

NAME: VEDANTH.V. BALIGA

SEC: 5<sup>th</sup> 'H'

USN: ENH20DS0044

COURSE: Data Warehouse And  
KNOWLEDGE Mining.

COURSE CODE:

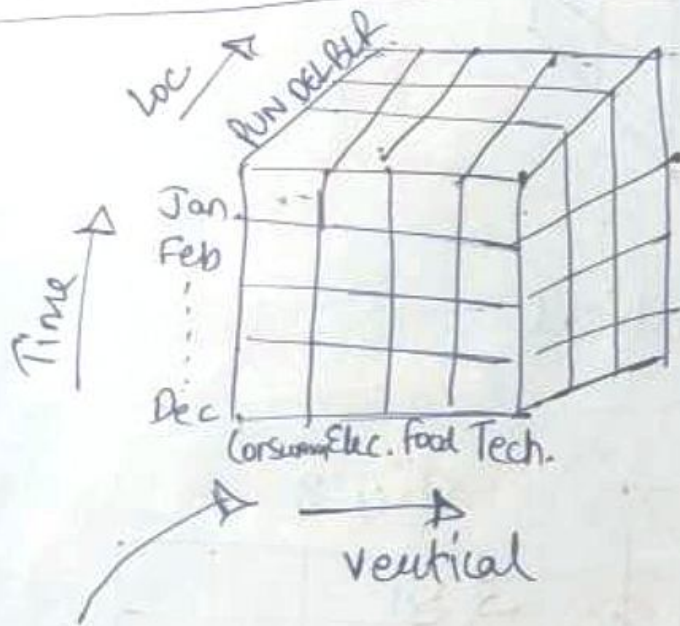
TOPIC: ASSIGNMENT-I

## Module 1

Choose a realtime domain & apply all OLAP operations on it.

A OLAP Stands for Online Analytical Processing.

OLAP operations on Retail data

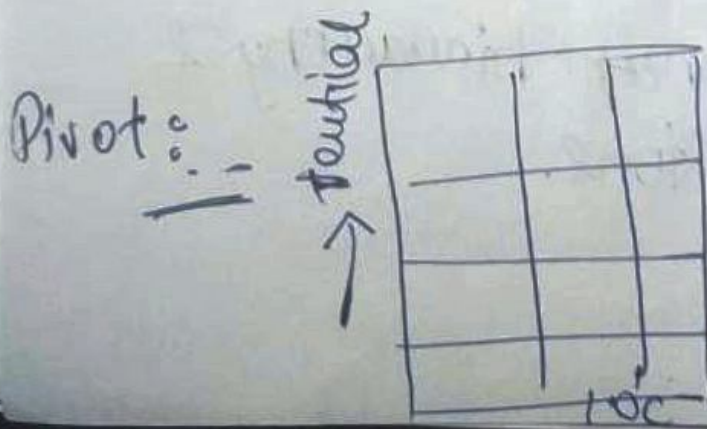
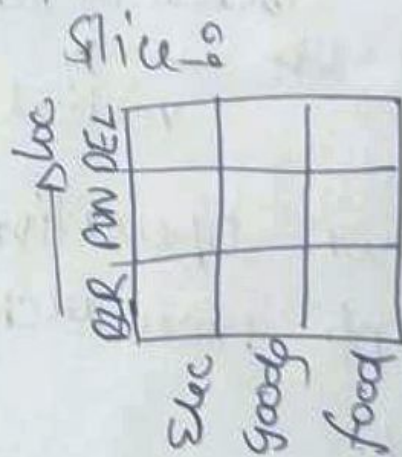
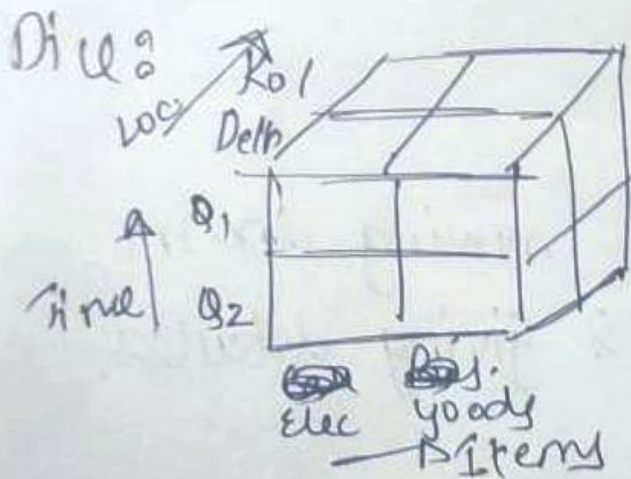
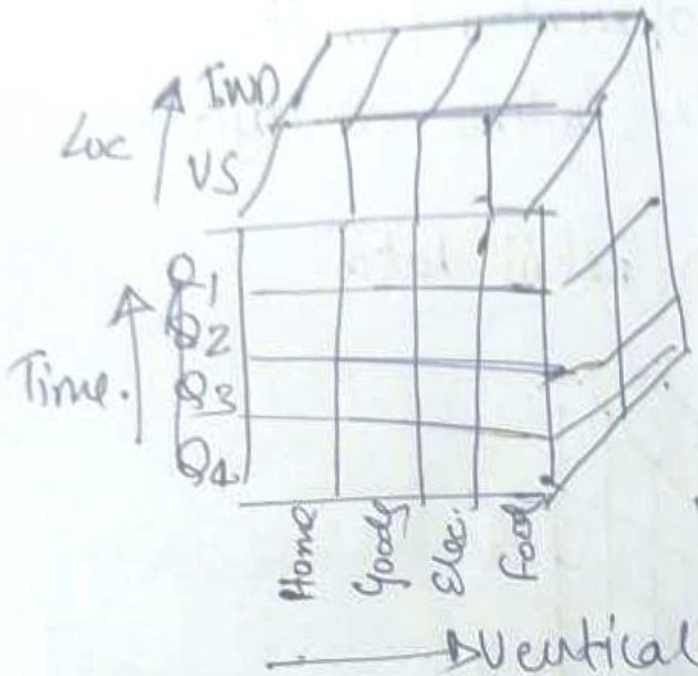


Drill Down operation : moving down a concept hierarchy & giving detailed summary.

