

12-8-20

Date:

Page No.

→ Data Systems Evolution

- Traditional Database Systems
- Indexing - to make the database systems fast
- Query languages - to make convenient to query the database
- Query Optimization - how it is optimally executed
- Transaction → happen according to ACID semantics.
 - A → Atomicity
 - C → Consistency
 - I → Isolation
 - D → Durability

} also allows multiple transactions to occur at the same time
- Recovery - If database is corrupted/power goes off, how will the database recover.

Relational Database Model → most famous

[Query Language → SQL
structured in nature.]

Graph Databases → Entities - Nodes
Relationships - Edges.

Functionalities:

i) Distributed & Heterogeneous

↓
distributed across various system, over cloud etc.
(not on single computer system)

→ different systems will store data in different way for eg: Excel sheets, JSON etc.

And such systems should have interoperability (able to communicate etc.)

Data Analytics

Data Warehousing

- i) Data is stored in one central place called as data warehouse.
- ii) Data needs to be collected from different servers in many different platforms & put it into data warehouse.
- iii) Data warehouse is not updated on a daily basis. Maybe weekly or monthly.
- iv) Purpose : Store all the data in your enterprise across many years to extract patterns.
- v) Data analysis is mainly done manually. using for example: OLAP.

Data Mining

- i) Data can be in the form of csv, data warehouse, etc.
- ii) Data is dumped into the data mining algorithms i.e. automated pattern discovery happens basically.

Data Science Life Cycle → KDD Life Cycle
 (Knowledge Discovery in databases.)

⇒ Types of Patterns

(1) Associations Eg: B/w ~~Sugar &~~ Coffee & Sugar

(2) Clustering Eg: Within a cluster, data or records are very similar to each other. & across the clusters they are very different.
here we do not give any label.

(3) Classification Eg: Predict / Estimate a class or label.

Association Rules

Is defined when we see something happening very frequently.

Eg: Tomato, Potato → Onion

If a person buys T & P, he is very likely to buy O.

$$\text{Confidence} = \frac{2}{3} \rightarrow \text{from below, how many buy O from those who buy T \& P}$$

$$= 66\%$$

$$\text{Support} = \frac{2}{4} = 50\%$$

Probability of RHS given LHS

i.e. $P(P|L)$

Probability of the entire item set occurring ~~at all~~ in the whole data set.

User gives min confidence & min support & based upon this, we need

to define such rules which gives more than the specified min^{sup} % of confidence & support.

Association Rule Application

① E-commerce: Recommendation system based upon associations. (eg: if we liked/bought this, we are likely to buy that.)



Types of Association Rules

① Boolean Association rules

→ if an item is present or not

Eg: Potato is ~~present~~ or not, Tomato is present or not (eg: of association rules slide.)

② Hierarchical Rules

→ **Doubt**

③ Quantitative & Categorical Rules

And other rules also. All of these can be represented as Boolean Association Rules with some tweaks.

⇒ Classification Problem \Rightarrow Supervised Learning

⇒ Clustering Problem \Rightarrow Unsupervised Learning

↓
One way is to find Euclidean distance and categorise based on the distance.
Eg Outlook / Temp / Humidity / Windy

We can treat Temp & Humidity as x & y coordinates & find Eucl. distance
And ~~hence~~ clustering will be done.

Clustering is great for exploration of dataset

\Rightarrow Classification

All the algorithms have a different kind of output. Thus the metrics to be considered is Accuracy, training time etc. and not only response time or memory usage.

(Accuracy) \rightarrow main metric

Accuracy \rightarrow Probability of success.

$$\text{Mean value} = p. \quad (p)$$

$$\text{Variance} = p(1-p)$$

Classification can be treated as Bernoulli process.

Discrete

classified correctly

~~series of~~ events where they are

classified correctly (p)

or ~~"~~ "incorrectly" ($1-p$)

random variable values.

$$X = 1 \quad (\checkmark)$$

$$X = 0 \quad (X)$$

- If there are n records, there are n independent & identically distributed bernoulli trials X_i where $i = 1, \dots, n$

Then the random variable $X = \sum_{i=1 \dots n} X_i$
 is said to follow a binomial distribution

Discrete

Total number of success
 at 1 is for success.

$$P(X=k) = {}^n C_k p^k (1-p)^{n-k}$$

Probability that number of
 correctly classified records are k .

