

7/1/2019

Efficient algorithm { Given T; generate all possible freq itemsets }

Apriori → too many repetitions → bottom up approach.  
fp growth → 2 scans.

Last step of Apriori is  $C_{i+1} = \emptyset$ .

- Upward closure - if an itemset is infrequent then its immediate supersets are infrequent.
- Downward closure - exact opp. of upward closure.

No of scans in Apriori =  $|C_2| + |C_3| + \dots + |C_n|$

Time complexity of apriori -

- How is fp-growth better?  
2nd scan is to generate fp-tree.
- Sequence pattern mining in Apriori → Generalised Sequence pattern min.

Apriori is adjustable in SPM domain.

fp growth is not suitable for SPM.

Fpgrowth was developed by Han + Team (2010).  
Fpgrowth is not for sequence pattern mining.

- Fpgrowth for sequence pattern → PrefixSpan.  
SPM used in web patterns. Generic to specific websites.

Yes ← Fpgrowth, does it require only 2 overall scans?

Tree is referred as the projection of the database.  
Projection is to support efficient counting.

DIC was about parallel counting.

Where do more countings occur in fp tree?

$P_i, S_j$  might occur in many paths. So the subset merging requires smaller scans (not the entire DB). So it suffers from repeated scans.

Subset merging:- Working at sequences where  $P_i$  is present.

Both Apriori & fp-growth are time complex.

DIC is space complex.

Fp-growth is a good balance b/w Apriori and DIC.

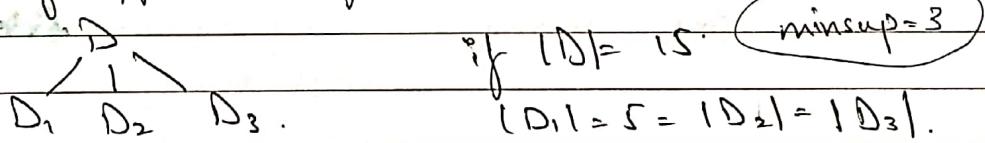
Pinell Search (Maximal); A-close. (Closed FP).

Apriori is good for sparse dataset. That means when there are many infrequent items.

- Can we try improving the efficiency of Apriori?

Apriori is a balance b/w space & time complexity.

- Partitioning approach of Apriori.



local mining approach is done. (local Apriori).

To arrive at the global count we need merging of  $D_1, D_2$  &  $D_3$ . (Union all three).

26/11/2019

In partitioning, something that is locally infrequent might be globally frequent.

Do union.

Hash pruning approach → in early levels  
for counting. Fit in the  $C_2$  level.

→ Why not in higher levels?

Occurance of the itemset with higher length  
will come down.

1, 2, 5

2, 4

2, 3.

Misup = 2.

referred

in ECLAT

1, 3.

1, 2, 3.

1, 3.

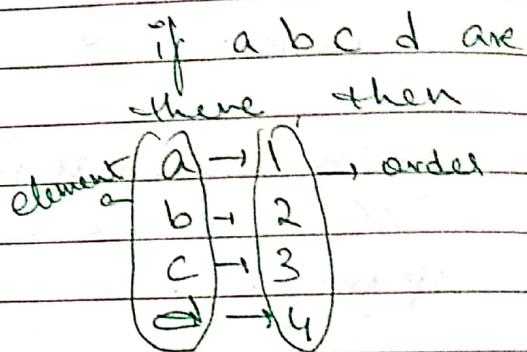
1, 2, 3, 5

1, 2, 3.

$L_1 = \{1, 2, 3, 4, 5\}$

$C_2 = L_1 \times L_1$

$C_2$   
 $L_2$      $C_2$   
freq infreq



1 1  
2 2  
3 3  
4 4  
5 5  
1 2 23 34 45  
1 3 24 35  
1 4 25  
1 5

Instead of going thru the whole dataset, we look at the hash table. In this way Apriori with hashing is better.

$C_2$  is for 2 itemsets

$$H(x, y) = (\text{order of } x \times 10^3 + \text{order of } y) \% 7$$

Bucket → 0 1 2 3 4 5 6      for (1, 4)

(1, 4) (1, 5) (2, 3) (2, 4) (2, 5) (1, 2) (1, 3)

Contents → (3, 5) (1, 5) (2, 3) (2, 4) (1, 2) (1, 3)  
(2, 3)

$$H(1, 4) = 10 \times 1 + 4 \\ = 14 \% 7 \\ = 0$$

2 2 4 2 2 4 4

Repeat for all. Collision happens only for infrequent.  
 $S_{C_2}$  is better.

In case of collision, refine the hash function

- Transaction Reduction:-

If  $\text{P}_2(1,2)$  is infrequent, then remove the transaction from the Database.

1, 3, 7       $M_S = 3$       1 (5)

2, 3, 7      2 (7)

1, 2, 3.      3 (9)

2, 3.      4 (3)

2, 3, 4, 5

2, 3.       $L_1 \neq C_1$

1, 2, 3, 4, 6.

2, 3, 4, 6.

1, 3.

1

If an itemset is infrequent in the earlier levels, remove it from the later levels.

5, 7, 6 are infrequent. Remove them from the Dataset.

So now dB is,

1, 3

2, 3.

1, 2, 3.

At least one character is being thrown out.

2, 3, 4.

(Time complexity)

2, 3.

The order might remain same.

1, 2, 3, 4.

2, 3, 4.

Apriori's complexity is exponential.

1, 3

1

9/1/2019

Transaction reduction is also level based.

Transaction reduction is not efficient as an algo, but implementation wise it is efficient.

ECLAT $\rightarrow$	$f_i$ 's
Aclose $\rightarrow$	$c_f$
Pincel $\rightarrow$	info

### \* ECLAT - Equivalence Class Lattice Travelal.

Breadth First

Apriori is BF traversal.

ECLAT  $\rightarrow$  vertical representation of DB  $\rightarrow$  focus on item

Horizontal focus on Transactions.

Where in item traversal	$I_1 \rightarrow \{T_1, T_4, T_5, T_7, T_8, T_9\}$
diff	$I_2 \rightarrow \{T_1, T_2, T_3, T_4, T_5, T_8, T_9\}$
Same for everything	$I_3 \rightarrow \{T_3, T_5, T_6, T_7, T_8, T_9\}$
	$I_4 \rightarrow \{T_2, T_4\}$
	$I_5 \rightarrow \{T_1, T_8\}$

MC = 2

'n' in order (intersection in order)

min sup

✓	$I_1, I_2 = \{T_1, T_4, T_8, T_9\}$ .
✓	$I_1, I_3 = \{T_5, T_7, T_8, T_9\}$ .
✗	$I_1, I_4 = \{T_4\}$ .
✓	$I_1, I_5 = \{T_1, T_8\}$ .
✓	$I_2, I_3 = \{T_3, T_6, T_8, T_9\}$ .
✓	$I_2, I_4 = \{T_2, T_4\}$ .
✓	$I_2, I_5 = \{T_1, T_8\}$ .
✗	$I_2, I_4 = 0$ .
✗	$I_3, I_5 = \{T_8\}$ .
✗	$I_1, I_5 = 0$

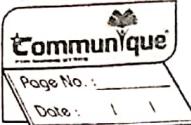
$\rightarrow I_2$  in terms of Apriori.