Roll up : climb up concept hierarchy & reduce the dimensions.

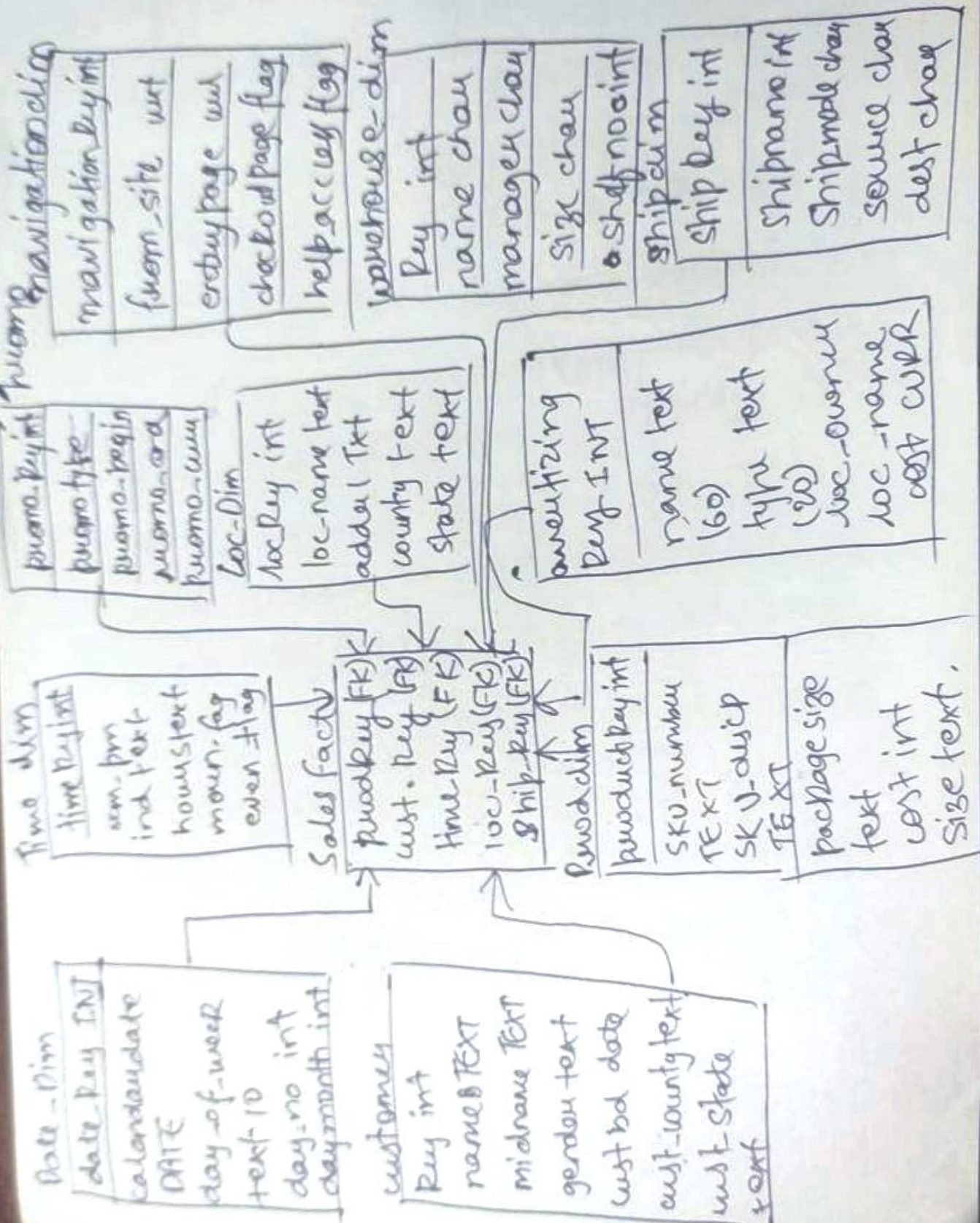


⑤ Design multidimensional multiterel pattern mining for Alpkant.



Verticals

Design a multidimensional data model for Amazon.





Design database for parallel processing.

Parallel execution: break down a task into subtasks & each task performed by many processors.

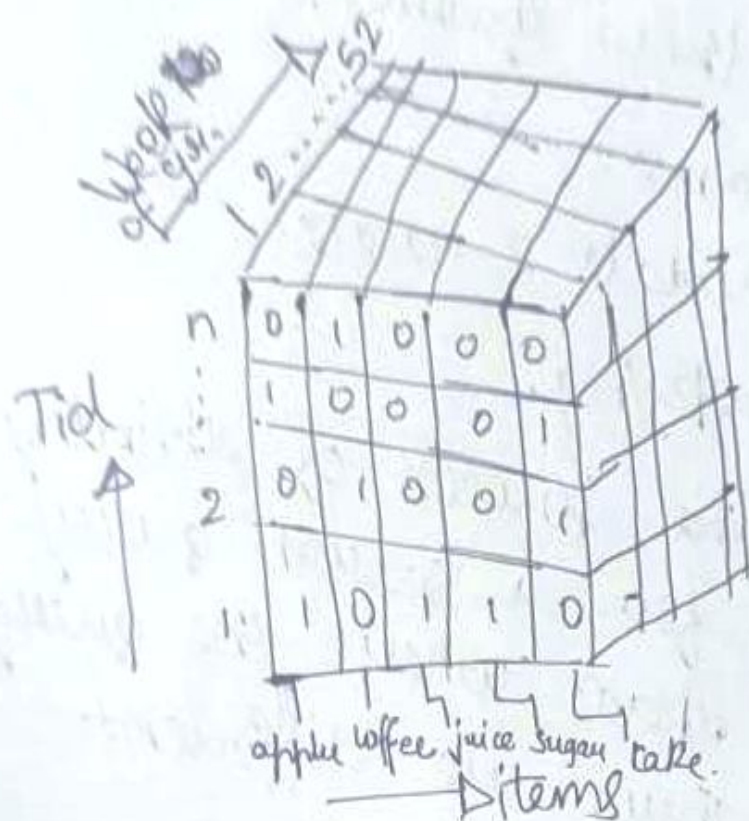
When to execute these queries.

