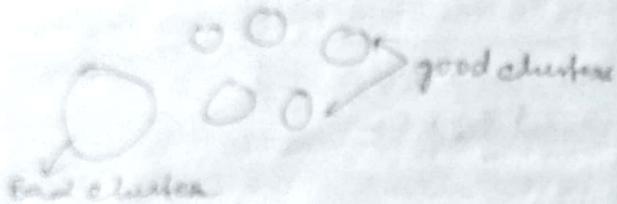
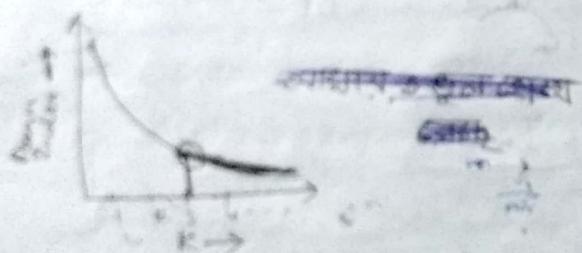


→ If $DN(c)$ is very high, then cluster is very good. Very low means N is high which mean inter-separation is more and D^* is less hence compactness is more.
 → Our obj is to maximize $DN(c)$



Diadv :-

① Instead of max in $\Delta(Dr, D^*)$, we can use avg value also when all clusters are good, if one is bad then if we take max value then it will give bad result.



Cluster Validation Indices

② Davies - Bouldin (DB) Index :-

$$\frac{\Delta(c_i) + \Delta(c_j)}{d(c_i, c_j)} \rightarrow \frac{\Delta(c_1) + \Delta(c_2)}{d(c_1, c_2)}$$

③ $\text{O} \text{ O} \text{ O} \dots \text{ O}$

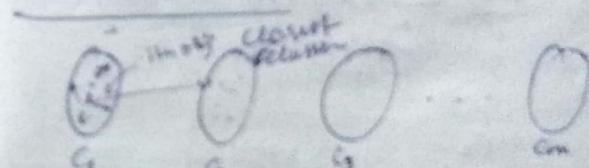
→ In the worst case,

$$\sum_{i=1}^m \max_{j=1 \dots m, i \neq j} \left\{ \frac{\Delta(c_i) + \Delta(c_j)}{d(c_i, c_j)} \right\} = DB(c)$$

- For better clusters, DB index must be low.
- Need to compute the K value for $K = 2, 3, 4, \dots$ to compute DB index.

- The lowest value of K that gives better DB index, is the ~~best~~ no. of clusters.
- If we don't know the actual requirement then we have to compute all the index values.
- If intercluster & ~~inter~~ intracuster are the criteria then we can use DB index.

Silhouette Index (SI)



we can calculate:

$SI(i) \rightarrow$ example

$SI(C_i) \rightarrow$ for cluster i

$SI(c) \rightarrow$ for all clusters

$a(i)$ = avg. dissimilarity of i th obj with the cluster in which it lies

→ then we have to find the closest cluster from i th object.

This can be done by:

- taking the centroid of other clusters
- or → calculating dist. with all objects and take their avg.

$b(i)$ = avg. dissimilarity of i th obj with ~~all~~ ~~closest~~ ~~closed~~ cluster. all objects in its closest cluster.

$$SL(i) = \frac{b(i) - a(i)}{\max\{b(i), a(i)\}}$$

$$-1 \leq SL(i) \leq 1$$

→ $a(i)$ must be low and $b(i)$ must be high.

if $a(i) > b(i) \Rightarrow$ compactness is less.
i.e.



if no. of objs in C_2 is very less than C_1
then $a(i) > b(i)$
↳ more values.

→ For good cluster, value must be nearer to 1.

If $SL(i) \approx -1$,

$$\Rightarrow b(i) - a(i) \approx \max\{b(i), a(i)\} \\ = -a(i)$$

$$b(i) \approx 0$$

That means i th obj is not properly placed.

If $SL(i) \approx 1$

$$b(i) - a(i) \approx b(i)$$

$$\Rightarrow a(i) \approx 0 \quad (\text{i.e. obj is in actual cluster})$$

If $SL(i) \approx 0$

$$b(i) = a(i)$$

→ Obj may be in any of the clusters.

$$SL(c_k) = \frac{1}{|C_k|} \sum_{i \in c_k} SL(i)$$

$$SL(C) = \frac{1}{m} \sum_{k=1}^m SL(c_k)$$

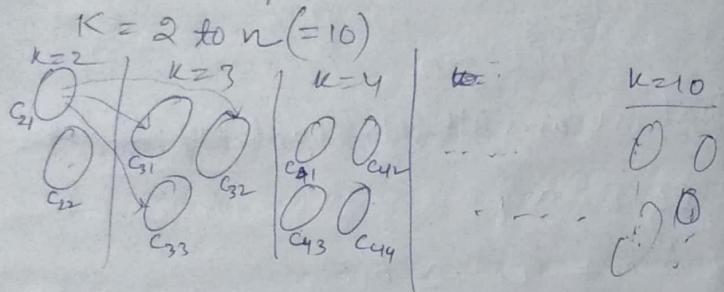
→ If $\underset{\text{no. of clusters}}{SL(C)}$ is higher, the cluster is better cluster

Requirement

→ Dissimilarity of obj. with centroid of other cluster or values of ~~all~~ objects of other cluster

Family and Lieu

→ Checks whether an object is stable within a cluster or not. If yes → keep, no → discard



$$\max \left\{ \frac{c_{21} \cap c_{31}}{|c_{21}|}, \frac{c_{21} \cap c_{32}}{|c_{21}|}, \frac{c_{21} \cap c_{33}}{|c_{21}|} \right\} = m_1$$

Similarly with $K=4, \dots, 10$

$$m_2, m_3, \dots, m_8$$

Then take min of all
i.e.