- Normal distribution \rightarrow Continuous

3rd
lecture \rightarrow



Estimating Accuracy

Hold out

Here ~~a~~ a record may never ~~be used~~ come in training or test | whatever be the case

k-fold cross validation

Here each record is treated as training or test at least once



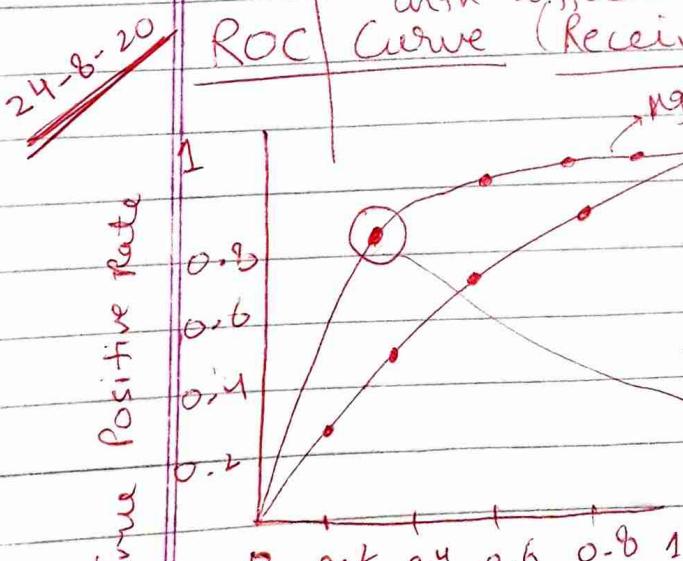
Classifiers are combined for many — reasons:

- i) Different strengths
- ii) Some are good for sparse data & some are good for closely related data.

No Free Lunch Theorem — There is no algorithm which can do the best on all datasets. An algo. may perform good on one dataset & bad on another.

each point basically is a run of ~~an~~ the algorithm with different parameter value.

ROC Curve (Receiver Operating Characteristic)



algo with different parameter settings

Best parameter setting bcz it is closest to 1 for true positive rate (which we wanna achieve)

Thus accuracy can be taken as area under the curve. More the area under curve, more accurate the algo.

→ Ensemble Classifiers - We take a particular classification algorithm and make tweaks in it & generate say k different classifiers (underlying algo is same. Just the different tweaks to same algo.)

e.g.: Now each classifier classifies the given record. And based on the k results, we decide which is the correct class.

Homogeneous classifiers - classifiers emerging from the same kind like decision tree etc. with small tweaks.

Heterogeneous classifiers - mean different types

In Bagging or Boosting, we have k number of classifiers, → they can be homo or hetero.

Majority

Bagging - The class which gets the majority votes by different ensemble classifiers is chosen

→ Boosting - Give ~~majority~~ weight

~~meritocracy~~

Outliers - points which are very ~~far~~ far away from others.

Date:

Page No.

→ Drawback of KNN

- Time Consuming because whole training data is the model itself. And when given a record, we need to search the whole training data to get K nearest neighbours.

~~lecture~~

→

~~Doubt~~

Reverse Nearest Neighbours - used specifically when there are outliers.

- ⇒ • K is typically chosen as an odd number
• K should be small.

→

Decision Tree

- Pure Node - where the class gets identified.
- Impure Node - where there are instances of different classes.

Q How to choose attributes for splitting?

① ID3 Algorithm/method is used.

G Use Info gain for the best split.

~~Info Gain = Entropy of the dataset - E.P(D, T)~~

how large is your partition wrt the dataset

$$\text{Info Gain} = H(D) - \sum_{i=1, -n} P(D_i)H(D_i)$$

↓
entropy of the whole dataset

$$\text{where } H(D) = H(p_1, p_2, \dots, p_m) = - \sum_{i=1, -m} p_i \log p_i$$

↓
summation over probabilities

entropy - measure of information

$$H(D) \rightarrow \text{entropy before split}$$

$$-\sum_{i=1,2,\dots,n} P(D_i) H(D_i) \rightarrow \text{expected entropy after split}$$

Thus info gain means how much info we are getting after the split. Thus, we need to maximize the info gain.

Thus we calculate for which attribute we are getting the max. info gain and take that for splitting.

$$\log 0 = -\infty \quad] \text{usual case}$$

$$0 \log 0 = \text{undefined}$$

But in the formula $0 \log 0 = 0$.

• 20-question game:

Is it a living thing - higher entropy, probability ~ 1

is it dust - small entropy value, prob ~ 0 .

∴ living thing ✓

(2) CART. \rightarrow split into 2 children always

$P(C_i|D_1), P(C_i|D_2)$ - Probability of particular class in partition 1 & 2 respectively.