- Large table scans
- Bulk insert, update & deletes
- aggregating & copying.

Partitioning is the process of dividing a large table for a search query to a smaller space where the query may be found. There are different types of query partitioning:

- Hash Partitioning
- List Partitioning.
- Composite Partitioning.
- Index Partitioning.

the folder "SQL Partitioning" in the zip code in Q3) module, file. The model is came as Design Data Cube for market basket analysis.



→ The x-axis represents items, y axis represents tid & z axis represents week of yr.

→ we represent week of yr & not here day to save space & the time is not an imp. factor here.

→ For every transaction we have a 1, 0 attribute. 1 represents item.



represented not purchase.  
This technique is called binary value  
indexing to same space.

In SQL the cube is built using CUBE  
operator. more about CUBE operator  
is in the SQL folder of attached zip  
file.

### CUBE OPERATORS

Retail Sales db

Syntax

Select

c1, c2, AGGR(c3)

From

table

Group by cube (c1, c2);

Table inventory table.

warehouse	product	model	quantity
SF	iPhone	6S	50
SF	iPhone	7S	10
SF	iPhone	X	200
SF	Samsung	Galaxy	200
SF	Samsung	Note	100
ST	iPhone	6S	100
ST	iPhone	7	50
ST	iPhone	X	50
ST	Samsung	AS	200
ST	Samsung	8	150

Task 1: Return total inventory of all inventory in warehouses.

SELECT

warehouse,  
sum(quantity)

ORDER BY  
warehouse

FROM Inventory

group by CUBE (warehouse)



Warehouse	QUANTITY
SF	560
SJ	650
(null)	1210.

Task 2: Group quantity by sum of the diff. brands.

SELECT  
warehouse,  
product,  
sum(quantity)

FROM  
inventory

group by  
warehouse, product)

order by  
warehouse,  
product;

warehouse	Product	Sum (quantity)
SF	Samsung	300
SF	iPhone	260
SF	(null)	560
ST	Samsung	350
ST	iPhone	300
ST	Null	650
null	Samsung	650
null	iPhone	560
null	null	1210.



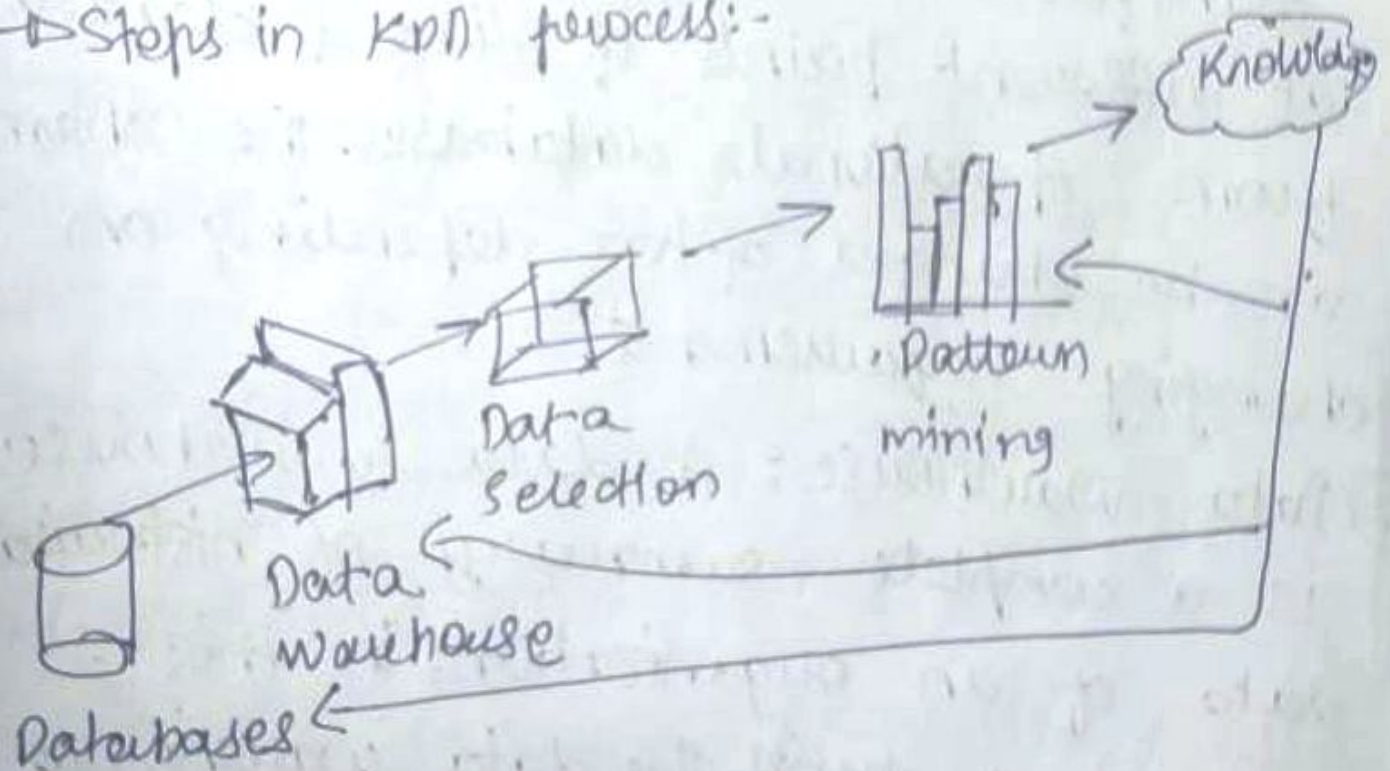
## Module 2

① Design a KDD model for credit card fraud detection.