ECLAT is better than Apriori and FP growth on counting.

No need to go to the dataset at all.

$I_1 I_2 I_3 \rightarrow S_1 S_2 S_3$

paper.  
C Baet Goethals.  
father of ARM



$$I_1 I_2 I_3 = \{T_8, T_9\}$$

$$I_1 I_2 I_5 = \{T_1, T_8\}$$

CHARM  
MAFIA

Read about them

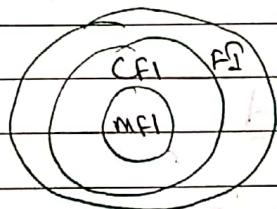
- What was the need for CF1, MF1 and others?  
Rare is not infrequent entirely. Sensitive & high utility itemsets are other types of itemsets.  
Hiding sensitive itemset  $\rightarrow$  Privacy preservation  
data mining.

16/11/2019

Final

A-close  $\rightarrow$  Apriori-close.

- Why do we need CF1, MF1 & LF1?



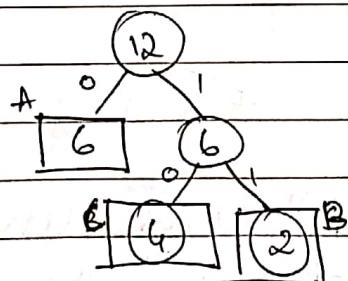
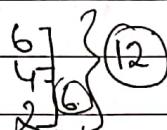
The application is not happy with so many FI's. So shrink it down.

Closed is from closure. The result should belong to the universe of Discourse. Here FI is the universe of discourse. CF1 is lawless, MF1 is lossy.  
Lawless  $\rightarrow$  from CF1 can we reconstruct FI? Yes..  
gzip, pkzip, RLE (run length encoding)  
Huffman trace was taught.

ABABABABCCCCC.

A  $\rightarrow$  6 B  $\rightarrow$  2 C  $\rightarrow$  4.

descending order of frequencies.



D  $\rightarrow$  text.

for A  $\rightarrow$  0.; B = 11.; C = 10.

The letter with more frequency gets the more level code length.

Biology + Data Mining = Bioinformatics.

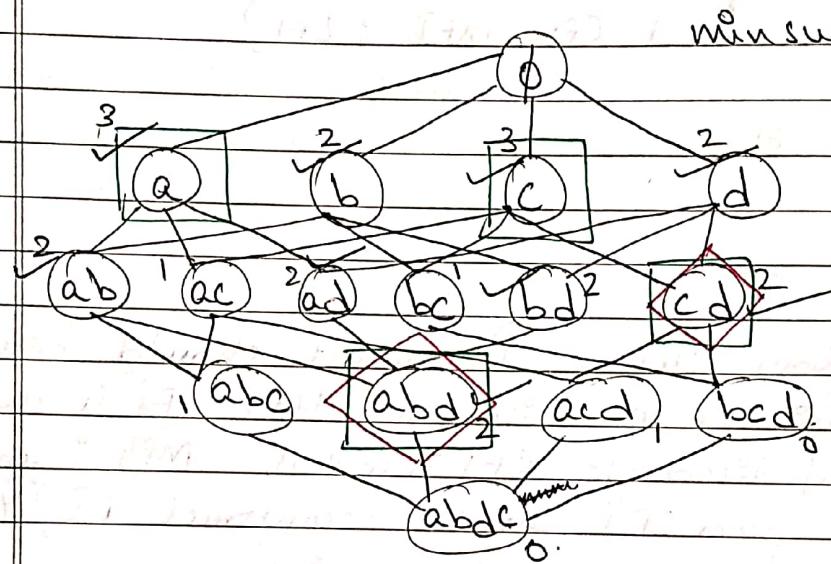
Decoding should be unambiguous. (Prefix free encoding)

(FF, MFI find applications in compression)

Voice and image compression can be lossy.

$$CFI = A(s) \cdot AC(s)$$

17/1/2019



✓ → FIS. (9)

◻ → CFI's (4)

◆ → MFI (2).

CFI is lossless, MFI is lossy.

are we getting the entire fi from cfi's.

MFI is more about IFFI (cardinality)

CFI is a tighter constraint on MFI.

MFI also comments on the subsets.

We will not be able to say the subset count from MFI. Confidence calculation won't be possible.

With CFP we can tell the support count of the subsets.

$x \rightarrow$  if there is no superset of  $x$  with the same support count as  $x$ . (or less support count)

$\rightarrow$  CFP

Superset ( $x$ )  $\leq$  SC( $x$ )

$x \rightarrow \nexists$  any superset of  $x$  which is frequent.

$\rightarrow$  MFP.

- Difference between closed itemset and closed frequent itemsets.

### \* A - "close" Algorithm :-

S items (A, C, T, D, W). A C T W.

$\Rightarrow$  all the transactions that contain. 2 C DW.

$t(x)$  itemset  $x$ . ~~if  $x \subseteq T$~~  3 A C T W.

4 AC DW.

$i(y) \cdot y \subseteq T$

5 AC D T W.

6 CD T.

$t(c) = \{1, 2, 3, 4, 5, 6\}$ .

$i(c)$  = undefined.

$i(T_1, T_2)$  = intersection of  $1 \notin 2$  = C, W.

$i(t(c))$  = C

$t(c, D) = \{2, 4, 5, 6\}$ .

$i(t(c, D))$  =

$t(A) = \{1, 3, 4, 5\}$ .  $\left\{ i(+A) \neq A \right.$  unlike the

$i(+A)) = \{A, C, W\}$  examples before this.

$i(+ACW) = \{A, C, W\}$ . it is equal to a greater set.

Do,  $i(+A, C, W) = \{A, C, W\}$ .

Do  $t(n)$  then  $i(+n)$ . If  $i(+n) \neq x$  but is equal to y, Do  $i(+y)$ . Here it has to stop.

1 level we have elements gaining.

- (i)  $i(t(x)) = x$ .
- (ii)  $i(t(x)) = y ; x \subset y \rightarrow$  elements gained.
- $i(t(y)) = y$ .

- Generators :-

Identify those elements that act as generators.

$x$  is gen(y) if  $x^+ = y$ .

generator of  $y$       closure of  $x$  yields  $y$

In the example above, A is a generator of ACW.

In CFP we are interested about minimal generators.  
There is no proper subset of  $n$  generates  $y$ .

What all generates ACW.

ACW and A.

We cannot have more than 1 minimal generator.

