

Practical - 6

Aim:- WAP a program for Apriori Algo in Python

Source code:-

```
import sys
from itertools import chain, combinations
from collections import defaultdict
from optparse import OptionParser
```

```
def subsets(arr):
```

```
    """Returns non empty subsets of arr"""
```

```
    return chain(*[combinations(arr, i+1) for i,
                    a in enumerate(arr)])
```

```
def returnItemwithminsupport (itemset, transactionlist,
                               minsupport, freqset):
```

```
    itemset = set()
```

```
    localset = defaultdict(int)
```

```
    for item in itemset:
```

```
        for transaction in transactionlist:
```

```
            if item.issubset(transaction):
```

```
                freqset[item] += 1
```

```
                localset[item] += 1
```

```
    for item, count in localset.items():
```

```
        support = float(count) / len(transactionlist)
```

```
        if (support >= minsupport):
```

```
            item itemset.add(item)
```

```
    return itemset
```

```
def joinset (itemset, length):
    return int ([i, union(j) for i in itemset for j
                in itemset if len([i, union(j)])
                == length])
```

```
def getItemSet TransactionList (data_iterator):
    transactionList = list()
    itemset = set()
    for record in data_iterator:
        transaction = for frozenset(record)
        transactionList.append(transaction)
        for item in transaction:
            itemset.add(frozenset([item]))
    return itemset, transactionList
```

```
def runApriori (data_iter, minsupport, minconfidence):
    itemset, transactionList = getItemSet Transaction
    List (data_iter)

    freqset = defaultdict (int)
    largeset = dict()
    assocrules = dict()
    onceset = return item with min support (itemset,
        transactionList, min
        min support, freqset)

    currentset = One (set
```

```
k = 2
while (currentset != set ()):
    largeset [k-1] = currentLset
    currentLset = joinset (currentLset, k)
    currentset = return item with min support
        (currentLset, transactionList, min support,
        freqset)
```


currentLset = currentCset
k = k + 1

def getSupport(item):

return float(freqset[item]) / len(transactionList)

to RetItems = []

for key, value in LargeSet.items():

to RetItems.extend([tuple(item), getSupport(item)]
-- for item in value])

to RetRules = []

for key, value in LargeSet.items()[1:]:

for item in value:

subsets = map(frozenset, [x for x in subsets
(item)])

for element in subsets:

remain = item.difference(element)

if (item(remain) > 0):

confidence = getSupport(item) / getSupport(element)

if confidence >= min confidence:

to RetRules.append(((tuple(element),
tuple(remain)), confidence))

return to RetItems, to RetRules


```
def printResults(items, rules):
```

```
    for item, support in sorted(items, key = lambda  
                                (item, support):
```

```
        print "item: %s, %3f" % (str(item), support)
```

```
    for rule, confidence in sorted(rules, key =  
                                lambda(rule, confidence): confidence):
```

```
        pre, post = rule
```

```
        print(str(pre), str(post), confidence)
```

```
def getDataFromFile(filename):
```

```
    file_iter = open(filename, 'r')
```

```
    for line in file_iter:
```

```
        line = line.strip().rstrip(',')'
```

```
        record = line.split(',')
```

```
        yield record
```

```
if name == '__main__':
```

```
    optparser = optparse.OptParser()
```

```
    optparser.add_option('-f', 'inputFile',  
                        (dest='input'),
```

```
                        help='filename containing csv', default  
                        = None
```

```
    optparser.add_option('-s', 'min support', (dest='min s',
```

```
                        help='minimum support value', default=0.15,  
                        type='float')
```

```
    optparser.add_option('-c', 'min confidence', (dest='min c',
```

```
                        help='minimum confidence level', default=0.6,  
                        type='float')
```

```
optparser.addoption('-c', 'min confidence', dest='minc',  
help='minimum confidence level', default=0.6,  
type='float')
```

```
(options, args) = optparser.parse_args()  
infile = None
```

```
if options.infile is None:
```

```
    infile = sys.stdin
```

```
elif (options.infile is not None):
```

```
    infile = dataFormatFile(options.infile)
```

```
else:
```

```
    print 'No dataset filename specified, system  
    with exit \n'
```

```
sys.exit("system will exit")
```

```
minSupport = options.minS
```

```
minSupport = options.minC
```

```
items, rules = runApriori(infile, minSupport,  
                           minConfidence)
```

```
printResults(items, rules)
```


Practical 7

Aim: Comparison b/w the classification results by J48 decision tree based method and naive Bayesian method

• Performance criteria for evaluating the classifiers:

There are different criteria for evaluating the performance of the classifiers:-

- i) classification Accuracy
- ii) specificity
- iii) sensitivity recall
- iv) Precision

i) (classification Accuracy:- It is the ability to predict categorical class labels. This is the simplest scoring measures. It calculates the proportion of correctly classified instances.

$$\text{accuracy} = (TP + TN) / (TP + TN + FP + FN)$$

ii) Specificity:- It relates to the classifiers' ability to identify negative rules.

$$\text{specificity} = TN / (TN + TP)$$

iii) Sensitivity / Recall: - Sensitivity is the proportion of actual positivities which are correctly identified positive by classifiers

$$\text{Recall} = TP / (TP + FN)$$

iv) Precision: - This is the measure of retrieved instances that are relevant

$$\text{Precision} = TP / (TP + FP)$$

True Positive (TP): If the instance is positive and is classified as positive

False Negative (FN): If the instance is negative but it is classified positive

False Positive (FP): If the instance is negative but it is classified positive.