Q. A What is KDD model?

KDD stands for Knowledge discovery in databases & it involves mining patterns from transactions, scientific & sensor data, pictures, text etc.

→ Steps in KDD process:-



## KDD Process in Credit Card Fraud De

Q.2  
A. ① Databases: A database usually stores data of the current activity & the data that powers the application. In the case of credit card fraud detection, a bank stores data of all its transactions including transaction id, amount & time. This is mostly transferred to data warehouses at different points of time & cleaned from operational databases. The schema can be changed often depending on changing requirements.

② Data warehouse: A data warehouse is a complete summary of historical data of an organization, in this case of a bank. The data warehouse can have different kinds of data features.



Some of them are cc num (customer number), card merchant, category of transaction, date & time, city, state, zip code & population. This data may or may not have missing values, so we need to preprocess it first.

Data Preprocessing: Since our data set is very imbalanced, we have to fix the balance first. On observing distributions of transaction amt & time, we can see that they are very skewed.

Some methods to solve this imbalanced dataset problem is subsampling, random undersampling & SMOTE.

Correlation Analysis: Correlation analysis helps us understand whether there is any relation between the various features of the dataset.

Boxplots & correlation plots help better understand distribution of data.

### Anomaly Detection

IQR: Any data beyond 75th percentile & below 25th percentile is an anomaly.

We should be careful about the threshold we set here to remove outliers.

### Feature Selection

Feature selection can be done on credit card fraud data using multiple algorithms. Some of them are PCA (Principal Component Analysis) where we reduce the data from high-dimensions to lower dimensions preserving as much info as possible. t-SNE preserves the neighbourhood embeddings based on the threshold of distance measured from the



important points.

Pattern mining: Pattern mining is the process of visualizing the different features of the dataset.

Some of the visualizations we can plot are:

- 1] Category of transaction v/s count
- 2] Gender of clients v/s count.
- 3] States / geographical area with most frauds.
- 4] No of credit card frauds by metro city.
- 5] No of credit card frauds by job of transaction.
- 6] How much skewness in the amt v/s density of amount.
- 7] We can use dob column to extract age & then devise age-groups v/s fraud amt.

2) Using date field, we can extract & devise month name v/s count of frauds plot.

Prediction using mined patterns:

There are many classifiers that can be used:

1) Logistic Regression: Since its a two class classfn, we can expect logistic regression to perform well. Other models like SVM & Naive Bayes can also be tested & parameters can be tuned.



Q. Explain the data preprocessing methods for Info Retrieval Applications.

Soln

- ① We can use info method to give statistical description of the data.
- ② missing data can be checked using missingno matrix
- ③ Correlation heatmap allows us to measure nullity correlations
- ④ Dendrogram show the hierarchical nullity relationship between columns. It uses a hierarchical clustering algorithm.
- ⑤ A simple numerical survey of data can show how many types of data counts we have in the dataset.
- ⑥ We can delete data to remove whole attributes or only a sample of attributes or data.

③ Evaluate the statistical descn for stock market analysis with data vizual

A ① Extract Tesla stock market data using BeautifulSoup.

② The data contains the date feature & Revenue column.

③ We can plot the date v/s prices graph

④ Plotting the data of years v/s revenue would give historical data of Tesla Stock.

⑤ Historical revenue of GameStop would be very uneven due to its downfall in 2020.



### Module 3

④ Find the frequent itemsets & strong association rules for the db using Apriori Algorithm. supcount = 2, conf = 50%.

Tid	items
T <sub>1</sub>	I <sub>1</sub> , I <sub>2</sub> , I <sub>5</sub>
T <sub>2</sub>	I <sub>2</sub> , I <sub>4</sub>
T <sub>3</sub>	I <sub>2</sub> , I <sub>3</sub>
T <sub>4</sub>	I <sub>1</sub> , I <sub>2</sub> , I <sub>4</sub>
T <sub>5</sub>	I <sub>1</sub> , I <sub>3</sub>
T <sub>6</sub>	I <sub>2</sub> , I <sub>3</sub>
T <sub>7</sub>	I <sub>1</sub> , I <sub>3</sub>
T <sub>8</sub>	I <sub>1</sub> , I <sub>2</sub> , I <sub>3</sub> , I <sub>5</sub>
T <sub>9</sub>	I <sub>1</sub> , I <sub>2</sub> , I <sub>3</sub>

A 1-itemset.

	sup
I <sub>1</sub>	6/9
I <sub>2</sub>	7/9
I <sub>3</sub>	6/9
I <sub>4</sub>	2/9
I <sub>5</sub>	2/9

1 frequent item set has all items.