21/1/2019

### A-Close Trace.

1 ACD

2 BEF       $mS = 3$ .

3 ABCE

4 BE

5 ABCE

$$C_1 = \{A, B, C, D, E\}$$

$$L_1 = \{A, B, C, E\}$$

$$C_2 = \{AB, AC, AE, BC, BE, CF\}$$

With CF we reconstruct the FI's but with closed itemset we can traceback to the transactions.

$$L_2 = \{ \underset{3}{AC}, \underset{3}{BC}, \underset{4}{BE}, \underset{3}{CE} \}$$

In any item having the same support as it's subset, here  $AC \not\subseteq A$  and  $BE \not\subseteq B$ . Throw them away.

$$L_2 = \{ \underset{3}{BC}, \underset{3}{CE} \}$$

$$C_3 = \emptyset$$

Since we left out  $AC$  &  $BE$  we must prove that  $AC$  can be reached from some other minimal set.

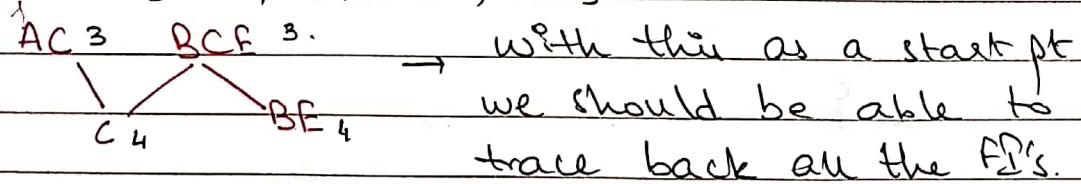
### (a) Generators

	$t(x)$	$\rho(t(x))$
A	1, 3, 5	$AC \rightarrow$
B	2, 3, 4, 5	$BE \rightarrow$
C	1, 2, 3, 5	$C$ We are able to
E	2, 3, 4, 5	$\times BE$ reach them here.
BC	2, 3, 5	$BC \rightarrow$
CE	2, 3, 5	$\times BCE \downarrow$

Whatever is left out is able to be reached from  $\rho(t(x))$ . Reaching out in mathematical terms is closure  $x^+$

Remove the duplicates in order.

$$CFF's = \{ AC, C, BCE, BE \}$$



By not retaining the higher cardinality sets, we are saving space. The cardinality of the Generators should be minimum.

One CFF alone cannot give all the F's. Go for the CFF with the large count.

topdown and  
bottom up.

\* Pincer Search.  $\uparrow \downarrow$  "Dataset is given"

Support = 20%, relative support = 15.

$$\therefore \text{min sup} = 3.$$

$\{1, 2, 3, 4, 5, 6, 7, 8, 9\}$  9 symbols.

Move in both bottom up and topdown.

maximal freq set  $\leftarrow \text{MFS} = \emptyset$

maximal freq candidate set  $\leftarrow \text{MFCS} = \{1, 2, 3, 4, 5, 6, 7, 8, 9\}$ . (Maximal frequent candidate set).

$S_i \rightarrow$  infrequent at  $i^{\text{th}}$  level.

$$L_1 = \{2, 3, 4, 5, 6, 7, 8\}.$$

$$S_1 = \{1, 9\}.$$

We can keep record of infrequent to see that they are not part of MFS.

Keep leveling MFCS

Trimming down MFCS is top down approach.

MFCS and  $S_i$  are helping to bring in the top down flavour to the trace.

updating MFCS.  $S_1 = \{1, 9\}$ . MFCS =  $\{1, 2, 3, 4, 5, 6, 7, 8, 9\}$ .

for  $\{1\}$  in  $S_1$  :  $\{2, 3, 4, 5, 6, 7, 8, 9\} = \text{MFCS}$   
 $\text{MFCS} \setminus S_1$   $\xrightarrow{\text{minus in set theory.}}$

for  $\{9\}$  in  $S_1$  :  $\{2, 3, 4, 5, 6, 7, 8\} = \text{MFCS}$   
 $\text{MFCS} \setminus S_1$

$$C_2 = L_1 \Delta L_1$$

$$L_2 = \{23, 24, 35, 37, 56, 57, 67\}$$

23/1/2019

$$S_2 = \{25, 26, 27, 28, 34, 36, 38, 45, 46, 47, 48, 58, 68, 78\}$$

Exhaust  $S_i$

2<sup>25</sup>s are infrequent together. → '26' is occurring together only here.

$$(25) S_2 \rightarrow MFCS = \{ 345678, 234678 \}$$

(26)

$$S_2 \rightarrow MFCS = \{ 345678, 34678, 23478 \}.$$

↑  
is already  
in

So do not consider '34678'

$$\therefore \text{updated MFCS} = \{ 2348, 345678 \}.$$

(27)

$$S_2 \rightarrow MFCS = \{ 3478, 2348, 345678 \}.$$

↓  
This is  
a subset  
of

$$\therefore MFCS = \{ 2348, 345678 \}.$$

(28)

$$S_2 \rightarrow MFCS = \{ 348, 234, 345678 \}.$$

$$\therefore MFCS = \{ 234, 345678 \}.$$

(29)

$$S_2 \rightarrow MFCS = \{ 23, 24, 45678, 35678 \}.$$

MFPI is  
a lossy  
reconstruction  
of FP. We  
do not have  
the sup count.  
So, we cannot  
have the confidence

Continue for the other  $S_2$ 's

(30)

$$S_2 \rightarrow MFCS = \{ 23, 24, 357, 5678, 8 \}.$$

$$\therefore C_3 = \{ 234, 3578, 567 \}.$$

Topdown will stop if  $S_2$  is nullset. There is no more trimming down of MFCS.

Here  $S_3$  is a null set.

## LAB - 2.

Implement A\* search and Page Search.  
Two more algos MAFIA & CHARM.

26/1/2019

- Difference b/w BigData and Data Mining  
How well the storage systems can handle the distributed data. Incremental learning (with added data), the system must be intelligent enough to relearn the patterns). Storage issues according to scaling up of data must be addressed. Scaling up is done across networks.

Online DB = Incremental learning.

Bucket Brigade → Genetic Algorithms

↓  
evolutionary computation

↓  
evolve for the good.

GA is one evolutionary computation strategy.

Tabu search, PSO, Ant colony are others.

particle  
swarm  
optimization.

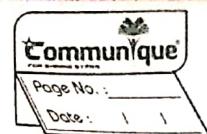
Reinforcement learning thrives on penalty & awards.

Most evolutionary computation algos work on initial population.

All the evolutionary computation methods are basically for optimization (on the computer front).

Search for optimum after seeing if the solution is feasible or not.

David Goldberg  $\rightarrow$  GA book [1 to 3]  
James Holland  $\rightarrow$  Inventor of GA.



- Enumerative search  $\rightarrow$  point by point search.  
Also called try and error. (but it is random).
- Operations Research  $\rightarrow$  mathematical optimization methods.
- Derivative based  $\rightarrow$  Gradient ascent or descent.  
Hill climbing is also another word. A slope should always exist. But real life data is not always continuous. Also, the data can be multimodal, the hill climbing gets stuck in local optimum.

