering types & how to implement the same.
- ii) use python libraries & appropriate data sets to perform clustering & visualize the same.

### \* Outcome:

Understood k-means, Hierarchical clustering & perform it on dataset cars data by brand.

### \* Theory:

- ① Clustering or cluster analysis is a task of grouping a set of objects in such a way that objects in the same group (cluster) are more similar to each other than to those in each other.
- ② Clustering is the main task of exploratory data mining & a common

technique for statistical data analysis, and is used in many fields.

- ① K-means clustering is a type of unsupervised learning, which is used when you unlabeled data (i.e. data without defined categories or groups). the goal of this algorithm is to find groups in the data, with the number of groups represented by variable k. the algorithm works iteratively to assign each data point to one of k groups based on the feature provided. data points are clustered based on feature similarity. the result of this algorithm is -

⇒ the centroid of k-clusters, can be used to label new data

⇒ labels from the training data (each data point is assigned to single cluster)

- ② Rather than defining groups, before looking at the data, clustering allows you to find & analyze the groups that have formed organically.

- ③ Hierarchical clustering is an algorithm that groups similar objects into groups called cluster. the end point is a set of cluster, where each cluster is distinct from each other, and the object

written each cluster are broadly similar to each other.

① Given a set of  $N$  items to be clustered and an  $N \times N$  distance matrix, the basic process of hierarchical clustering is:

- assign each item to its own cluster,  
So now there are  $N$  clusters; let the dist. between the cluster equal the dist. (similarities) between the items they contain.
- find the closest (most similar) pair of clusters and merge them into a single cluster
- compute the distances (similarities) between the new cluster & each of the old cluster
- Repeat above 2 steps until all clusters are clustered into a single cluster of size  $N$ .

② The default distance measured is the Euclidean dist. which is the Sq. root of the sum of the differences.

③ In the agglomerative clustering approach, there are 4 possible methods, Ward's method being one of them. It says that the dist. between two clusters, A & B, is how much the sum of squares will increase when they are merged.

- ① The dataset used was 'mall-customers' which contains various parameters.
- ② these entries were clustered using the algorithm.

### \* Conclusion:

K-means and hierarchical (agglomerative) clustering techniques were understood, successfully implemented, & the req. output was obtained.

## Assignment No. 3

### \* Problem Statement:

Apply a-Priori algorithm to find freq. occurring items from given data and generate strong association rules using support & confidence thresholds.

### \* Objective : Model associations between products by determining sets of item freq. purchased together & building association rules to drive recommendations.

### \* Outcome : Demonstrated market basket analysis using a-priori algorithm to find freq. occurring items from given data & generate strong association rules using support & confidence algo.

### \* Theory :

Association rule mining finds interesting associations and relationships between large sets of data items . this rule shows how frequently an item occurs in a transaction . It is defined as an implementation expression of the form  $x \rightarrow y$  where  $x$  &  $y$  are any 2 item sets.

Market Basket analysis is one of the key techniques used by large retailers to show association between items. It allows retailers to identify relationships between the items that people buy together freq, so that items can be strategically placed next to each other to boost sales.

Some definitions:

- ↳ Support count: freq. of itemset occurrence
- ↳ freq. Itemset: Itemset whose support  $\geq$  minsup threshold

The association rule evaluation metrics are

- ① Support (S): The number of transactions that include items in the  $S \times S$  part of the rules as percentage of total number of transaction.
- ② confidence (C): It is the ratio of number of transactions that includes all items in an itemset  $S \times S$  as well as the transaction that include all items in  $S \times S$  to the number of transaction that include all items in  $S \times S$ .
- ③ Lift (L): The lift of the rule  $x \rightarrow y$  is the confidence of the rule divided by

the expected confidence, assuming  $S_{X1} \& S_{Y1}$  are independent of each other. Expected confidence is confidence divided by freq.  $S_{Y1}$ . Itemset occur together

- (i) as expected if  $\text{lift} = 1$
- (ii) more than expected if  $\text{lift} > 1$  and
- (iii) less than expected if  $\text{lift} < 1$

### Apriori Algorithm

→ It is an algorithm for freq. itemset mining & association rule learning over relational databases, first proposed by Agrawal & Srikant in 1994.