### 2-itemset

Item	Support
$I_1 I_2$	4/9
$I_1 I_3$	4/9
(x) $I_1 I_4$	1/9
$I_1 I_5$	2/9
$I_2 I_3$	4/9
$I_2 I_4$	2/9
$I_2 I_5$	2/9
(x) $I_3 I_4$	0
(x) $I_3 I_5$	1/9
(x) $I_4 I_5$	0

### 2-freq itemset

Item	Sup
$I_1 I_2$	4/9
$I_1 I_3$	4/9
$I_1 I_5$	2/9
$I_2 I_3$	4/9
$I_2 I_4$	2/9
$I_2 I_5$	2/9

### 3-freq itemset

Item	Sup
$I_1 I_2 I_3$	2/9
$I_1 I_2 I_5$	2/9
$I_1 I_2 I_4$	1/9 (x)
$I_2 I_3 I_4$	0
$I_2 I_3 I_5$	1/9 (x)

### 3-freq Set

Item	Sup
$I_1 I_2 I_3$	2/9
$I_1 I_2 I_5$	2/9



nd valid association rules.

$$I_1 \rightarrow I_2 I_3$$

$$\text{Support} = \frac{\sigma(I_1 I_2 I_3)}{|\mathcal{H}|} = \frac{2}{9}$$

$$\begin{aligned} \text{Confidence} &= \frac{\sigma(I_1 I_2 I_3)}{\sigma(I_1)} = \frac{2/9}{6/9} \times 100 \\ &= \frac{2}{6} \times 100 \\ &= 33.33\% < 50\% \end{aligned}$$

$$I_2 \rightarrow I_1 I_3$$

$$\begin{aligned} \text{Conf} &= \frac{\sigma(I_1 I_2 I_3)}{\sigma(I_2)} = \frac{2/9}{7/9} \times 100 \\ &= 28.57\% < 50\% \end{aligned}$$

$$I_3 \rightarrow I_1 I_2$$

$$\begin{aligned} \text{Conf} &= \frac{\sigma(I_1 I_2 I_3)}{\sigma(I_3)} = \frac{2/9}{6/9} \times 100 \\ &= 33.33\% < 50\% \end{aligned}$$

$$I_1 I_2 \rightarrow I_3$$

$$\begin{aligned} \text{Conf} &= \frac{\sigma(I_1 I_2 I_3)}{\sigma(I_1 I_2)} = \frac{2/9}{4/9} \times 100 = 50\% \\ &\quad \checkmark \quad (= 50\%) \end{aligned}$$

$$I_1, I_3 \rightarrow I_2$$

(✓)

$$\text{conf} = \frac{\sigma(I_1, I_2, I_3)}{\sigma(I_1, I_3)}$$

$$= \frac{2/9}{4/9} \times 100 = 50\%$$

50%

$$I_2, I_3 \rightarrow I_1$$

(✓)

$$\text{conf} = \frac{\sigma(I_1, I_2, I_3)}{\sigma(I_2, I_3)} = \frac{2/9}{4/9} \times 100$$

$$= 50\% < \underline{\underline{50\%}}$$

$$I_1 \rightarrow I_2, I_5$$

$$\text{Support} = \frac{2/9}{9} \times 100 = 22.22\% > \underline{\underline{20\%}}$$

$$\text{conf} = \frac{\sigma(I_1, I_2, I_5)}{\sigma(I_1)} \times 100$$

$$= \frac{2}{9} \times \frac{9}{6} \times 100 = 33.33\% < 50\%$$

$$I_2 \rightarrow I_1, I_5$$

$$\text{conf} = \frac{\sigma(I_1, I_2, I_5)}{\sigma(I_2)} \times 100$$

$$= \frac{2}{9} \times \frac{9}{7} \times 100 = 28.57\% < 50\%$$



$$I_3 \rightarrow I_1, I_2$$

$$\textcircled{v} \quad \text{conf} = \frac{\sigma(I_1, I_2, I_3)}{\sigma(I_3)} = \frac{2}{9} \times \frac{9}{2} \times 100 = 100\% \geq 50\%$$

$$I_1, I_2 \rightarrow I_3$$

$$\textcircled{v} \quad \text{conf} = \frac{\sigma(I_1, I_2, I_3)}{\sigma(I_1, I_2)} = \frac{2}{9} \times \frac{9}{4} \times 100 = 50\% \geq 50\%$$

$$I_1, I_3 \rightarrow I_2$$

$$\textcircled{v} \quad \text{conf} = \frac{\sigma(I_1, I_2, I_3)}{\sigma(I_1, I_3)} = \frac{2}{9} \times \frac{9}{2} \times 100 = 100\% \geq 50\%$$

Rules  $I_1, I_3 \rightarrow I_2$

$$I_2, I_3 \rightarrow I_1$$

$$I_3 \rightarrow I_1, I_2$$

$$I_1, I_2 \rightarrow I_3$$

$$I_1, I_3 \rightarrow I_2$$

① Apply Apriori Algorithm for Market Basket Analysis.

Soln • Dataset used has 20,000 entries over 9000 transactions

- Coffee & bread are most frequent items

Afternoon is peak hours of sales.

- most sales are on Fri, Sat & Sun.

- most productive month is march

- The support for Diptyques & Bread is highest.

- Bread → Pastry, Cake → Coffee, Coffee → Cake, Cake → Tea are

- Seeing connection graph, coffee has most connections: & followed by cake with 2.

- The rules are refined further to being Tea → Bread & Coffee → Bread as frequent patterns.



Design FP Tree for breast cancer detection.

Soln The FP Tree growth algorithm is used mainly to decrease the number of features in the dataset to feed in to the neural network. After dim. reduction, the features are fed to NN.

Some terms: Confidence  $(x \rightarrow y) = \frac{\text{Support}(x, y)}{\text{support}(x)}$

Database:

#Records	#Malignant	#benign
699	241	458

Summary of db characteristics:

code	Desp	mean	Std
I	Thickness	...	...
II	Uniformity of cell size	...	...
...	...	...	...
IX	mito ses	...	...

③ Apply vertical data partitioning domain.

Now our task is to bring down the number of features to be fed to the ANN.

1st method: Suppose we have association rule  $X \rightarrow Y$  where  $X$  can take up multi itemsets &  $Y$  also. If the set items in  $X$  pass the minSup & minConf levels, we can remove  $Y$ .

eg  $\{I, III, VIII, IX\} \rightarrow \{II\}$

Here  $I, III, VIII, IX$  have values  $\geq$  minSup & minConf, so  $II$  can be removed.