In the search for optimal value, the path towards optima is also important.

Where will GA fit in?  $\rightarrow$  Classification.

If we are not able to give unique class labels in DTree, it is non feasible.

Convergence  $\rightarrow$  how quickly will I reach the leaf node.

Efficiency vs efficacy.

28/1/2019

- Why EVOLUTIONARY?

Can my overall objective keep bettering. GA, Tabu Search etc are search algos that search for optimal.

Fitness of a Solution: - how good is a solution.

Iteration = Generation.

As we do more generations, the fitness should improve.

Darwin's theory of evolution forms the basis of the GA's.

GA → search algo based on natural selection + Genetics.

The power of the evolutionary algos comes in if the search base is really large and the function is too complex.

- Objective function :- it is to be identified.  
This is the first step.
- Initial population :- random collection of solutions

$(2, 6, 8) \rightarrow 3$  random search  
 "Roll of a die" points from search space.  
 Can be one approach  
 of constructing the initial population.

Using the initial population can we do some Genetics (biological ops).

$$f(x) = x^2$$

$$[0, 31]$$

Initial popn = 2, 6, 8  $\rightarrow \max x^2 = 64$ .

Using genetic = 5, 9  $\rightarrow \max x^2 = 121$ ,  
 we have to reach 31.

Human "flair" of a flexible search procedure.

Biological System solutions are far more robust than the existing search procedures.

So, can the innovative search brought into GA?

GA

Selection  
Crossover  
mutation.

know atleast 1 way  
of each.

Application specific operators have improved  
GA over time

Search space = 10k.

larger the search

Initial popn = 20 pts. (randomize the initial popn)

space, larger can

Selection opn :- out of 20 pick only those  
which are fitter ones.

Allele → attributes → chromosomes

Initial popn :- is a collection of chromosomes.

In the probability of fitter solution. Cross parents  
who have fitter characteristic.

After selection of the fit ones, crossover.

Mutation is hardly recommended. It is done when  
stuck in the local optima.

Roulette wheel → simple selection (random)

minimizes the owner's loss.

Initial popn = 13, 24, 8, 19.

$$f(x) = x^2$$

169 576 64 361.

2

3

4

← Chromosomes.

better fit

Sum of fitness = 1170. (wheel dimension) (sample space)

Probability of selection =  $f_i / \sum f_i$ .

$$= \frac{13}{1170} \quad \frac{24}{1170} \quad \frac{8}{1170} \quad \frac{19}{1170} \rightarrow \text{Sum is 1}$$

1      2      3      4.

13 → percent of the wheel is occupied by 13 on the wheel.

0.14	0.49	0.05	0.30
1	2	3	4
14%	49%	5%	30%

Assume that the population size is 4 (fixed).  
So cross over a population of size 4.  
We don't want any copy of 8.

11/11/2019

- How does GA gel well with classification?  
Information gain using entropy, gini index etc. Choosing the optimal path in Dec. Tree. There might be more than one DecTree for one classification. We must be able to choose the optimal path from so many trees.

Expert system forms the base for Bucket Brigade classifier.

Expert system :- Control is given by domain experts ELIZA, ALEXA.

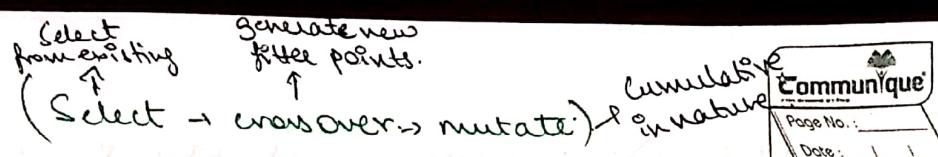
Expert systems = Training model.

Expert systems are basically classifiers. Paths in DecTree is rules in expert systems.

The fitness of the path is important.

GA strikes the balance between efficiency and efficacy.

Efficacy + effective → closer to feasibility.  
The solution must be robust too.



Efficacy in searching is Relevance. It is qualitative measure.

One solution is a chromosome in GA terminology.  
Bucket Brigade is a multi-classified approach.

Look at the past fit solutions. Discourage unfit solutions, encourage fit solutions.

<u>Index</u>	<u>x.</u>	<u>Chromosome</u>	<u>fit.</u>	$\rightarrow$ fitness.
1	13	0 1 1 0 1	169	
$\max(x)$	24	1 1 0 0 0	576	
[0, 31]	8	0 1 0 0 0	64	
4	19	1 0 0 1 1	361	

Result of 20 coins tossed.  
 $(PP = 4)$

Let frame the PP.

$x$  can be precision, recall, F1 score, anything

$\bar{f} = 293 \rightarrow$  avg fitness of popn.  
Keep track of current maximum  $\rightarrow$  here 576

$\rightarrow$  prob of selection on roulette wheel.

<u>Index</u>	<u><math>f_i / \bar{f}</math></u>	<u>Expected Count (<math>f_i / \bar{f}</math>)</u>	<u>Round.</u>
1	• 14	• 58	1
2	• 49	• 1.94	2
3	• 06	• 0.22	0
4	• 31	• 1.23	1

Flipping the wheel is selection.

If Round = 1 then we got it once.

2 then we got it 2ce. (2 copies)

Nett generation is the output of selection.

Population size is the same.

0 1 1 0 1

1 1 0 0 0

1 1 0 0 0

1 0 0 1 1

The essence of selection is  
duplicate fitter ones, throw  
away the unfit ones.

1 (single point) crossover. Identify the crossover site.

Crossover site →

011011	Completely random - crossover is done.
110010	Identify a crossover site at
111000	random
101011	01100 11001 11011 10000

1-1 crossover site possible.

Portion after crossover site gets swapped.

Jumped from	Index.	Chromosome.	$f(x)$ .	$\frac{f(x)}{\sum f_i}$	$\frac{f(x)}{F}$	Round
1	12	01100	144.			
24 to	2	25	11001	625		
27	3	(27)	11011	729.		
	4	16	10000	256.		
				1754.		

$$\bar{f} = 438.5$$

You will get the fitter solutions no matter where you cross and whenever you cross.

31/1/2019

Big data is not about the volume. <sup>only</sup>

Collection of databases or data warehouses. How is the data organised? → Big data.

Structured data contains a schema.

Unstructured, Semi-structured, Quasi-structured, location. <sup>↓</sup> weblinks.

Big data deals with all the above types of data.

Because of IoT, the amount data generated is increasing.

We have ~~gigabytes~~ of data.

## The V's.

Volume → how huge the data is

Velocity → streaming data (weather, tweets)

Variety → structured, unstructured, etc

Business insight → business.

Value → historic data's importance.

Variability → error in data, how can it be avoided

## \* Challenges of Big data.

Capturing

Curation.

Storage.

!

Hadoop is a framework.

Scale up, and scale down (out)

SPoF  $\leftarrow$  single point of failure.

how many machines can be added.

Architectures:- MapReduce, HDFS.