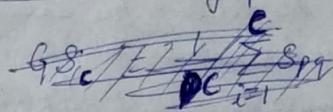
$$S(c_{21}) = \min \{ m_1, m_2, \dots, m_8 \}$$

Stability of c_{21}

$$S_{21} = \min_{i=3 \text{ to } n} \max_{j=1 \text{ to } i} \left\{ \frac{|c_{21} \cap c_{ij}|}{|c_{21}|} \right\}$$

$$S_{cl} = \min_{i=K+1 \text{ to } n} \left\{ \max_{j=1 \text{ to } i} \left\{ \frac{|c_{cl} \cap c_{ij}|}{|c_{cl}|} \right\} \right\}$$

General stability for a particular K

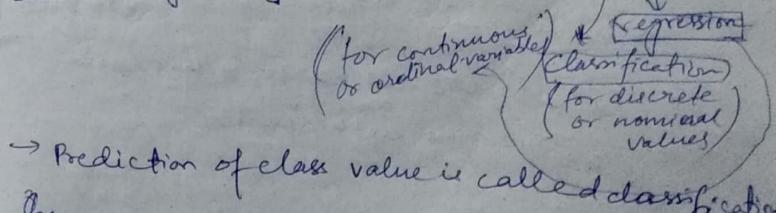


$$GSC = \frac{1}{C} \sum_{i=1}^C S_{ci}$$

Classification

- Model ~~testing~~ training (Model construction)
- Testing (Prediction and Validation)

Issues regarding classification and prediction



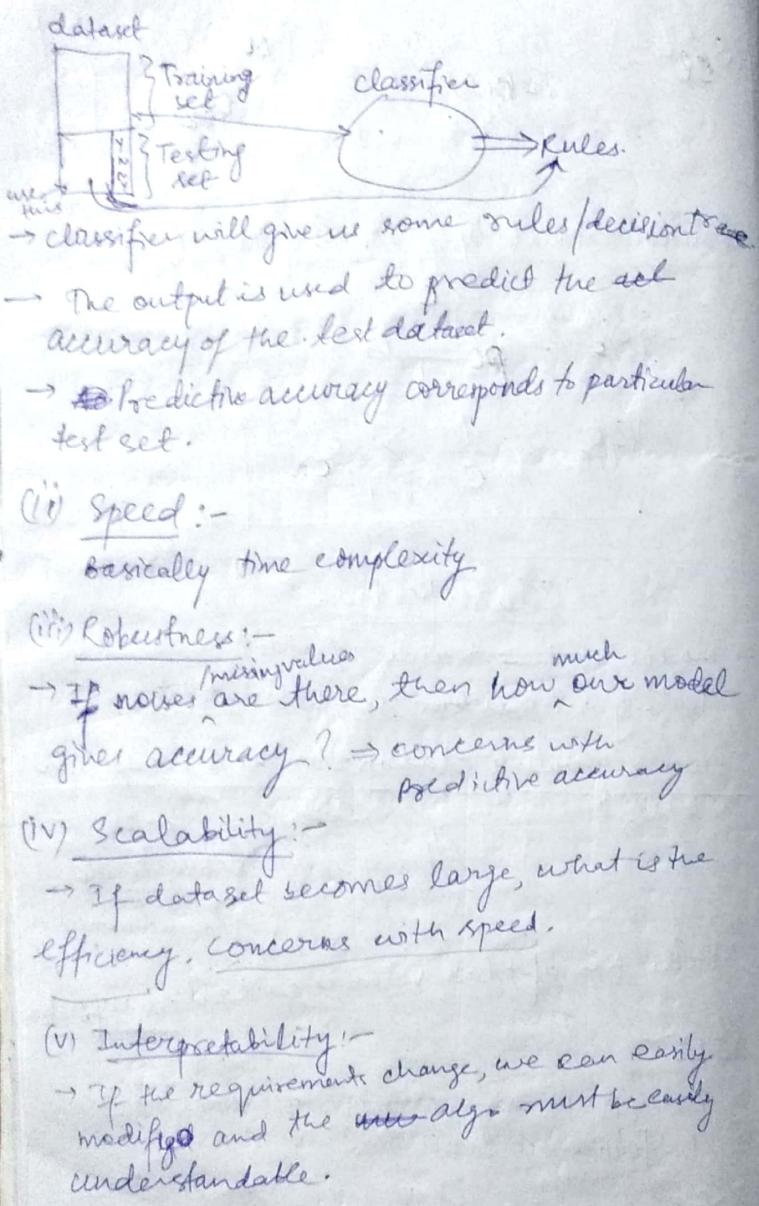
→ Prediction of class value is called classification

Steps :-

- ① Data preparation (All preprocessing tasks)
- ② Performance evaluation

Performance Evaluation :-

- Predictive Accuracy



Accuracy :-

(i) Holdout Method :

Training :- $\frac{2}{3}$ of the dataset (randomly select)

Test set :- $\frac{1}{3}$ of the dataset → check the accuracy.