④ Max Miner method

Here we ~~remove~~ choose only supersets of features that ~~are~~ do not have subsets that are frequent. This will help reduce the # of features.



1 class benign (value = 1):

Input: {II, VIII, IX}

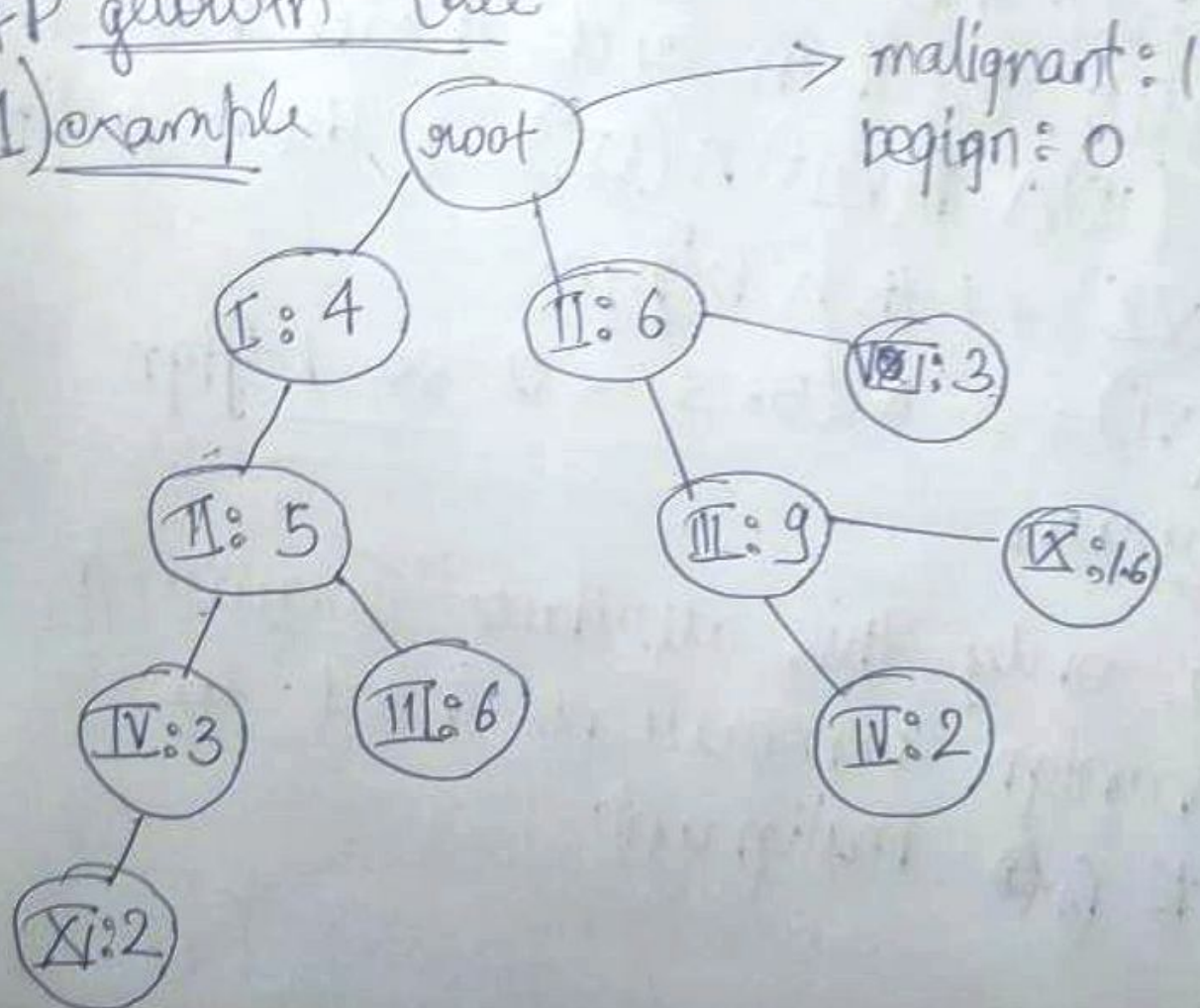
Class malignant (value = 0):

Input: {~~VI~~}

This implies classes II, VIII, IX enough  
for malignant class detection & VII  
enough for benign class detection.  
This can be fed to ANN.

FP growth tree

1) example



To solve above example, first  
 let's find the means of each of  
 the features that are important.