4/2/2019

## Recap

- What is the main aim of Roulette wheel?
- Disadvantages of Roulette wheel:-
  - Monopoly: fitter solutions occupy major part of the wheel.
  - We may sometimes lose the optimum solution.

$f$  of the population should be increasing to get the single solution.

GA leads to the path of convergence on every iteration.

Based on increase in avg, we will eventually come up with better results. - Mathematical function application.

Schemata is a pattern of chromosomes.

1 0 0 0

1 0 0 1

1 1 0 0

The pattern in the above data is MSBs are 1.

Pattern is 1\*\*\*

where \* can be either 0 or 1.

Need for schemata - we can reach the optimum solution faster.

Choosing the pattern after two or three result we need to choose the best pattern.

single pattern.

Fundamental Theorem of Schemata:-

Choosing the better pattern may result to the optimum solution. - Mathematical proof.

fitness of schema is average of population.

We are concerned about the increase in fitness from previous generation. (Betterment).

Schema  $\rightarrow$  H

Length of the schema  $\delta(H)$  :- Difference b/w 1st fixed and the last fixed bits of schema.

### Example:-

$$H_1 = * \overset{2}{\underset{1}{|}} * * * * 0$$

$$\delta(H_1) = 5.$$

$$H_2 = \text{*** } 10 \text{ ***} \\ \delta(H_2) = 1$$

Selection :- Selecting the pattern based on schema pattern.

Crossover is a better solution than Mutation

GA is more probabilistic, because of randomness

Schemas are fundamental theorem of f.A.

- Short defined length  $\delta(t)$
  - low  $O(t)$ .
  - Above Average.

Selection Schema Theorem :- survival.

$$m(H, t+1) \geq m(H, t) \frac{f(H)}{T} \left\{ 1 - \frac{p_c f(H)}{L-1} - p_m \theta(H) \right\}$$

↓                      ↓                      ↓                      ↓  
 Count function.      fitness.      popn.      Destruction  
 ↓                      ↓                      ↓                      ↓  
 of schema i.e., Avg.      of schema      mutation

Schema :-

$$H_1 = 1 \text{ 光 沙 沙 沙}$$

$$H_2 = *10*$$

$$f_2 = 1 \text{ 公里} \text{ } 0$$

## Initial population :-

13 : 01101 → String 1

24 : 11000 → String 2

8 : 01000 → String 3

19 : 10011 → String 4

$H_1$  schema is the string representation of 2, 4.

$H_2$  schema is the string representation of 2, 3.

$H_3$  schema is the string representation of 2.

Why are we keen on string representations.

$$M(H_1, t+1) = \frac{2 \times 468.5}{293}$$

$$= 3.19 \cdot \{1 - 0\}$$

$$= 3.19$$

for  $H_2$

$$M(H_2, t+1) = \frac{2 \times 320}{293} = 2.18$$

i. There are 2 copies of  $H_2$  in popn.

$$2.18 \cdot \left\{1 - \frac{1}{4}\right\}$$

$$= 1.63$$

## Benefit of schema based GA:-

- It is the parallel processing capability of GA.
- Scheme based processing is parallel.
- ∴ processing is fast.

## Main objective of Schema theorem:-

In selection we are concerned with betterment

fixed positions survival is exponentially increasing.  
are near:

fixed positions exponential decay or distinctness  
are apart:

11/2/2019

GA classifier : entire concept is ~~based~~ based Bucket brigade

Efficient path which leads to the off faster is  
done by GA

Rule + Msg system - in the domain of expert system  
Bidding results in

① Penalty

② Reward

## Credit Apportionment System

good classifier = +ve credit

bad classifier = -ve credit

What is encoding in GA context?

Soln:- It is used to extract patterns like MCBS etc.

Using Decimal numbers extracting pattern becomes difficult.

In travelling salesman problem.

The feasible solution is cost of visiting each node.

Each solution is called chromosome in GA.

1 2 3 4 5 6 ) swap.  
6 3 4 5 1 2.

Not feasible

Fitness of unipolar mean?

Strength of unipolar which gives better fitness.

Five alarm example

In GA  $O(n^3)$  schema is effectively useful where  $n$  is the population.

GA is inherently parallel.

Among the different paths, the poor paths are not neglected straight away but they are exponentially decreasing in fit.

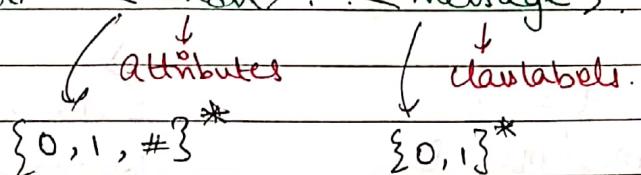
When  $m$  classlabels are there, it is less time consuming than binary class label because getting patterns becomes difficult in binary tree.

18/2/2019

3 aspects of Bucket Brigade :-

- (1) Rule - Meg system.
- (2) Apportionment of credit.
- (3) Genetic Algo.

Rule :- <Classifier> = <condn> : : <message>



Fitness value = strength.

Increasing or decreasing the strength of classifier  
in apportionment of credit.

Table to simulate credit sharing.

	Index	Classifier	Env	Meg	Mtch	Bid	S	Meg	Mtch	Bid	S	Meg	Mtch	Bid	S	Meg	Mtch	Bid
	In	Classifier	t = 0				t = 1				t = 2				t = 3			
1	01##	: 0 0 0 0	200				E	20	180	0 0 0 0					220			220
2	00#0	: 1 1 0 0	200					200			1	20	180	1 1 0 0				218
3	11##	: 1 0 0 0	200					200				200			2	20	180	1 0 0 0
4	# # 00	: 0 0 0 1	200					200			1	20	180	0 0 0 1	2	18	162	0 0 0 1
															3	16.2		

Env = 20.

Bid = 10%.

$\Sigma = 0111 \rightarrow$  environment variable

First strength value is random.

Action or Payoff scheme  $\rightarrow$  submit bit.

One classifier reduces its strength so that the others can get triggered in the subsequent generations.

The Bid is added onto the other classifiers.  
The same path should not keep getting the priority in each generation.

GA kicks in after one run.

New classifier paths are not explored by bucket brigade.

S	t = 4.
220	Match Bid.
218	
196	
146	0001 terminal.

Bucket Brigade is not a contribution of GA

Crossover b/w ①  $\Rightarrow$  ④.

$$\begin{array}{l} 01 \# \# \\ \# \# ; 0 \quad 0 \end{array}$$

$$= 01 \quad 0 \quad 0 \\ \# \# \quad \# \# .$$

21/2/2019

## TUTORIAL ①

- ① A decision tree is being error pruned by collapsing of nodes. For the particular node, on the left branch there are 3 training points [5, 7, 9.6] and for the right branch there are 4 training points [8, 7, 9.8, 10.5, 11].
- (i) What were the original responses along left, right branches?

(i) what is the new response after collapsing?

(2) Say True / False - Reason also.