→ If accuracy is high, we can rely on the training.

Drawbacks

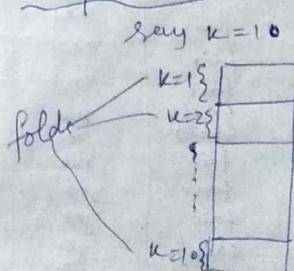
If there are not sufficient objects available in a class, then training will not be correct and hence misclassification will occur

$$\text{class} \rightarrow C_1 \Rightarrow \{ \text{sufficient data} \}$$

$$C_2 \Rightarrow \{ o_1 \\ o_2 \\ o_3 \}$$

$$C_3 \Rightarrow \{ \text{sufficient} \}$$

(ii) K-fold cross validation:



→ Take avg. of all k accuracies, and that is taken as the accuracy of the classifier.

(iii) Bootstrapping:-

- Used when dataset is very small.
- Say test data contains only 1 or 2 objects.
- Repeat this process (for no. of obj)
- Take the first one, in next iteration the other one, then finally avg

Overfitting of Classifier:-

→ Training gives ~~test~~ good ~~good~~ accuracy, but test data will not give ~~good~~ ^{avg} accurately.

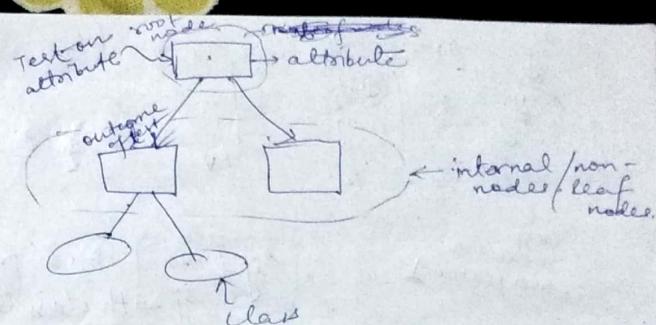
This is because:

- There are some noisy data in training set.
- Insufficient no. of features
- Irrelevant features

27/3/18

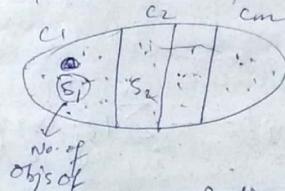
Decision Tree based classifier (Rule Based)

	age	income	student	credit rating	buys comp. (class)
1	≤ 30	high	No	fair	no
2	≤ 30	high	No	Excellent	no
3	31...40	high	No	fair	yes
4	> 40	medium	No	fair	yes
5	> 40	low	yes	fair	yes
6	> 40	low	yes	excell.	no
7	31...40	low	yes	excell.	yes
8	≤ 30	medium	No	fair	no
9	≤ 30	low	yes	fair	yes
10	> 40	med	yes	fair	yes
11	≤ 30	med	yes	excel.	yes
12	31...40	med	No	excel.	yes
13	31...40	high	yes	fair	yes
14	> 40	med	No	excel.	no



→ Information gain for each attribute is calculated

$C_1, C_2, \dots, C_m \Rightarrow m\text{-class}$

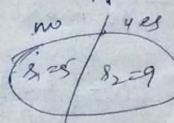


$$\text{No. of objects of class } i \text{ in } C_1 = \frac{n_i}{S} = p_i$$

$$I(S_1, S_2, \dots, S_m) = -\sum p_i \log p_i$$

expected info to classify a sample

e.g.



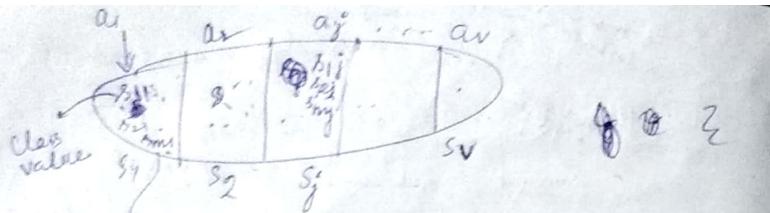
$$I(S_1, S_2) = -\frac{5}{14} \log \frac{5}{14} - \frac{9}{14} \log \frac{9}{14} = 0.940$$

Consider a conditional attr. A

→ Check how many distinct values are there corresp to A.