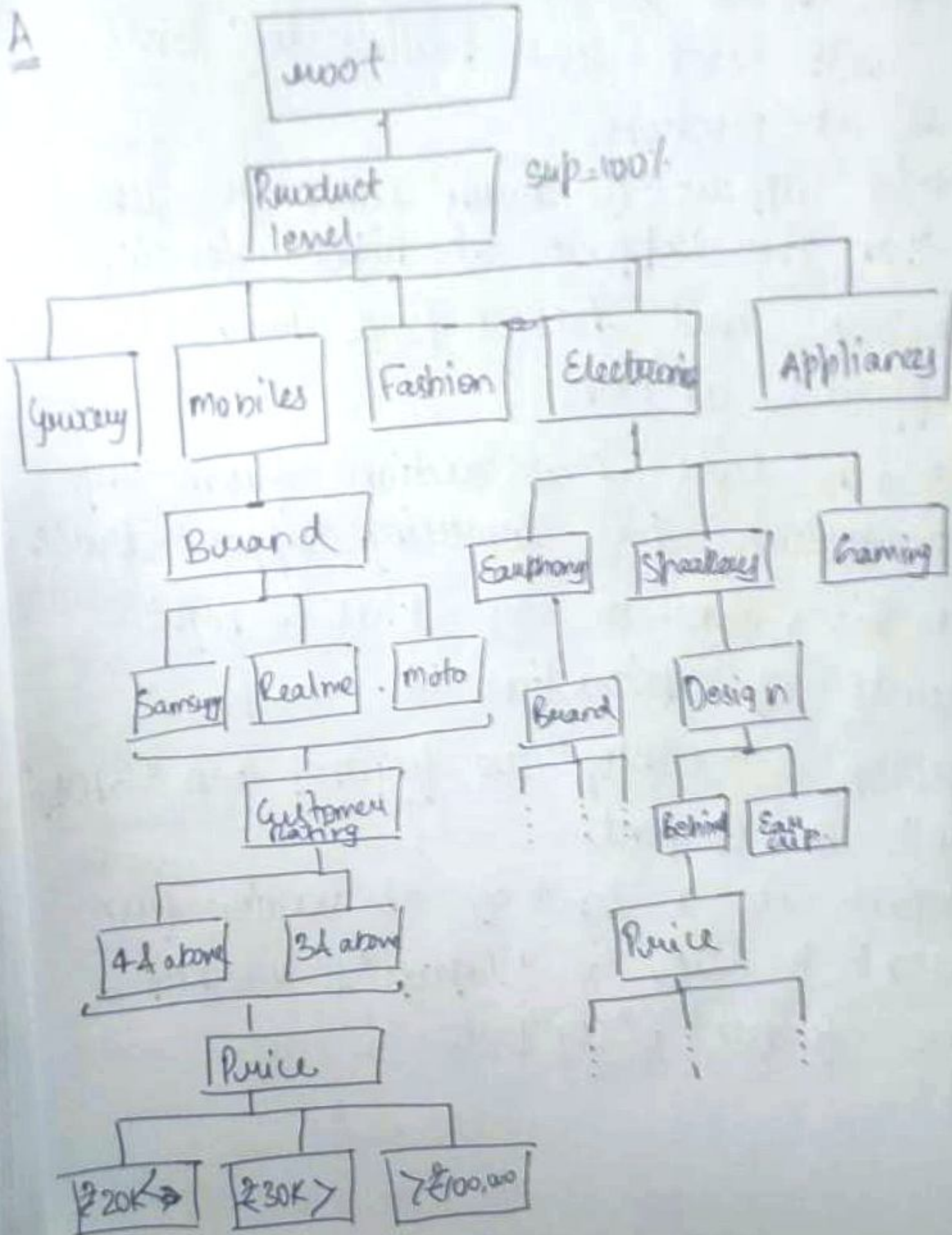
II	3.1	→ Uniformity of Cell
VIII	2.9	→ Normal Nucleoli
IX	1.6	→ Mitoses.
VI	3.5	→ Bore Nuclei

Now these means act as minSup.  
 we back track from leaf to  
 root. If the value  $>$  minSup assign  
 1 else 0. Then and these values.  
 If  $(II) \wedge (VIII) \wedge (IX) = 1$  then malignant  
 &  $VI = 1$  then benign.  
 $\rightarrow (VI) = 3 < 3.5 = 0$  so benign  
tumour.

Since only this attribute required  
 for benign tumour no need to  
 check for malignant.



Designs - multilevel multidimensional  
pattern mining for PipRant.



The previous page shows a pattern mining format for Flipkart.

→ I have built a multilevel model with each level having diff. level of abstraction.

→ The support @ lower levels is less than the support at higher levels.

→ Each level has a fixed threshold support & confidence.

→ Each level of abstraction from top to bottom has lowering support levels.

→ This means to say that @ higher levels of abstraction if "computer", "laptop" & "desktop" are frequent then "tabletop" will be frequent.

→ But at a level of abstraction from parent to leaf, if "laptop" & "desktop" are frequent, "computer" is not.



## Module 4

1. Design DT for play tennis example.

Outdoor	Temperature	Humidity	Windy	Play
Sunny	Hot	High	False	N
Sunny	Hot	High	True	N
Overcast	Hot	High	False	Y
Rainy	Mild	High	False	Y
Rainy	Cool	Normal	False	Y
Rainy	Cool	Normal	True	N
Overcast	Cool	Normal	True	Y
Sunny	Mild	High	False	N
Sunny	Cool	Normal	False	Y
Rainy	Mild	Normal	False	Y
<del>Overcast</del>	Mild	Normal	True	Y
Sunny	Mild	High	True	Y
Overcast	Mild	Normal	False	Y
Overcast	Hot	High	True	N
Rainy	Mild	High	True	N

$$\text{Entropy (D)} = -\frac{9}{14} \log \frac{9}{14} - \frac{5}{14} \log \frac{5}{14}$$

$$= 0.4097 + 0.530$$

$$= 0.940$$

E(Temp)

$$E_{\text{hot}} = (2+, 2-) = 1$$

$$\begin{aligned} E_{\text{mid}} &= (4+, 2-) = -\frac{4}{6} \log \frac{4}{6} - \frac{2}{6} \log \frac{2}{6} \\ &= 0.389 + 0.528 \\ &= 0.917. \end{aligned}$$

$$\begin{aligned} E_{\text{cool}} &= (3+, 1-) = -\frac{3}{4} \log \frac{3}{4} - \frac{1}{4} \log \frac{1}{4} \\ &= 0.311 + 0.5 \\ &= 0.811 \end{aligned}$$

$$E_{\text{rot}} = (2+, 2-) = 1$$

$$\begin{aligned} \text{Gain (temp)} &= 0.940 - \frac{4}{14} \times 1 - \frac{6}{14} \times 0.917 - \\ &\quad \frac{4}{14} \times 0.811 \\ &= \boxed{0.029} \end{aligned}$$

$E(\text{humidity})$ .