• Confusion Matrix: - In the field of machine learning, and specification, a confusion matrix, also known as an error matrix, is a layout that allows visualization of the performance of an algorithm, typically a supervised learning

Comparison b/w Decision tree & Naive Bayes
for diabetes Dataset:-

Confusion Matrix in case of J48:-

a	b ← classified as	
407	93	a = tested-positive
108	160	b = tested-negative

$$\begin{aligned}\text{Accuracy} &= \frac{TP + TN}{TP + TN + FP + FN} \\ &= \frac{407 + 160}{407 + 160 + 93 + 108} = \frac{567}{768} = 0.7382\end{aligned}$$

Confusion matrix for Naive Bayes:-

a	b ← classified as	
422	78	a = tested positive
104	164	b = tested negative

$$\text{Accuracy} = \frac{TP + TN}{TP + FP + FN + TN} = \frac{586}{768} = 0.7630$$

$$0.7630 > 0.7382$$

Therefore, accuracy of Naive Bayes is better than as compared to J48 for given diabetes dataset.

Practical - 8

Aim:- Demonstration of k-means clustering using weka tool.

Clustering:- Clustering is a collection of data object with the properties -

- a) Similar to one another within the same cluster
- b) Similar to the object in other cluster.

Cluster Analysis:- Cluster analysis is finding similarities b/w data according to the characteristics found in the data and grouping similar data object into clusters

Major clustering approaches are:-

- 1) Partitioning 2) Hierarchical 3) Density based
- 4) Grid based 5) Model based

K-means Clustering:- K-means is one of the simplest unsupervised learning algorithm that solve the well known clustering problem.

The procedure follows a simple and easy way to classify a given dataset through a certain number of cluster (assume k clusters) fixed a priori. The main idea is to define k centroids, one for each cluster. These centroids should be placed in a cunning way because of different location causes different result. So, the choice is to place them as much possible far away from each other.

The next step is to take each point belonging to a given dataset and associate it to the nearest centroid.

A loop has been generated. As a result of this loop we may notice that the k -centroids change their location step by step until no more changes are done. In other words centroids do not move. Finally this algorithm aims at minimizing an objective function, in this case a squared error function. The objective function

$$J = \sum_{j=1}^k \sum_{i=1}^n \|x_i(j) - c_j\|^2$$

where $\|x_i(j) - c_j\|^2$ is a chosen distance measure b/w a data point $x_i(j)$ and the cluster center c_j is an indicator of the distance of the n data points from their respective cluster centers.

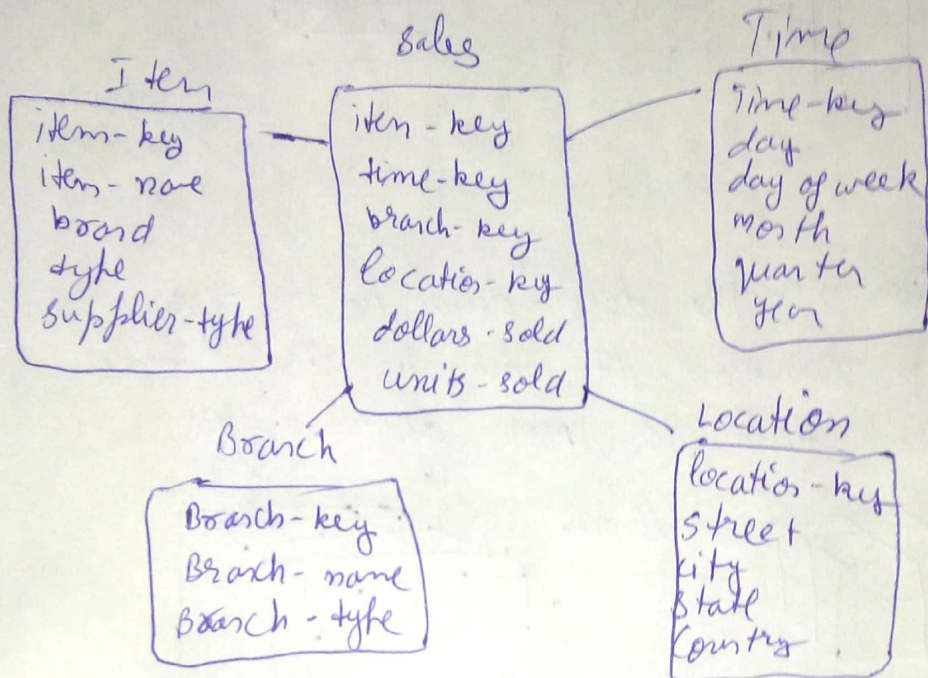
Algorithm :-

- 1) Place k points into space representing by the objects that are being clustered. These points represent initial group centroids
- 2) Assign each object to the group that has the closest centroid
- 3) When all objects have been assigned, recalculate positions of k centroids
- 4) Repeat steps 2 & 3 until the centroids no longer move.