$$A = \{a_1, a_2, \dots, a_n\}$$

discretization



$$I(S_1j, S_2j, \dots, S_mj) = ?$$

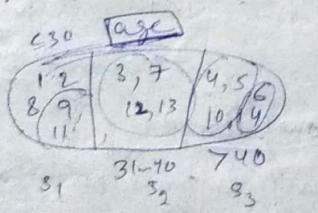
$$P_{ij} = \frac{S_{ij}}{|S_j|}$$

$$I(S_1j, S_2j, \dots, S_mj) = - \sum_{i=1}^m P_{ij} \log P_{ij}$$

$$E(A) = \sum_{j=1}^v |S_j| I(S_1j, S_2j, \dots, S_mj)$$

$$\text{Information gain}(A) = I(S_1, S_2, \dots, S_m) - E(A)$$

→ Corresponding to each attr. calculate $I(A)$, and then take the max attr. with max $I(A)$ and stop.



$$E(\text{age}) = \frac{5}{14} \times 0.971 + 0 + \frac{5}{14} \times 0.971 \\ = \frac{5}{7} \times 0.971 = 0.694$$

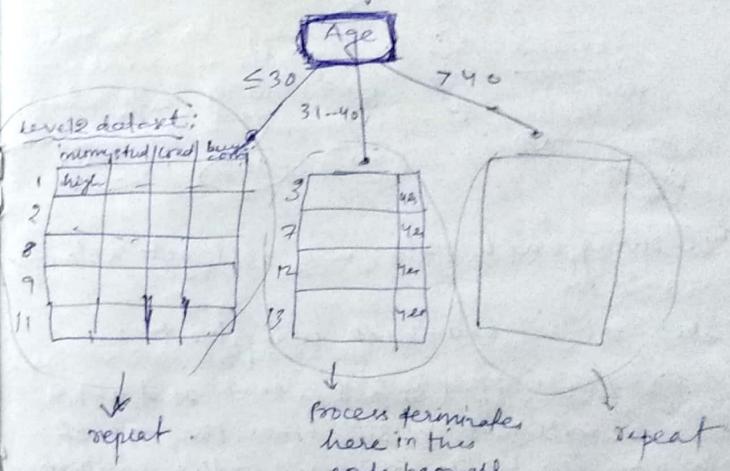
$$\text{gain}(\text{Age}) = 0.940 - 0.694 \\ = 0.246$$

Compute other attributes:

$$\text{gain}(\text{income}) = 0.029$$

$$\text{gain}(\text{student}) = 0.151$$

$$\text{gain}(\text{credit_ratio}) = 0.048$$



→ It is greedy approach bcoz when max inf. gain is there we are following top down approach

Termination:

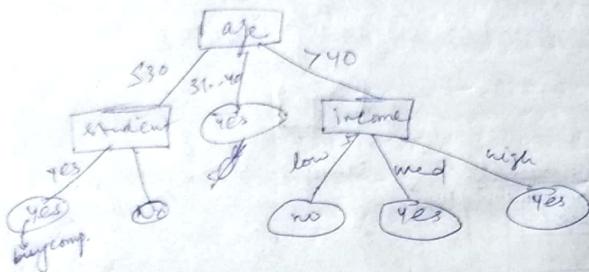
→ we have splitted based on all attrs, we don't have any other attr. to split. (only one attr. condition attr. with class label)

→ e.g.

credit class	
1	
2	
3	
4	
5	
6	
7	
8	
9	
10	
11	
12	
13	
14	
15	
16	
17	
18	
19	
20	
21	
22	
23	
24	
25	
26	
27	
28	
29	
30	
31	
32	
33	
34	
35	
36	
37	
38	
39	
40	
41	
42	
43	
44	
45	
46	
47	
48	
49	
50	
51	
52	
53	
54	
55	
56	
57	
58	
59	
60	
61	
62	
63	
64	
65	
66	
67	
68	
69	
70	
71	
72	
73	
74	
75	
76	
77	
78	
79	
80	
81	
82	
83	
84	
85	
86	
87	
88	
89	
90	
91	
92	
93	
94	
95	
96	
97	
98	
99	
100	

corresponds to 3 values
1 2 3
attr 1
attr 2
leaf node

- No more samples are there, so we can terminate.
- If all objects are in same class, say decision tree is:



→ Following each path, we will get a classification rule.

→ Take test set and check using the tree.

→ Using k folds, if $1 \leq k$ fold is test and $k-1$ folds are training, we will get K diff. trees, the avg. of accuracy of all K trees is the accuracy of the classifier.

Decision Tree Classifier :-

28/3/18

↓
Overfitting problem

comes to any tree,
training accuracy is 98% (say)
and another 70%. Testing accuracy
another } overfitted.

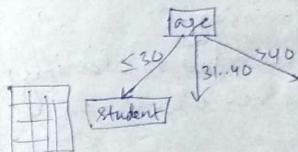
$T_1 \rightarrow 90\%$ } not overfitted.
 $T_2 \rightarrow 85\%$ }

- leads to outliers.
- tries to irrelevant noisy data.
- solving Overfitting :-
- Remove some branches from tree.
- ↳ faster classification
- ↳ More accuracy
- This is called Pruning.

PRUNING

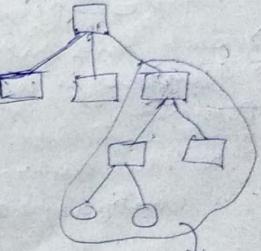
PRE PRUNING

before splitting check whether splitting is reqd. or not (done by checking the strength of splitting).