It is designed to operate on databases containing transactions, with each transaction being a set of items (itemsets).

Given a threshold C, the algo identifies the itemsets, which are subsets of atleast C transactions in database

It uses a bottom-up approach, where frequent subsets are extended one at a time (called candidate generation), and groups of candidate are tested against the data. The algo terminates when no other successful extension are found.

Apriori used BFS & a hash tree to count candidate item sets efficiently. It generates K-itemsets of length  $K=1$ ; then it prunes the candidate, which have

classmate

Date \_\_\_\_\_  
Page \_\_\_\_\_

in-frequent sub pattern. After that it scans the transaction database, to determine freq. item sets among the candidate.

\* Conclusion:

Thus frequently occurring items from given market basket dataset and strong association rules using support & confidence thresholds found using a-priori algorithm.

# Assignment No. 4

## \* Problem Statement :

Consider a suitable text dataset. Remove Stop words, apply Stemming & feature Selection techniques to represent documents as vector, classify documents & evaluate Precision, recall.

## \* Objective : Learn how to tokenize & filter a document into its diff. words and then do word count for each word in doc. Apply stemming, feature Selection on doc (Ex)

## \* Outcome : demonstrated text processing using nltk, understood vectorizing, removing Stop words.

## \* Theory :

① NLP is a Subfield of Linguistic computer Science & AI concerned with interaction b/w computer & Human language, in particular how to program computer to process & Analyse large amt of Natural language dt.

② In computing Stop words are words that are filtered out before or after the natural language data (text) are

processed. While 'stop words' typically refer to the most common word, there is no universal list of stop words.

① STEMMING  $\Rightarrow$  for grammatical reasons, documents are going to use diff. forms of a word (e.g. organized, organizing; etc.). Additionally, there are families of derivitatively related words with similar meaning like democracy, democresy & democratization.

$\rightarrow$  The goal of stemming (and lemmatization) is to reduce inflectional forms & sometimes derivationally linked forms of a word to its common base form.

am, are, is  $\rightarrow$  be

car, cars, car's, car's  $\rightarrow$  car

$\rightarrow$  when applied to a doc, result will be somewhat like this

the boy's cars are diff. color )

the boy car be differ color ↴

$\rightarrow$  Stemming is a more cuttable process, a heuristic that chops off the end of a word in the hope of achieving the goal correctly more of time, and often includes removal of derivational affixes.

$\rightarrow$  Feature selection is a process of selecting a subset of the terms occurring

in the remaining set and using only this subset as feature in text classification.

- It serves 2 purposes. first, it makes training & applying a classifier more efficient by decreasing the size of vocab. this is of particular importance to classifier that are expensive to train.
- Second, feature Selection often increases classification accuracy by eliminating noise features. These features when added to docs representation increases the classification error on new data.

→ Vectorization is a process of converting text into machine readable form. words are represented as 'vectors' (numerically)

- Count vectorizer (one-hot encoding) involves counting the number of occurrence for each word occurring in the doc.
- Idea behind it is simple. vector is created having as many dimensions as there are distinct words in text/doc/collection of docs being used. each unique word has a unique dimension & is represented by a 1 in the dimension with 0s everywhere else. This results in huge & sparse vectors that capture no relational data.

- If IDF vectors are related to one hot encoded vectors, but instead of just feature a count, they feature numerical representation, where the words aren't just present or not present instead, they are expressed by their term freq. multiplied by inverse doc. freq.
- In simpler terms, words that occur everywhere should be given very little weight on significance; coz they don't provide a large amt. of value. However if a word appears very little or freq frequently but only in specific place, then they are possible of higher significance.
- The downside is there is no capture of Semantic relatedness, this is solved with the co-occurrence matrix and a ~~stochastic~~ probabilistic model.



#### Conclusion :

Text docs were tokenized, filtered, vectorized, & thus successfully processed. basics of NLP were understood.