If their difference is good enough \rightarrow then good split.

~~lec 5~~
~~doubt~~

most pure

Date:

Page No.

Shannon's Entropy

③

C4.5

- Handles missing data → can put average of the nearby/related values in case of numerical attributes; and can put majority vote/answer from neighbouring value in case of categorical attributes → just like KNN
- Handles numerical attribute split
Sort the values. And then use binary search on the sorted attribute values to find the best split point.
- Pruning
- * Prepruning - Helps to avoid overfitting of data. The split stops when the nodes become almost pure. To avoid the decision tree from doing split based on specific attributes like name, roll number etc, prepruning is done.

If it is not done, accuracy will be very high on training set but bad on testing.



Time consuming ↘
Post pruning - In this, the Decision Tree is made ~~to~~ until we achieve all the pure nodes. ~~Now~~ Here, the dataset is divided into training, validation & test. DT is made from Training. Now, accuracy/error rate is checked before removing a pure node & after removing it. If error rate/acc. doesn't change much, it is removed. And it continues until accuracy is good.

- Rules - Instead of DT, rules shall be made for each path (from root of DT to all the leaf nodes) in words/text.
- Goodness of split:
 - Info gain favours attributes with many values. So might split based on roll number, names etc. Thus ~~info~~ gain-ratio should be considered.

DT + Bagging/Boosting → one of the best classifiers.

→ next page

Q) How to get ensemble DT classifiers:

Incorporate randomness (for eg.) in the following ways:

- i) Instead of picking up the best attribute for split, choose a random "n" among the good attributes (based upon some value like if $\text{info gain} > n$, then those attributes are good for split)
- ii) Take ^{smaller} random dataset ^{from entire dataset} & repeat the process to get many DTs.

Above points will help us get a number of DTs with small tweaks - Random Forest.

And now Bagging / Boosting can be applied.

Q) DT Complexity : $O(N \times f \times d)$

for visiting all the nodes, we traverse all the features, getting each probability/distribution of samples. And then we go through nodes, we go through all the samples and traverse all the features of the tree (until we get pure/ almost pure nodes).

N = Number of samples
 f = number of features

d = depth of tree

(\Leftrightarrow no. of interior nodes basically)

very expensive

~~31-8-20
Sec 6:~~

~~Doubt~~ \Rightarrow SLIQ

Date

Page No.

Date:

Page No.

Suppose keyword is "free". It will have some prob. in spam & some prob. in not spam. Date: Prob(free | spam). Prob(free | not spam). But Prob(not free | spam) > Prob(not free | not spam). Prob(free | not spam). And more such words. And hence

Bayesian Methods

Naive Bayes

overall prob. for a record $P(C_i | X)$ gets affected accordingly.

$$X = (x_1, x_2, \dots, x_m)$$

Consider the classification as spam & not spam email.

Let x_1, x_2, \dots be words in the email like free, congratulations etc.

$P(C_i | X) \Rightarrow$ Probability of a word being a class (spam or not spam) given a record X .

$$P(C_i | X) = \frac{P(X | C_i) P(C_i)}{P(X)}$$

will be same class + for every class + such doesn't depend up so remove proportion

$$P(C_i | X) \propto P(X | C_i) P(C_i)$$

$$P(C_i | X) \propto P(x_1 | C_i) P(x_2 | C_i) \dots P(x_m | C_i)$$

probability of a keyword on a particular class

we can multiply like this

on the assumption that all the keywords x_1, x_2, \dots, x_m occur independently of each other.

- Naive Bayes is supposed to work nicely for text classification.
- Calculations might be wrong due to the above assumption but classification is good.

~~bees are newsworthy~~

The assumption of keywords being independent is bad. It gives approximated prob. values.

Eg: University & degree are 2 different words yet dependent.

To overcome this, we have

⇒ Bayesian Belief Networks → ^{mathematically robust}

Exponential problem → every keyword dependent on every other keyword (N^N)

We create conditional Probability Tables for the ~~other~~ nodes to which we have incoming edges.

The table is filled in 2 ways:

- i) Either ~~look for~~ calculate using the given dataset (if it is proper)
- ii) Domain experts can give values and we can use them directly.

$$P(X) = P(x_1 | \text{Parents}(x_1)) P(x_2 | \text{Parents}(x_2)) \dots P(x_m | \text{Parents}(x_m))$$

$$P(X) = \left(\begin{array}{l} \text{Eg: } P(\text{Positive XRay}, \text{Dyspnea}) \\ \quad \quad \quad \downarrow \quad \quad \quad \downarrow \\ \quad \quad \quad x_1 \quad \quad \quad x_2 \end{array} \right)$$

~~Now~~ when we get $P(X)$, we consider then calculate $P(c_i | X)$ as calculated in ~~vec 2~~ Naive Bayes Method.