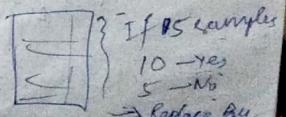


income, student, credit rating,
buys-computer
calculate IG and compare
with threshold. If less, stop.
Eg. If for student IG is high
but < than threshold.
Then Stop, don't classify

Decision Rule



Take away the whole subtree and replace by a leaf node.
leaf node construction is based on the table for the branch.



If 15 samples
10 - yes
5 - no
⇒ Replace by (yes)

If error is more, then don't prune this node.

- Go to next node.
- Now we have several pruned trees, from all these pruned classifiers take the best classifier and that can be done by taking a test set which is independent of train set. The one which gives more accuracy is the best pruned decision tree.

Merits

- Learning technique is easy.
 - Simple: Divide and conquer.
 - Each time only 1G is needed to classify.
 - Accuracy is comparable to other classifiers.
- Demerit :-
- Without discretisation, can't be classified for continuous attributes, intervals, dataset.

Naive Bayes Classifier :-

- Statistical classifier where prob. of different samples, classes, etc. are to be calculated.

X : Object

H : Hypothesis

- Say. H says: X is in class C and X is the object whose class level is not given.

say attributes be - colour and shape
and class : fruit type

colour	shape	class
Red	round	apple

$$X : (\text{Red}, \text{Round})$$

$$H : X \text{ is in class } C$$

$P(X)$: Prob that X is either red and round.

$P(H)$: Prob that X is in class apple.

~~$P(H/X)$~~

$P(X/H)$: What is the prob. that X is red & round given X belongs to class apple

$$P(H/X) = \frac{P(X/H) P(H)}{P(X)}$$

$$X = \{x_1, x_2, x_3, x_4\}$$

$H = \{C_1, C_2\}$ different states attributes.
(student, age, name, ...)

Training sets:

	A_1	A_2	A_3	A_4	A_n	Class
X	$x_1 \dots x_2$	$x_2 \dots x_n$	x_n			C_i
	x_1	$x_2 \dots x_k \dots x_n$	x_n			C_j
	$x_1 \dots x_m$	$x_m \dots x_n$	x_n			

$$X = \{x_1, x_2, \dots, x_n\}$$

$$C = \{C_1, C_2, \dots, C_m\}$$

Q. The object X that is to be taken should be in set X

$$P(c_i/x) = \frac{P(x/c_i) P(c_i)}{P(x)}$$

compute $\forall i = 1 \dots m$.

If $P(c_i/x) > P(c_j/x) \forall j = 1 \dots m$, then
 $\Rightarrow X$ is in class C_i

If there are 's' samples, and 's' corresponds to x , values x occurs.

i.e. $x = \{x_1, x_2, \dots, x_s\}$

$$P(x) = \frac{s}{N}$$

$\Rightarrow x$ is a tuple in the dataset.

\rightarrow If probabilities of different classes are not given, then take $P(c_i)$ as equally likely.

$$\text{Now, } P(c_i/x) = K \cdot P(x/c_i) P(c_i)$$

$$K = \frac{1}{P(x)} \quad (\text{For all}),$$

so we can ignore.

$$P(c_i/x) = P(x/c_i) \cdot P(c_i)$$

$$P(c_i) = \frac{s_i}{N}, \quad \begin{aligned} &\text{Btw no. of sample belong} \\ &\text{to class } C_i \\ \hookrightarrow &\text{total no. of samples.} \end{aligned}$$

Assumption :-

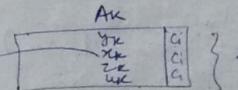
① Attributes are class conditional independent
 \rightarrow i.e. the attributes are not dependent on each other. (Not always true)

$$P(x/c_i) = \prod_{k=1}^n P(x_k/c_i)$$

\rightarrow This is possible when attributes are discrete values.

Cases :-

Case I :- If attribute A_k values are distinct discrete.

A_k  take the values which belong to class c_i

$s_{ik} = \text{no. of samples in class } C_i \text{ with } A_k \text{ value } x_k$

$$\therefore P(x_k/c_i) = \frac{s_{ik}}{s_i}$$

\hookrightarrow No. of samples in class C_i

Case II :- If A_k values are continuous.

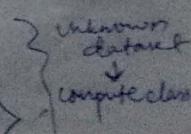
\rightarrow we assume the values follow normal or Gaussian distribution.

$$P(x_k/c_i) = g(x_k, \mu_{ci}, \sigma_{ci}) = \frac{1}{\sqrt{2\pi} \sigma_{ci}} e^{-\frac{(x_k - \mu_{ci})^2}{2\sigma_{ci}^2}}$$

\rightarrow Real values/ordinal values.

\Rightarrow By this compute $P(c_i/x)$ for all classes and take the max value.