(a) u. can list all fi's & their support if u know the MFS - B(F)

(b) to perform hierarchical clustering., we do not need to know the coordinates of elements, only their mutual distances are enough (F)

(c) If one used cosine measure to measure interestingness of A + B. its value doesn't change if we add a new transaction. that doesn't contain A or B.

T

$P(X \cdot Y)$

(3) Given T[A .... J]

$\checkmark P(X) \cdot P(Y)$

$\begin{smallmatrix} A & B & C & D & E & F & G & H & I & J \\ 1 & 1 & 1 & 0 & 0 & 0 & 0 & 1 & 1 & 1 \end{smallmatrix}$

$2 \rightarrow 1010110111$

$\frac{1}{\sqrt{X \cdot Y}}$

$3 \rightarrow 1001001111$

$4 \rightarrow 0110101100$

$\sqrt{|X| \cdot |Y|}$

$5 \rightarrow 0111101000$

$6 \rightarrow 1100111010$

$7 \rightarrow 1100000010$

$8 \rightarrow 1110110010$

$9 \rightarrow 1100000011$

$10 \rightarrow 0011110000$

(i) Find all fi's with support  $\geq 5$ .  $\{A, B, C, F, I\}$ .

(ii) which in (i) are closed?

(4) Given T(A,B,C,D) =  $\{1110, 1110, 0111, 0110, 0111, 1100, 1100, 1110, 1110, 0011\}$  (10 transactions).

(a) What is the confidence of.  $(A, B) \rightarrow \emptyset$

(b) Given itemset  $X = \{A, B, C\}$ , List all rules of form.

$Y \rightarrow Z ; Y \neq \emptyset ; Y \subseteq X, Z = X \setminus Y$ . with support  $\geq 4$  and confidence. = 2/3.

(5) Classes .

Buys - Comp = Yes.

Buys - lamp = NO.

Buys - comp = Yes.

6 9 5 4

4 6

Buys - comp = NO

4 1 2

2 5 8 8

Given above, calculate. Precision, Recall, sensitivity, specificity, Accuracy.

⑥ Trace A-Close, pince Search for the data labelled X.

$$X = T_1 \quad \{M, O, N, K, E, Y\}$$

$$T_2 \quad \{D, O, N, K, F, Y\}$$

$$T_3 \quad \{M, A, K, E\}$$

$$T_4 \quad \{M, D, C, Y\}$$

$$T_5 \quad \{C, O, O, K, E, F\}$$

Support = 2

4/3/2019

Neural network is a composition of layers.

Input variables are nodes in network topology.

Back propagation of errors.

Neural network is supervised learning. It has training and testing.

What is error in neural network?

How much prediction is deviated from expectation.

Inputs are fed forward. If the error is greater than error threshold, they are fed back.

"Adjust weights".

Neural network is non linear.

One neuron = node in network.

Firing of neuron

If we are not happy with the o/p, then add hidden layers.

Activation function: the o/p from hidden layer.

$$\sum w_1 + \sum w_2 \text{ if } 0 < = 0.5$$

OR Gate.

$\leq 0.5 \Rightarrow 0$
$> 0.5 \Rightarrow 1$

$I_1$	$I_2$	Output	$x$	0	or output
0	0	0	0	$0.5 \frac{1}{1+e^{-x}}$ → 0	
0	1	1	1	0.7 → 1	
1	0	1	1	0.7 → 1	
1	1	1	1	0.88 → 1	

OR even  
non linear  
combination  
of inputs due  
to activation  
function.

$$I_1, w_1 \rightarrow \textcircled{1} \rightarrow \textcircled{2}$$

$$I_2, w_2 \rightarrow \textcircled{2}$$

Depending on the  
weights, the inputs  
are taken

#### ACTIVATION FUNCTION

$$\text{Sigmoid}(t) = \frac{1}{1 + e^{-x}}$$

$x \rightarrow$  error threshold.

The quality of NN depends  
also on the activation  
function.

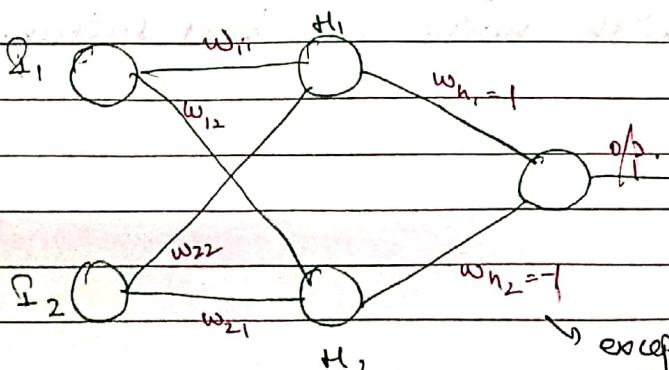
In this network activation function is applied  
in the output layer.

Where does exclusiveness come in XOR?

XOR needs atleast one hidden layer.

11/3/2014

for XOR ① Hidden Layer ② Output neuron.



except for this  
set of the weights  
are 1

overall input to the hidden layer  
 $H_1$

$$\text{Sig } H_1(x) = \frac{1}{1 + e^{-(x - 0.5)}}$$

$$\text{Sig } H_2(x) = \frac{1}{1 + e^{-(x - 1.5)}}$$

$$\text{Sig } D/P(x) = \frac{1}{1 + e^{-(x - 0.2)}}$$

$I_1$	$I_2$	$x$	$H_1$	$H_2$	$D/P$	out.
0	0	0	0.3775	0.1824	0.4998	0
0	1	1	0.6225	0.3735	0.5712	1
1	0	1	0.6225	0.3735	0.5712	1
1	1	2	0.8176	0.6225	0.4998	0

put these in  
(0.1951, 0.2450, 0.19)

$H_1 + H_2$

< 0.5 then 0  
> 0.5 then 1

You can't have a neural network learning correctly without a hidden layer.

Neural network  $\rightarrow$  supervised.

- Disadvantage:-

For live inputs there are too many wt adjustments. That means long training time. Interpretability is less than the previous algos.

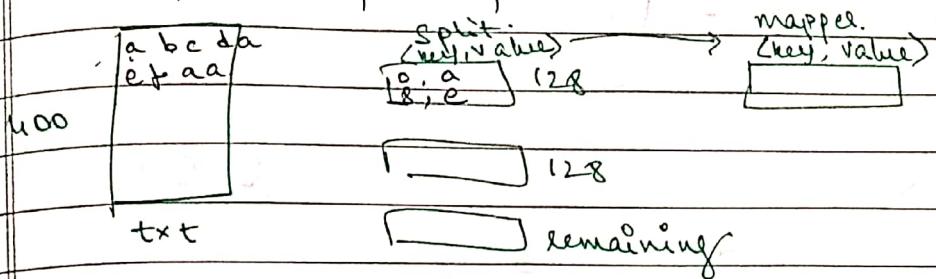
- Advantages

Non linearity of model.

Face recognition, etc where human like learning is required. in terms of precision, recall, etc.