What is this here?

Vec: 2-9-20
Doubt

2017
04 09 21

Date:

Page No.

Maximum Entropy Approach

CLUSTERING

↳ needs a similarity function

Measurement of similarity

→ Categorical variables

$$d(x, y) = 1 - \frac{m}{n}$$

↓
a record ↓
some other
record!

Jaccard coefficient
/ Jaccard similarity.

$$d(x, y) = \text{similarity b/w } x \text{ & } y$$

Jaccard similarity - Given 2 sets.

Take intersection of the 2 sets divided by union of the two sets i.e.

$$\frac{|A \cap B|}{|A \cup B|}$$

tells overlap.
how many total variables in 2 features

→ Numeric variables - Euclidean distance

$$\text{can be used } d = \sqrt{(x_1 - x_2)^2 + (y_1 - y_2)^2}$$

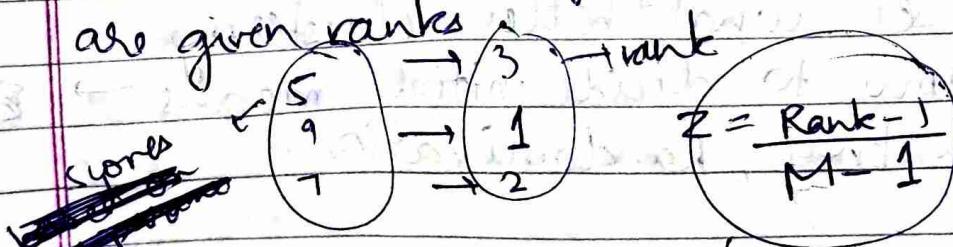
$$\text{Manhattan dis} = |y_2 - y_1| + |x_2 - x_1|$$

$$\text{Minkowski dis} = \text{n}^{\text{th}} \text{root}(|y_2 - y_1|^n + |x_2 - x_1|^n)$$

if $n=1 \Rightarrow$ Manhattan distance

$n=2 \Rightarrow$ Euclidean distance

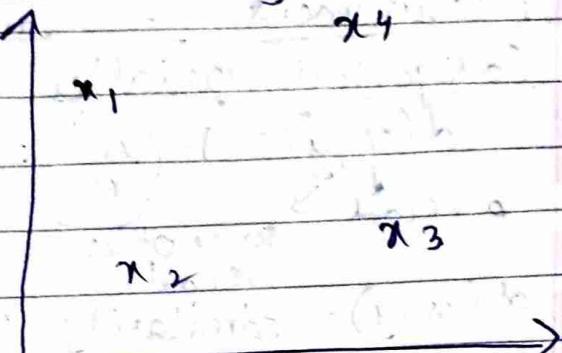
Ordinal variables - Where the exact values are not meaningful. These numbers are given ranks



and then we can use Euclidean or any other dist. formula.

z -score \rightarrow b/w -1 to 1.

\Rightarrow K-Means (Partitioning Based)

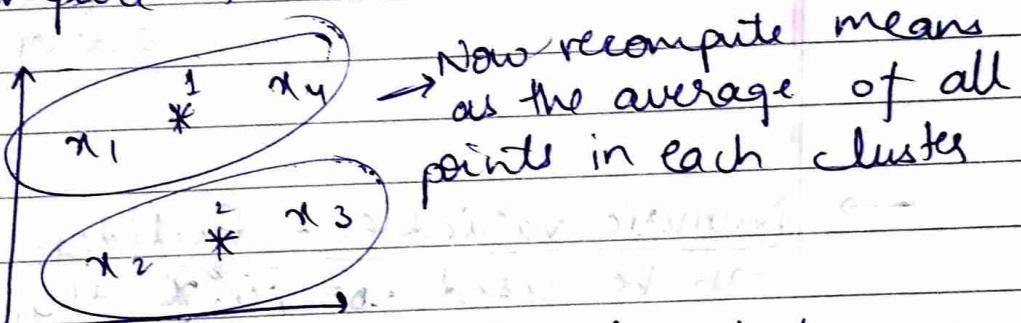


choose $k = 2$. Initial means $\Rightarrow x_1$ & x_2 .