Example :-

$X = \{ \text{age} = " \leq 30 ", \text{income} = \text{medium}, \text{student} = \text{yes}, \text{credit_rvew} = \text{fair} \}$ 

$C_1 = \text{Buy computer} \quad \text{Yes}$
 $C_2 = \text{Buy computer} \quad \text{No.}$

$$P(C_1/x) = ? , P(C_2/x) = ?$$

$$P(C_1/x) = P(x/C_1) P(C_1)$$

$$P(C_2/x) = P(x/C_2) P(C_2)$$

$$P(C_1) = \frac{9}{14} \quad \left\{ \text{From prev dataset} \right.$$

$$P(C_2) = \frac{5}{14} \quad \left\{ \quad \right.$$

$$P(x/C_1) = P(x_1/C_1) P(x_2/C_1) P(x_3/C_1) P(x_4/C_1)$$

$$\text{Here } x = \{x_1, x_2, x_3, x_4\}$$

$$= \boxed{\quad} \quad \left\{ \begin{array}{l} 9 \text{ samples} \\ \text{with class yes} \\ x_1 = \text{"Age} \leq 30\text{"} \\ \quad \quad \quad \downarrow \\ \quad \quad \quad 2 \text{ samples} \end{array} \right.$$

$$\therefore P(x_1/C_1) = \frac{2}{9}$$

$$P(x_2/C_1) = \frac{4}{9}$$

$$P(x_3/C_1) = \frac{6}{9}$$

$$P(x_4/C_1) = \frac{6}{9}$$

$$\therefore P(x/C_1) = \frac{288}{9^4} = \cancel{\frac{288}{6561}} 0.044$$

$$P(C_2/x) = P(x/C_2) P(C_2)$$

$$= 0.044 \times \frac{1}{14} = 0.003$$

similarly,

$$P(C_2/x) = 0.007$$

$$\therefore P(C_1/x) > P(C_2/x)$$

so, the class for the dataset is Yes

Demerit:

- It's expected to give good result because it has the critical background. But ~~not~~ we have assumed that attributes are independent.

K-nearest Neighbour Classifier :-

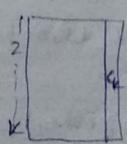
→ density based classifier.



every object is n-dimensional and belongs to space (assumed)

To get the class value of object x

- Take measure euclidean distance of x with all objects and take k neighbours.



→ If neighbours belong to class C_k , then x 's class value is C_k

(max of the neighbours)

→ It is called the lazy learner.

→ Training time less, classification time is high.

→ whenever we are given taking euclidean dist.

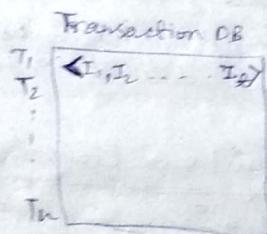
→ we aren't giving importance to different attributes (all attributes are of equal weightage)

- this may effect the classifier
- For each sample calculate the classifier.
- Take $k = \text{even odd}$, because if same value comes in two parties, the vote goes $\frac{1}{2}$ to any one.

4/4/18

Association Rule Mining :-

- Relationship among the items in a database.
- Useful bcz:



Market Basket Analysis:

- Those items which will be in the basket which the customer has the tendency to buy.
- Based on Trans. DB.

$$I_1 \wedge I_3 = I_5$$

- If items I_1 and I_3 are consumed by a customer, then I_5 will also be consumed by him.

Measures of Association rules:-

- ① $A \Rightarrow B$
- ② Support($A \Rightarrow B$)
- ③ Confidence($A \Rightarrow B$)

$$\text{① } \text{support}(A \Rightarrow B) = P(A \cup B) = \frac{\# \text{ of transactions containing both A \& B}}{\text{Total \# of trans}}$$

by

$\langle I_1, I_2, I_3 \rangle$
 $\langle I_1, I_3, I_5 \rangle$
 $\langle I_2, I_5 \rangle$
 $\langle I_2, I_3 \rangle$

$$\text{support}(I_1 \Rightarrow I_3) = \frac{2}{4} = \frac{1}{2}$$

- If $A \cup B \neq \emptyset$
- If no of transactions containing Both A and B is very less than no. of transaction, then support($A \Rightarrow B$) will be very small value which means it is noise.

∴ some threshold must be set i.e.
min support threshold.

$$\text{② } \text{Conf}(A \Rightarrow B) = P(B/A) = \frac{\# \text{ of trans. containing A \& B only}}{\# \text{ of trans. containing A only}}$$

- frequent item sets are within same transaction

Apriori Alg

Apriori Algorithm

- Step 1 : Join step
- Step 2 : Pruning step

	Transac. DB				
T ₁₀₀	I ₁	I ₂	I ₃	I ₄	I ₅
T ₂₀₀		I ₂		I ₄	
T ₃₀₀		I ₂		I ₃	
T ₄₀₀		I ₁	I ₂	I ₄	
T ₅₀₀	I ₁	I ₃			
T ₆₀₀		I ₂	I ₅		
T ₇₀₀		I ₁	I ₃		
T ₈₀₀	I ₁	I ₂	I ₃	I ₄	I ₅
T ₉₀₀	I ₁	I ₂	I ₃		

Find no. of candidate items : I₁, I₂, I₃, I₄, I₅
By scanning the DB.