$$\begin{aligned} E_{\text{high}}(3+, 4-) &= -\frac{3}{7} \log \frac{3}{7} - \frac{4}{7} \log \frac{4}{7} \\ &= 0.5283 + 0.461 \end{aligned}$$

$$\begin{aligned} E(6+, 1-) &= 0.924 \\ &= -\frac{6}{7} \log \frac{6}{7} - \frac{1}{7} \log \frac{1}{7} \\ &= 0.190 + 0.401 = 0.591 \end{aligned}$$



$$E_{\text{normal}}(6+, 1-)$$

$$= -\frac{6}{7} \log \frac{6}{7} - \frac{1}{7} \log \frac{1}{7}$$

$$= 0.401 + 0.190$$

$$= 0.591$$

$$E_{\text{high}}(3+, 4-)$$

$$= -\frac{3}{7} \log \frac{3}{7} - \frac{4}{7} \log \frac{4}{7}$$

$$= 0.523 + 0.461$$

$$= 0.984$$

Gain humidity  
(~~perp~~)

$$= 0.940 - \frac{7}{14} \log \times 0.591 - \frac{7}{14} \times 0.984$$

$$= 0.940 - 0.2955 - 0.47$$

$$= \boxed{0.1745}$$

e(wind)

$$E_{\text{True}}(3+, 3-) = 1$$

$$E_{\text{False}}(6+, 2-) = -\frac{6}{8} \log \frac{6}{8} - \frac{2}{8} \log \frac{2}{8}$$

$$= 0.311 + 0.5 = 0.811$$

$$\text{Gain (wind)} = 0.940 - \frac{6}{14} \times 1 - \frac{8}{14} \times 0.811$$

$$= 0.940 - 0.428 - 0.463$$

$$= \boxed{0.049}$$

Outlook  
E(Outlook) =

$$E(\text{Rainy} | (2+, 3-))$$

$$= -\frac{2}{5} \log \frac{2}{5} - \frac{3}{5} \log \frac{3}{5}$$

$$= 0.233$$

$$E(\text{Outlook}) = 0.940$$

= Overcast (4+, 0-)  
= 0.

$$E(\text{Rainy}) (3+, 2-)$$

$$= -\frac{3}{5} \log \frac{3}{5} - \frac{2}{5} \log \frac{2}{5}$$

$$= 0.442 + 0.528 = 0.970$$

$$\text{Gain}(\text{Outlook}) = 0.940 - \frac{5}{14} \times 0.233 -$$

$$- \frac{5}{14} \times 0.970$$

$$= \boxed{0.510}$$

$$\text{Root Node} = \begin{array}{c} 0.029 \\ T \end{array} \bigg| \begin{array}{c} 0.049 \\ W \end{array} \bigg| \begin{array}{c} 0.510 \\ O \end{array} \bigg| \begin{array}{c} 0.1745 \\ H. \end{array}$$

Outlook is root node





age	income	student	credit	com
		N	fair	N
<=30	high	N	exc	N
<=30	high	N	fair	Y
31..40	high	N	fair	Y
>40	medium	Y	fair	Y
>40	low	Y	excl	N
>40	low	Y	excl	yes.
31..40	low	N	fair	N
<=30	medium	Y	fair	Y
<=30	low	Y	fair	Y
>40	medium	Y	excl	Y
<=30	medium	N	excl	Y
31..40	medium	Y	fair	Y
31..40	high	N	excl	N
>40	medium	$P(\frac{N}{Y}) =$		

age	income	student	credit
<=30	high	2/9	2/5
31..40	medium	4/9	2/5
>40	low	3/9	1/5
		Y	N
Student		6/9	2/5
Y		4/5	3/5
N		3/9	



Apply SVM for breast cancer domain.

Preprocessing: ① Check for null values.

EDA: Check for value counts of benign & malignant tumours here. Box Plots & Violin Plots allow us to see ~~which~~ what distributions the data follows.

Classification Models: Split the data into train & test sets.

KNN gave an accuracy of 92%. If there are normalization problems, we can scale the data.

We can further use SVM to classify the type of cancer & tweak with multiple kernels. SVM gave accuracy of 98%.

Logistic Regression gave an accuracy of 0.92.

So SVM can be the best model out of all of them.

④ Apply outliers detection in t-doom Indu

① Tukey's IQR method.

① Find  $Q_1$  ② Find  $Q_3$  ③  $IQR = Q_3 - Q_1$

④ Define the normal data range with lower limit  $Q_1 - 1.5 IQR$  & upper limit as  $Q_3 + 1.5 IQR$ .

Any point beyond this range is outlier.

② Standard deviation

68% of data lies within 1 std deviation of mean.

95% of data lies in 2 std dev of mean.

99.7% of data lies in 3 std dev of mean.

As the std dev increases the outliers increases.

Z score method: We rescale & center the data & look for pts which are too far from zero.

modified Z score value is robust to outliers



## Module 5

Analyse the sentiments & polarise into  
+ve & -ve sentiments

① I have used nltk to classify tweets.

② Check the value - counts of +ve, -ve & 0.

③ Remove null value tweets

④ Remove URLs from tweets

⑤ Tokenize text, Remove emails, newline  
characters & punctuation signs.

⑥ Lowercase all text.

## Model Building

~~And~~ Some models I built were single LSTM  
with 0.74 accuracy.

Bidirectional LSTM with 0.709 accuracy

One D CNN with 0.59 accuracy

Tune parameters to increase accuracy.

② Analyse sarcasm using clustering approach on Twitter dataset.

A ① Drop article link attribute

② Check value counts of sarcastic & non sarcastic tweets

③ Plot word cloud of tweets.

④ Split data into train & test set.

⑤ Use tfidf to vectorize the text.

⑥ Use logistic regression  $\rightarrow$  0.76 f1 score.

Use Naive Bayes  $\rightarrow$  0.784 f1 score

⑦ Voting classifiers  $\rightarrow$  0.73 f1 score.

⑧ Linear SVC gave the best results

⑨ RNN gave accuracy of 77%.