Now find dist of x_4 with x_1 & x_2 .

Now find dist of x_3 with x_1 & x_2 .

Compare the distances & cluster.



Now here for all the 4 points, we will check the dist & from 1 & 2.

make 2 clusters accordingly.

If cluster formation does not change, we stop.

How to choose k ? \rightarrow Try different k &

see what gives better clustering

How to divide initial means \rightarrow Nit & try; Randomization.

There is a theorem which states the algorithm converges. But if there are many iterations before converging,

Suppose there is an outlier & we force it to be a part of a cluster. Then quality will decrease as the distance of that point with the mean will be large & thus it will be magnified. And we can put an upper limit to the kind of number of iterations.

Now to know we get the best clustering? \rightarrow Evaluate clustering quality.

$$E = \sum_{i=1}^N \sum_{x \in C_i} d(x, m_i)^2$$

squared error. \rightarrow we need to minimize this for better clustering. intuitively a measure of how far away things are from their means.

K-Means works only for numerical attributes.

Work around:

i) Drop categorical & do clustering based only on numerical

(doubt) ii) Take numerical separately & categorical separately. calculate distance for both. & then give different weights to each of them (if required).

Then combine these two distances with a weighted summation.

→ stored as bit strings.

One hot encoding is used for converting categorical to numerical because it are computationally very efficient. (inbuilt hardware operations etc. are computationally ~~not~~ very fast.)

\rightarrow 1 → records which have Red
 \rightarrow 1 → " Daten Green
 Page No.

eg (Red: 11010011000)
 Green: 01010001001

Red & Green → have both R & G.

→ This Bitwise operation (AND) is very fast.

This one hot encoding works nicely for smaller sized datasets.

In case the datasets are large, this notation can be converted into only representing 1s i.e [1, 2, 4, 7, 8, 12] records that have Red.

⇒ Hierarchical

① Agglomerative → each cluster is a part of some other cluster.

Q How to merge 2 close clusters & ie how to find distances b/w clusters?

1. Single link - $O(N^2)$ → dist b/w 2 points & that are closest in 2 clusters
2. Complete Link - $O(N^2)$ → " " " farthest "
3. Average link - $O(N^2)$
4. Mean Link - $O(N)$ → mean of Dist b/w means of 2 clusters.

~~In agglom.~~

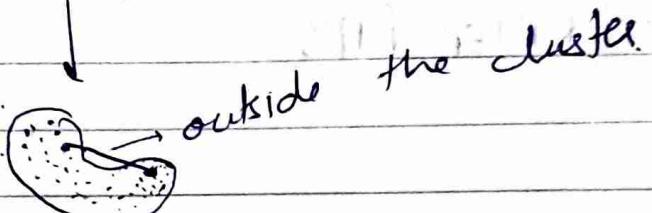
In agglomerative or divisive, we can set thresholds to stop in p/w rather than getting a single large cluster or all single points as clusters respectively.

⇒ Density Based → dense region should become a cluster.
DBSCAN

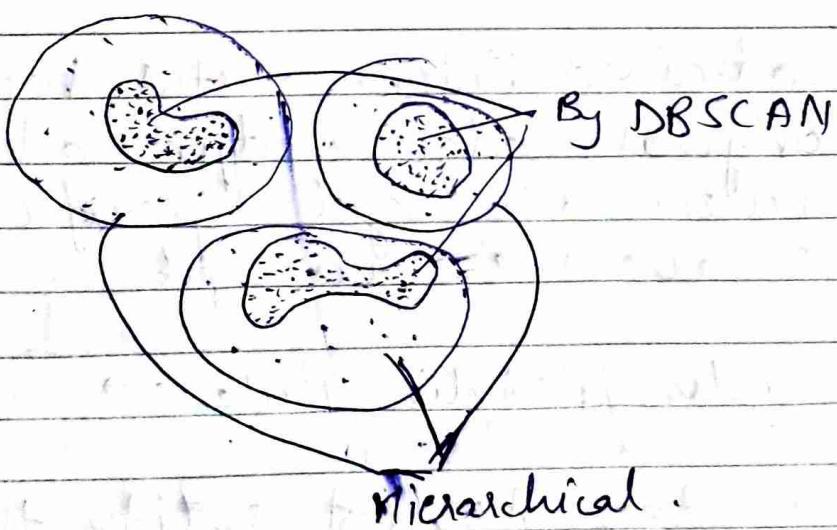
K-Means & Hierarchical clustering Algorithms gave spherical / convex clusters (the line b/w 2 points lie in the cluster itself).