Domains with noise are best suitable for NN.

Aim of Hadoop is parallelism.



Before going to mapper, it will go to split.  
Done in terms of lines.

$\langle \text{key}, \text{value} \rangle$   
 ↓      ↓  
 pos    string.

Each mapper works on a single block.

Mapper.

$\langle a, 1 \rangle$

$\langle b, 1 \rangle$

$\langle c, 1 \rangle$

$\langle d, 1 \rangle$

$\langle a, 1 \rangle$

18/03/2018]

Back propagation possible only if error is there.

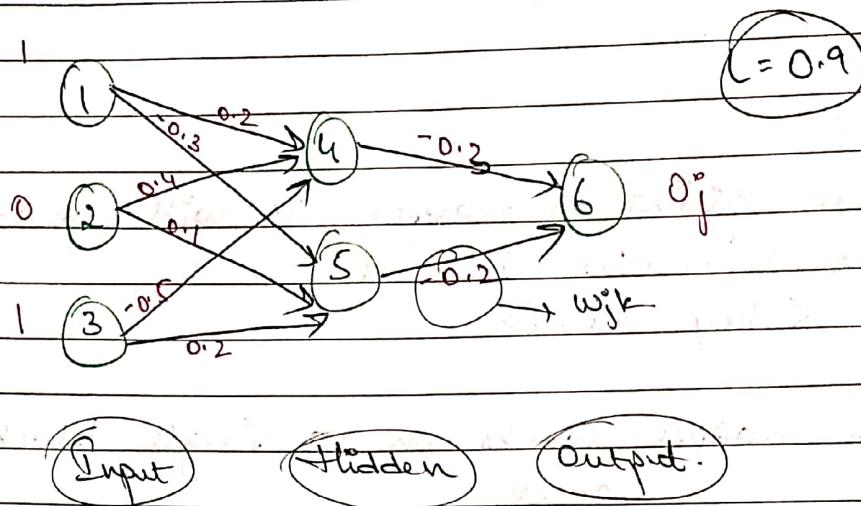
Target off  $\neq$  Predicted off.

$$\text{Error} = T_j - O_j, > \begin{array}{l} \text{Accepted.} \\ \text{Error Threshold.} \\ \text{possible with.} \\ \text{then back.prop.} \\ \text{training Data.} \end{array}$$

Increase or decrease weights depends on the updated error.

$\langle a_0, a_1, a_2 \rangle \rightarrow$  binary (1/0) in nature

$$\boxed{\langle 1 \ 0 \ 1 \rangle - 1}, T_j$$



Network topology  $\equiv$  design of network

BPN was deemed fit for non linearly separable data.

SVM is for linearly separable.

There is no loop in the network. So, back propagation happens as weight update.

$\theta$  - Bias  $\rightarrow$  (0<sub>p</sub> nodes)

$\lambda$  - learning rate

(how quickly  $O_j$  converges to  $T_j$ )

Here  $\lambda = 1$

(t)  $\rightarrow$  at any given point of time, the no. of iterations done is 't'.

$$O_j = \sum_i w_{ij} O_i + \theta_j$$

i  $\rightarrow$  input

j  $\rightarrow$  output.

$$O_j = \frac{1}{1 + e^{-\theta_j}}$$

$$Err(O_j) = O_j(1 - O_j)(T_j - O_j)$$

comparison with overall output and target output

$$Err(+1) = O_j(1-O_j) \sum_k Err_k w_{jk}$$

this is where  
back propagation  
happens.

$$O_4 = -0.4.$$

$$O_5 = 0.2.$$

$$O_6 = 0.1.$$

also gets updated.

$$0.2 + 0 - 0.5 - 0.4.$$

Unit

$$8/p$$

$$\text{similar} \quad O/p (1/(1-e^{-2}))$$

$$4$$

$$-0.7$$

$$0.332.$$

$$5$$

$$0.1$$

$$0.525.$$

$$6.$$

$$-0.3(0.332).$$

$$0.474$$

$$-0.2(0.525).$$

$$+0.1$$

$$= -0.105$$

If error threshold...

ERROR.

$$6. = 0.474(1-0.474)(1-0.474).$$

$$= 0.1311$$

$$7. S = 0.525(1-0.525) \times (-0.2 \times 0.1311) \\ = -0.0065.$$

Similarly

$$8. 4 = -0.0087.$$

$$9. 0.332(1-0.332) \times (-0.332 \times 0.1311)$$

output of 'i' scaled by error of 'j'.

$$\text{New } w_{ij} = \text{Old } w_{ij} + L \times Err_j O_i$$

$$\text{New } O_j = L \times Err_j$$

DEAP



$$w_{46} = -0.3 + 0.9(0.1311 \times 0.332)$$

$$\theta_6 = 0.1 + 0.9(0.1311)$$

25/3/2019

Distance computation in clustering is good for numerical data.

But what do we do for categorical data?

TUTORIAL (As the population is thin)

26/3/2019

Matrix multiplication in Hadoop.

$$A \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix}$$

$$B \begin{bmatrix} 5 & 6 \\ 7 & 8 \end{bmatrix}$$

Input split

A, 0, 0, 1

A, 0, 0, 2

A, 1, 0, 3

A, 1, 1, 4

B, 0, 0, 5

B, 0, 1, 6

B, 1, 0, 7

B, 1, 1, 8

Mapel ( $\langle i, (m, j, v) \rangle$ )

(PHASE-1)

i	Value
0	(A, 0, 1)
0	(A, 1, 2)
1	(A, 0, 3)
1	(A, 1, 4)

i	Value
0	(B, 0, 5)
0	(B, 1, 6)
1	(B, 0, 7)
1	(B, 1, 8)

B matrix is transposed for some reason in the mapper. (Vasanti knows!) To do the multiplication in Row major fashion

Do not shuffle "i" values.

PHASE - II

(key)

(value)

	Value
0	(A, 0, 1)
0	(A, 1, 2)
1	(A, 0, 3)
1	(A, 1, 4)
0	(B, 0, 5)
0	(B, 0, 7)
1	(B, 1, 6)
1	(B, 1, 8)

key of phase - 2.  
corresponds to the  
output indices.

28/3/2018

New matrices are:-

$$\begin{bmatrix} 1 & 2 & 3 & 4 \\ 5 & 6 & 7 & 8 \end{bmatrix}_{2 \times 4} \quad \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 7 & 8 & 9 \end{bmatrix}_{3 \times 3} \quad \boxed{\quad}_{4 \times 3} \quad \boxed{\quad}_{3 \times 2}$$

Do the input split -  
Then mapper phase - I

0 A 0 1	0 B 0 1
0 A 1 2	0 B 1 2
0 A 2 3	0 B 2 3
0 A 3 4	1 B 0 4
1 A 0 5	1 B 1 5
1 A 1 6	1 B 2 6
1 A 2 7	2 B 0 7
1 A 3 8	2 B 1 9
	2 B 2 9
	3 B 0 1
	3 B 1 2
	3 B 3 1

## Mapper phase 2 :-

from  
phase 1  
only

1	$m^o$
0	A 01
0	A 12
0	A 23
0	A 34
1	A 05
1	A 16
1	A 27
1	A 38
0	B 01
0	B 04
0	B 02
1	B 12
1	B 15
1	B 18
1	B 12
2	B 23
2	B 26
2	B 29
2	B 21

shuffled

copied  
from j

## Reduced

00

(A 01) (A 12) (A 23) (A 34)  
(B 01) (B 04) (B 07) (B 01).