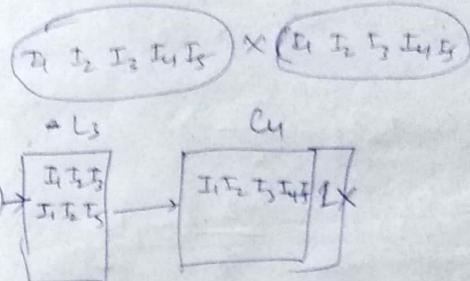
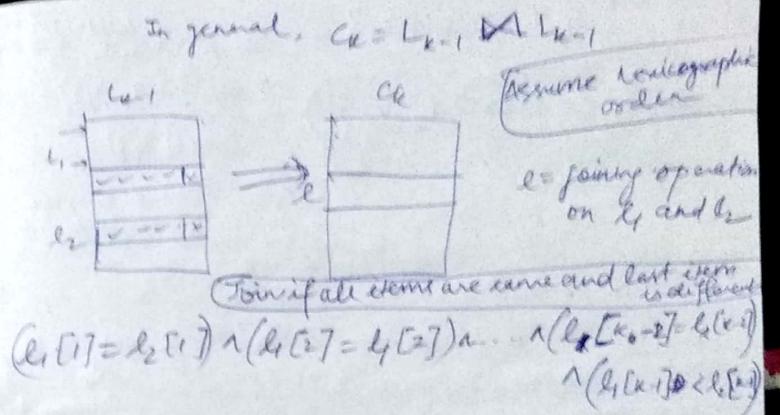
Scan the database
D to count
each candidate

min-support count
will be given

$$\text{min-sup} = \frac{2}{9} \times 100\% = 22\%$$

C ₁	Item	support count	L ₁
	I ₁	6	
	I ₂	7	
	I ₃	6	
	I ₄	2	
	I ₅	2	

frequently occurring element
remove elements which are less than min support count



→ so L₃ is considered as the 3-frequent items sets (2 tuples)

→ Joining step is used to find the candidate items and pruning step is used to remove the ~~unnecessary~~ required tuples.

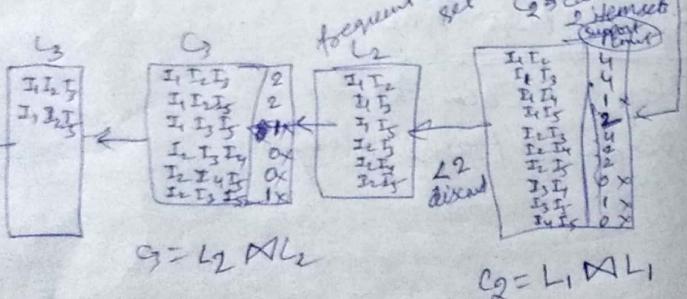
Now

Take < I₁, I₂, I₃ >
subset of this

{I₁} {I₂} {I₃} , {I₁, I₂} {I₁, I₃} {I₂, I₃}

Generate the rules:

$$\begin{array}{l|l} \{I_1\} \Rightarrow \{I_2, I_3\} & I_3 \Rightarrow \{I_1, I_2\} \\ \{I_2\} \Rightarrow \{I_1, I_3\} & \{I_1, I_2\} \Rightarrow I_3 \\ \{I_3\} \Rightarrow \{I_1, I_2\} & \{I_1, I_3\} \Rightarrow I_2 \\ \{I_1, I_2\} \Rightarrow \{I_3\} & \{I_1, I_2\} \neq \{I_2, I_3\} \end{array}$$



$$\{I_1, I_2\} \Rightarrow \{I_1\}$$

No. of trans. where
\$I_1, I_2, I_3\$ occur

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

No. of trans.
where \$I_1\$ occurs

$$= \frac{2}{6} = \frac{1}{3}$$

$$= 33\% \approx \text{Not interesting}$$

↓
Discard

$$\text{conf}(\{I_2\} \Rightarrow \{I_1, I_3\}) = \frac{2}{7} < 70\%$$

$$\text{conf}(\{I_3\} \Rightarrow \{I_1, I_2\}) = \frac{2}{6} \times$$

$$\text{conf}(\{I_1, I_2\} \Rightarrow \{I_3\}) = \frac{2}{4} = \frac{1}{2}$$

$$\text{conf}(\{I_1, I_3\} \Rightarrow \{I_2\}) = \frac{2}{4} = \frac{1}{2}$$

$$\text{conf}(\{I_2\} \Rightarrow \{I_1\}) = \frac{2}{4} = \frac{1}{2}$$

None is relevant

Go for. $\langle I_1, I_2, I_3 \rangle$

$$\text{conf}(\{I_1\} \Rightarrow \{I_1, I_2, I_3\}) = \frac{2}{6}$$

$$\text{conf}(\{I_2\} \Rightarrow \{I_1, I_2, I_3\}) = \frac{2}{7}$$

$$\text{conf}(\{I_3\} \Rightarrow \{I_1, I_2, I_3\}) = \frac{2}{6} = 100\% > 70\%$$

$$\text{conf}(\{I_1, I_2\} \Rightarrow \{I_2\}) = 100\% \quad \checkmark$$

$$\text{conf}(\{I_1, I_3\} \Rightarrow \{I_1\}) = 100\% \quad \checkmark$$

Interesting Measures of a rule :-

① Simplicity :- the structure defining the rule must be as much simple as possible.
e.g. in decision tree if more levels are there then rule length is more.