Whereas DBSCAN ~~helps~~ is a way in which all clusters need not be convex.



Some enhancements to DBSCAN:
 DBSCAN + Hierarchical.



⇒ Mining Outliers using Clustering

↳ can be found by simply running a clustering algo. If any point or a smaller cluster is left, treat it as outlier.

There are techniques to mine outliers which does not require clustering.

21-9-20

Date:

Page No.

⇒ Association Rules

④ Naive algo ~~for~~ for frequent itemset
↓

Keep counters. Then again
scan & find support.

Exponential search ~~is~~ $\rightarrow 2^n$
where $n = \text{no. of items}$.

Doubt

⇒ KILLER APPS

- Constrained Rules - useful bcz the original algos output a lot of rules. So to get specific rules, we use this concept
- Cyclic / Periodic Rule → Time associates

we try to first satisfy the support & then try to increase the confidence.

Doubt

Interest & Conviction

Q How to choose minimum support value?

→ Choose different values say 10%, 20% -- 50%. etc. on X-axis.

And on Y-axis, output the ~~frequency~~ number of itemsets.

Low min sup → Higher num. of itemsets

High " " → Lower " "

* As support gets closer to 1, we get 0 number of itemsets (simply means all the ~~itemsets~~ records containing all itemsets which is almost impossible.)

A knee ^{point} in the graph can be a good minsup.

Q How to modify datasets to ~~avoid~~ maintain data privacy

Ans: One way can be to remove attributes like name, id, phone number etc. because they anyway does not help in data mining

- ① Generate frequent items using Algo
 ② Generate association rules

Date:

Page No.



Frequent Itemset Algorithms

Apriori Algorithm

Idea: An itemset can be frequent only if all its subsets are frequent.

in slides

$D \rightarrow ABC \rightarrow$ negative borders itemsets.
 i.e. their subsets are frequent itemsets but they are not.

∅

A, B, C, AB, AC, BC.

subsets

minimal

infrequent itemsets.

(i.e. they are infrequent but could

proceeds level wise.
 start with level 1.
 so candidates of have been frequent &
 generate level 2. so they also need to be
 (generate level 2. check)
 then make a scan for level 2 & generate for level 3.

Apriori Algorithm → Bottom-up.

i.e. first check for ←
 singletons and then look at its supersets & so
 on

Pruning step - If immediate subsets
 are not in f, prune Z.

$$F = \underline{ABC} \subset \underline{ABD}$$

prune bcz ACD
~~ABD~~ ABCD → is not in F.

here imm. subsets of ABCD will be the subsets of length 3. We need not check for levels below 3 bcz we are already proceeding in bottom-up approach checking at every level.

This algo is fairly efficient bcz we are not generating any candidate → It generates the minimal candidates.

While implement, at the 2nd scan, the algorithm becomes very slow (pruning cannot happen) suppose F has 600 singletons. Now $F \cup C_1$ will be new candidates generated and they ~~cannot~~ ~~will~~ not be pruned.

So for 1st pass: Use a bucket/count sort technique to store the frequency of singletons.
They ~~will~~ will be pointing to a triangular matrix which is to be used for 2nd pass

in slides: White ones \rightarrow freq.
violet ones \rightarrow infrequent.

$A[b_i]++$.
 $\hookleftarrow e \in t$ (transaction)

$t' \Rightarrow$ will have all frequent singletons.

Singletons - 1D Array
2 items in set \rightarrow 2D Array - lookups for fast lookups
But for 3 items in a set like (ABC, BCD) etc.

we cannot have a 3 dim. array bcoz now that structure will be sparse enough and thus waste of memory So we use hash trees from 3rd scan onwards.

Hash Tree Algo \rightarrow do from slides.

②. The Partition Algorithm

- In 2 DB scans, we get frequent itemsets.
- Divide the DB into some n partitions such that each partition can fit in the main memory (along with the data structures needed to implement this algo.) Thus, there will be no disk access & hence efficient.

on getting freq. itemsets from each partition, do $F = F_1 \cup F_2 \cup F_3 \dots$

in the end, only those F_n .

itemsets which are frequent in the whole database are left.

Thus this algo claims that an itemset infrequent in each partition is infrequent in entire database.

Proof : in slides.

(3)
Doubt

Sampling → ~~get~~ got an approximation of frequent itemsets instead of the exact freq. itemsets because we consider a random sample of the entire dataset & find F_s (freq. itemsets) & N_s (negative borders/minimal infrequent itemsets).