01

(A 01) (A 12) (A 23) (A 34).  
(B 12) (B 15) (B 18) (B 12).

02

(A 01) (A 12) (A 23) (A 34).  
(B 23) (B 26) (B 29) (B 21).

10

(A 34) (A 05) (A 16) (A 27)  
(B 01) (B 04) (B 07) (B 01),

11
 $(A_{34}) (A_{05}) (A_{16}) (A_{27})$ 
 $(B_{12}) (B_{15}) (B_{18}) (B_{12})$ 
12
 $(A_{34}) (A_{05}) (A_{16}) (A_{27})$ 
 $(B_{23}) (B_{26}) (B_{29}) (B_{21})$ 

Now,

00

$$(1 \times 1) + (2 \times 4) + (3 \times 7) + (4 \times 1)$$

$$\Rightarrow 1 + 8 + 21 + 4 = 34$$

for all

the

combinations

1/4/2018

(Manhattan distance too)

 Euclidean distance  $\rightarrow$  Numeric Attributes

(Binary Variables, Category Variables, Text Docs).

Similarity Measures for Text Docs :- Hamming wds

K-means evolved as K-medoids.

Interval Scaled Data:-

One centre differs from the other centre.

in terms of units. (kg, g)

Getting data from multiple sources

kg changing to mg (or) mg to kg ?

Impact magnitude of data.

You can't live with existing similarity measures unless there is some amount of standardization.

In J&amp;F it's Normalization.

### ① Mean Absolute Deviation :-

$x_1, \dots, x_n$

$$m_f = \frac{x_{1f} + \dots + x_{nf}}{n}$$

frequency.

$$\hat{x}_{if} = \frac{x_{if} - m_f}{s_f}$$

normalised.

$$S_f = \frac{|x_{1f} - m_f| + \dots + |x_{nf} - m_f|}{n}$$

Mean Abs Dev deals with outliers very well

Dissimilarity  $[0, 1]$

↓  
max.  
div.

Some clustering algs work with  $[-1, 1]$   
Negative correlation.

Targeted Advertisements in mail in clustering  
grouping of related products together.

Near Duplicate detection:- Old  $\not\approx$  new webpage  
Identify duplicates  $\not\approx$  eliminate near duplicates  
Don't get the two links to the same document

Image Segmentation today is a clustering problem

Name a famous partitioning algo:-

Cosine Similarity Measure :- If docs are orthogonal then least similar.

2/4/2019

Jaccard's Similarity.	}	Categorical.
Hamming		
Levenshtein	}	Strings.
Winkler		
Tan's similarity.	}	Euclidean.
Minkowski		Num.
Manhattan		

+ How does clustering help in searching?  
Cluster related customers together. Generate interesting rules.

Clustering followed by ARM  
Instead of searching 1,00,000 URLs search 100 URLs.

Document Vectors

$D_1 \rightarrow$  apple released new ipod.  
 $D_2 \rightarrow$  apple released new ipod.  
 $D_3 \rightarrow$  New apple pie recipe.  
 $D_4 \rightarrow$  phi releases new book with apple pie recipe

Distance  $\rightarrow$  scalar.  
What order are the words in, etc ; - vector.

$$\text{Jaccard's similarity} = \frac{|D_1 \cap D_2|}{|D_1 \cup D_2|}$$

$$\left. \begin{array}{l} D_2, D_3 \\ D_1, D_3 \end{array} \right\}_2$$

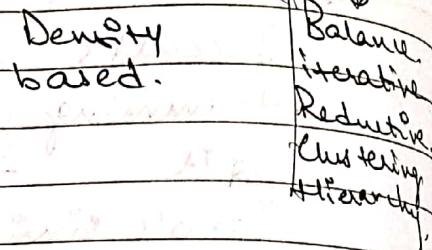
Quantum of deviation.

3 out of 5 words

$$J(D_1, D_2) = 3/5 \text{ are common}$$

Similar not just in terms of words but also ordering.

Other clustering option are DBSCAN & BIRCH.



Choosing the right similarity measure, will lead to better clustering

~~\*\*\* TREC conference~~

V2DB

	Apple	MS.	Trump	Police
Term document frequency matrix:	D <sub>1</sub>	10.	20	0
	D <sub>2</sub>	30.	60.	0
	D <sub>3</sub>	0	0	10

Apple occurs in D<sub>1</sub> 10 times.

Stop & Stem words help in identifying real duplicates. Get unique solutions.

→ Singular Value decomposition (SVD).

→ What is LSI? Latent Semantic Indexing

"T-Tdf is transpose of Tdf"

Taxerds needs what are the 10 words

$D_1, D_2$  look ~~more~~ closer to each other and are far from  $D_3$ .

$$D_3 \perp (D_1, D_2)$$

$D_1, D_2$  growing in the same direction,  
 $D_3$  is growing orthogonally.

\* Similarity score is between 0 & 1.

\* H-distance :-

no. of bits of difference.

$$d_1 = \text{Peter} \quad \begin{matrix} \downarrow \\ \uparrow \downarrow \downarrow \end{matrix} \quad d_2 = \text{Pedro}, \quad \begin{matrix} \downarrow \\ \uparrow \downarrow \downarrow \downarrow \end{matrix}$$

How many bits are different in the two.

$$Hd(d_1, d_2) = 1+1+1 = 3.$$

$$P_1 = 10101.$$

$$P_2 = 10011$$

$$Hd(P_1, P_2) = 1+1 = 2$$

Do  
XOR

No of one's  
not equal gotten in the  
H-d.

\* Edit distance :-

Another view of H-distance.

Paul. Peter  $\rightarrow$  Tardis is better here.  
 Ab Sarah  $\circlearrowleft$  Paul. Pedro  $\circlearrowright$

Weird  
Wierd. ] typo X  
error.

Edit - Insert Delete Modify

What are the no of edits required to transpose / transform one string onto the other

Peter

Peteo.

Peter.

Ped ro.

Paul.

Pial.

req only  
1 transpose  
opn

No new replacement  
Only transposition  
(swapping)  
Less costly operation

Cost here is how many edits to transform along with swapping

Swap is only for a pair

Choose the one that minimizes the edit distance.

Some characters have more priority.

EXAMPLE:-

D<sub>1</sub>. Jones  
D<sub>2</sub>. Johnson

D<sub>1</sub>. to D<sub>2</sub>. → 3 delete & 5 add.

This is feasible, but is it optimal.