② Certainty :- confidence value must be high

③ Utility :- usefulness of the rule. It is denoted by support value. Higher support value \rightarrow more utility.

④ Novelty :- If greater the generated rule, but we aren't getting any info

⑤ Novelty :- if any rule is giving some info but is redundant. (It must be eligible to be eliminated). But we keep it.

e.g. $\text{loc}(x, "WB") \Rightarrow \text{buys}(x, "SONY-TV") : [81, 70]$

$\text{loc}(x, "KOL") \Rightarrow \text{buys}(x, "SONY-TV") : [21, 72]$

Redundant.

\Rightarrow used in finding concept hierarchy rule

Drawbacks of Apriori algo :-

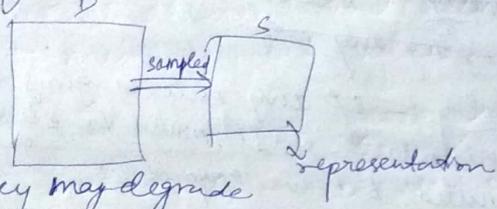
→ If DB size is big, everytime scanning the DB is time consuming.

LK-1

$C_k = \{I_1, I_2, \dots, I_k\}$

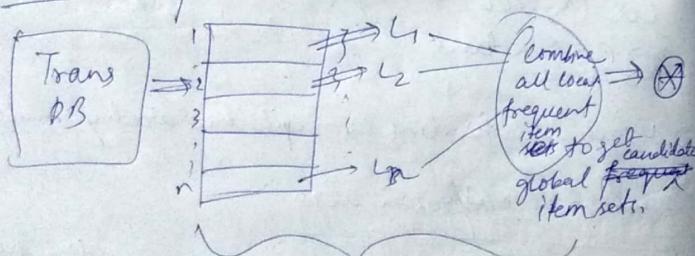
If all the subsets of this are frequent, then it is frequent otherwise not frequent. So we don't need to check scan the database. So we reduce the database.

② Sampling :-



→ Accuracy may degrade
→ Obj - increase accuracy

③ Partitioning :-



Phase I

↓
we are scanning
database once.

→ Partitioning reduces the dataset so that it can perfectly fit in the MM.

the dataset size must be

Phase II

Find Global frequent item sets

Frequent items in DB

Check the support value

④ Hash Table Technique :-

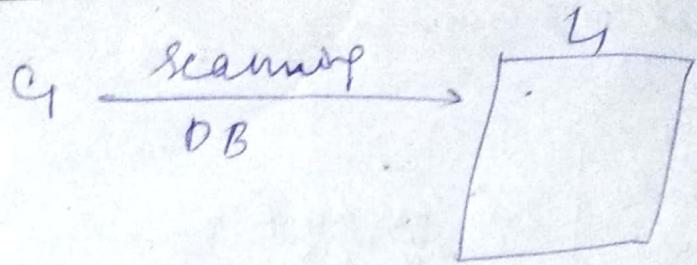
⑤ Transaction Reduction :-

→ If a transaction is not growing frequent item sets then it won't give frequent item set for that transaction. Mark this transaction.

⑥ Hash Table Technique :-

First tuple
 $\{I_1, I_2, I_3\}$

Bucket address	0	1	2	3	4	5	6
Bucket count	2	1	1	1	1	2	
Bucket contents	$\{I_1, I_2\}$	$\{I_3\}$	$\{I_2, I_3\}$	$\{I_1, I_3\}$	$\{I_1, I_2\}$	$\{I_1\}$	$\{I_1, I_2\}$



Generate two item sets corr. to each tyle
and generate hash funcn for these items.

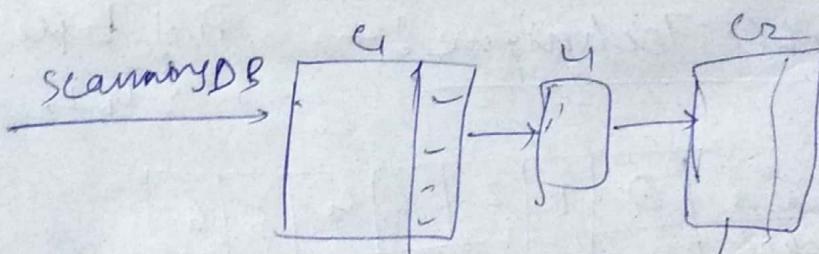
e.g. $\{I_1, I_2\}$

$$h(I_1, I_2) = (1 \times 10 + 2) \bmod 7$$

Consider each element and check the
element min^m support value

↓
more ✓ (take)
less ✗ (forget)

→ In the next step (C_2) → increase the
min^m support value by 1.



No need
to scan the DB,
from the hash table
we will get the
 L_2 and C_2 sets.