But if we want to get the exact answers using F_s and N_s :

Here also we have 2 DB scans / Pass.

(4)

Incremental Mining Technique

→ 1 Scan Algorithm

Date:

Page No.

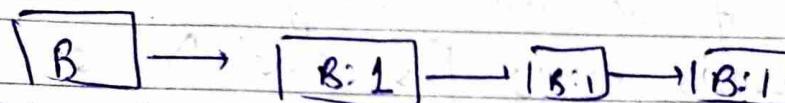
→ FP-Growth → Does data mining without candidate generation.

Follows a trie kind of tree structure.

The header table has 2 things:

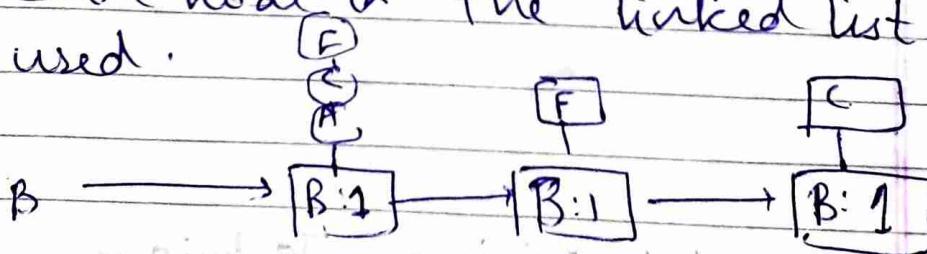
- What is the count of each letter
- Pointer to the first node and so on.

Like link all the similar ~~items~~ items.



FP-Tree (Frequent Pattern tree) → Tree + Linked list.

Conditional Pattern Base → Go through the linked list & the prefix of each node in the linked list is used.



So Pattern base of B is

FCA : 1 → values of these B.

F : 1
C : 1

~~Denote dataset~~ →

In the FP-Growth algo, we have conditional pattern base of m levels like first for singletons. Then for pairs & so on as this algo is recursive.

→ does counting amazingly.

Date:

Page No.

So somehow we are generating the candidate sets.

Dense dataset - when there is a lot of sharing of prefixes or commonality.
(Long prefixes get shared.)



we have frequent itemsets. We are just counting

→ ORACLE / ~~A~~

DAG structure: has 2 links from which a node is formed:

Tidlist: All the transactions that contain of an itemset this itemset in the current partition.

Eg: for A, → AB, AC, AD, ABC -- etc

Eg: for AB → we'll take tidlist(A) ∩ tidlist(B)

if eg: AB

In intersection, we convert one into tid vector (1 represents if A is present in that i^{th} record (ith bit) & B is represented as a tid vector).

so the total cost = $O(|\text{Tid-vector}|)$

but we scan this & directly check for that i^{th} position

in $O(1) \rightarrow$ look up is

$O(1)$

⇒ ARMOR

L, we might add or remove items
in the candidate → meaning of Update
C in first Pass algo part.

We have account of frequent itemsets
along with their partial support values
~~as soon as~~ as soon as an item becomes
part of candidates set.

Why second Pass? Suppose an item
~~came into candidates~~
~~became frequent~~ from partition
number 3. But it might be
present in 1 & 2 also.

~~It wasn't added until~~
~~it would be having~~
~~support > minsup.~~

So now we will traverse
the whole database & to
get the actual count of such itemsets.

And if the actual count ~~> minsup~~ & add
~~it to f else delete it.~~

So C will convert into f

candidates

frequent itemset

⇒ Implementation

Concise Mining

⇒ Post-mining Rule Pruning Schemes

• Improvement: Eg. $ABC \rightarrow D$
 $AB \rightarrow D$. → like a
If the confidence of sub-rule of the
these 2 are very similar (very less above.
difference), then we need not report
both of them.

⇒ Constrained Association Rules are good if the user has some constraints in mind. Else it is not good.

⇒ Maximal frequent Itemsets

↳ If we are o/p these positive borders,
then all its subsets will also be frequent → so do not o/p them. But this will have drawbacks.

Suppose ABCD is max. freq itemsets and the rules are formed by breaking them say $AB \rightarrow CD$. But if we remove the subsets, then we will lose their support & could not make these rules.

⇒ Closed Itemsets

$\text{Support}(A) = \text{Support}(AB)$ means both are appearing in the same records.

→ Mining Real World Data

Binning - Try to make sure that bins are equi-width & equi-depth.
→ same length range.
similar frequency of ~~the~~ values in
of every range.