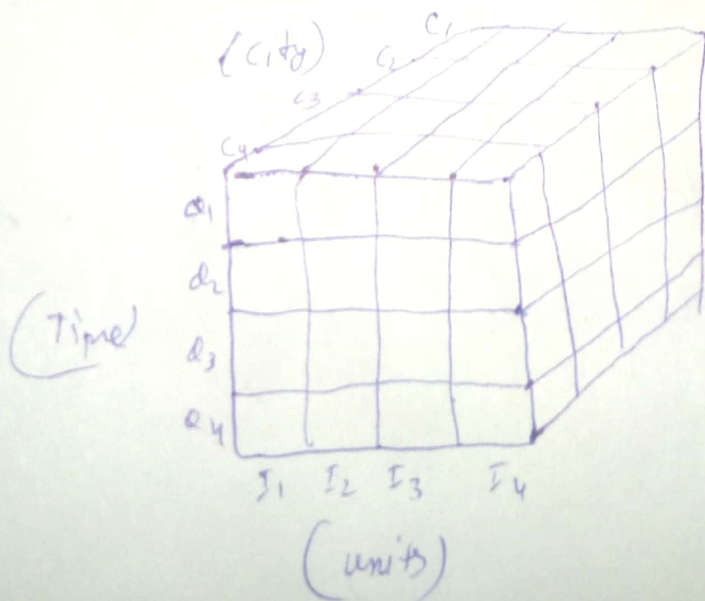
Practical-9

Aim:- Demonstration of OLAP operations on multidimensional data with example.

Let us perform OLAP operations on the given Sales schema

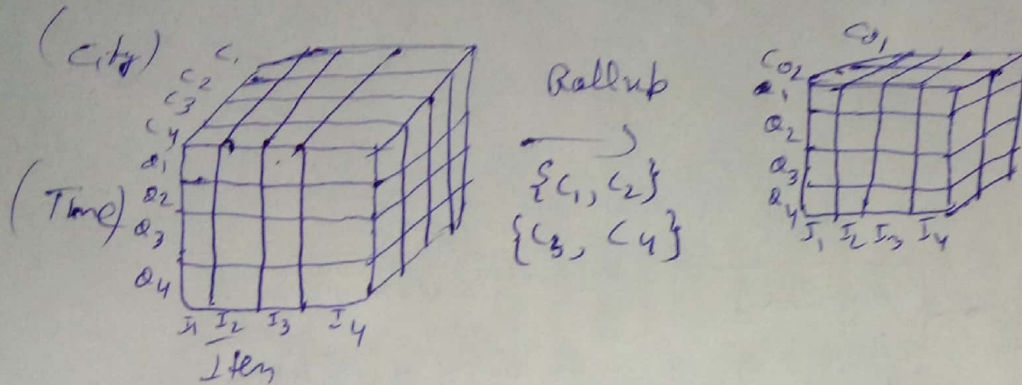


The data cube can be represented as:-

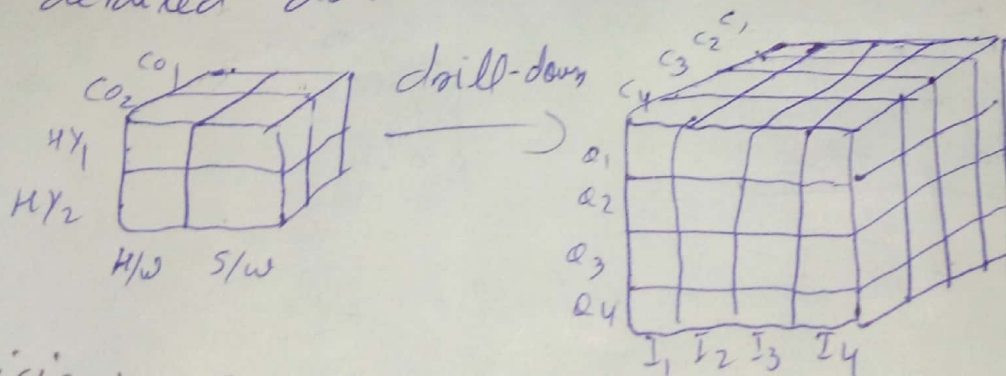


Various OLAP operations are:-

1) Roll up:- The roll up operation performs an aggregation on data cube either by climbing up a concept hierarchy for a dimension.



2) Drill Down:- The drill-down is the reverse of roll up operation. It navigates from less detailed data to more detailed data.



3) Slicing:- The slice operation performs a selection on one dimension of the given cube

4) Dicing:- The dice operation defines a subcube by performing a selection on one or two dimensions