Relational → directly suitable
for clustering and classification
Transactional → directly suitable
for association rules mining.

But to convert transactional to relational,
we would want to normalize ~~so~~
④ (3NF & 4NF forms) & then could
be converted to relational.

For classification or clustering,
we need to do joins of the various
relational tables. ④ to get a single
table.

Multi-dimensional -

Distributed - Like relational but the data is distributed across multiple nodes (maybe across data centres.)

- Maybe you can bring all the scattered data to a central location & do data warehousing.
- Maybe mining separately & doing analysis by getting the results to a central location (just like partition.)

Distributed

different data

at diff.
locations
(enterprise)

Homogenous

Run like partition

Heterogeneous

Do join & bring
to central location

some
data at
all locations

⇒ IR Concepts

Boolean queries → OR Query, AND Query etc.

Ranked queries → Among the retrieved query documents, which are top ranked and so on. (Trying to rank the output.)

→ Text & Web data models → semi-structured or unstructured data models;

Q How to do data mining on above data model?

Ans i) Bag of Words:

word 1: count 1, word 2: count 2,
word 3: count 3.

Bag of words model

ii) $[count_1, count_2, \dots, count_n]$ where
↳ Vector space Model

words are features
and their order is
fixed so we need
vector of n dimension to just keep the counts.

iii) ~~Cosine~~ Cosine similarity - To find similarity b/w 2 documents.

$$\frac{d_1 \cdot d_2}{\|d_1\| \|d_2\|} \rightarrow \begin{array}{l} \text{dot product of} \\ \text{2 vectors} \end{array}$$

dot product:

$$x = (x_1, x_2, \dots)$$

$$y = (y_1, y_2, \dots)$$

$$x \cdot y = (x_1 y_1 + x_2 y_2 + \dots)$$

$$\text{norm} \Rightarrow |x| = \sqrt{(x_1)^2 + (x_2)^2 + (x_3)^2 + \dots}$$

iv) Entropy - Treat count of words as probabilities

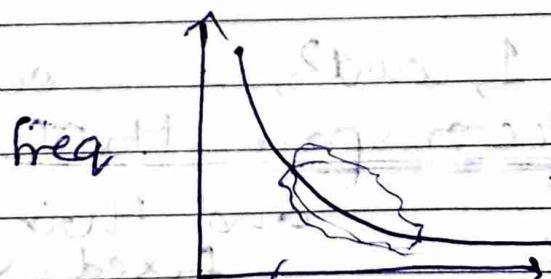
$$H = - \sum p_i \log p_i$$

v) Stop words - Words like "and", "the", "to" occur too frequently and are irrelevant in classification or clustering. So we need to remove them.

Zipf Law \rightarrow Most freq word \rightarrow Rank 1

Universal law Second " " \rightarrow Rank 2

Strong universal
for all languages.
freq



So now we can use some threshold based on freq &

Rank to decide what all words to consider.

Mostly the middle words are used. which lie in this range maybe

vi) Stemming - Take the root word & stem out others.

Eg: Sing, sung, sang, ~~sing~~, simply
Just take ~~sang~~.

helps in reducing vocabulary space.

vii) TF-IDF - Tells about what are the most important ^{best} words in the documents.

viii) Inverted Index - for each word, we store what all documents contain it.
 Eg: word 1 → doc1, doc3, doc4, etc.

(Inverted index usually like store what all words are present in each doc.)

⇒ Search Engine Architecture

i) Crawling - getting links from links and follow those documents.

Normalized ↑ Tokens

ii) Normalization → same word can be written in different formats like U.S.A, USA, United States of America etc. so we need to normalize & get them. We get single token ^{for} from them & use it.

iii) Ranking

a) TF-IDF → Term frequency — Inverse document frequency.

~~mean(tf)~~ ~~max(fif)~~ \rightarrow freq of all words.

DATE	_____
PAGE NO.	_____

TF - num. of times the word appears in the doc.

df - how many docs contain that word.

good keyword - Appear in this doc many times (tf high) but should not occur too many times in other docs. (df is low)

b) Page Rank

if ~~all~~ the docs linked to say doc D are highly ranked, then D is highly ranked.

Talk about inlinks.

$$\text{Page-Rank} = \text{SUM} \left[\frac{\text{rank}(y)}{\text{# links}} \right]$$

all the inlinks / outlinks

of that particular doc.

Initially, all docs have ~~rank~~ rank = 1.