

# Web Systems and Algorithms:

## Unit 1:

- introduction
- basic element (content, data structure, algo)
- what application can benefit from intelligence
- eight fallacies (data is reliable, size of data doesn't matter, inference happens instantaneously, apply the same algo. everywhere.)

## Unit 2:

- searching
- indexing
- spiders
- link analysis
- click analysis
- ranking
- precision, recall

## Unit 3:

- recommendation engine
- collaborative (user, item)
- content based

## Unit 4:

- clustering (grouping the things together)
- K-means / K-medoids
- DBSCAN
- Hierarchical

8

Unit 5:

- classification (placing thing where they belong)
- KNN
- Decision Tree
- Bayesian
- Back propagation

Unit 6:

- combining the classifiers
  - don't tie your ship in a single rope
  - don't tie your life in a single hope
  - comparing multiple classifier on same data.
- Bagging and Boosting

Unit 7:

- semantic web
- motivation
- automation
- personalization
- IR
- Data Reuse (Linking of Data)
- RDF (Reuse, Distribution Framework)
- Web Ontology

## 2 Intelligent Web:

- implementing the art of intelligence in web application.

Eg: food ordering system

- machine ask → "What do you like to order today?"
- if every Wednesday you are ordering fish then in Wednesday;  
machine ask → "Would you like fish today?"
- human have learning capability then why don't machine.
- if the machine can't think, we can make them learn.

Eg: Social network offering fact checking, here machine can ensure whether your statements are consistent with your previous message.

Eg: Google Search Engine

- Here the machine take the advantage of interconnected documents to rank them.

Eg: Netflix (machine knows about the movie you enjoy)

- Amazon knows about your buying behaviour.

## 1 Basic elements of Intelligent Applications:

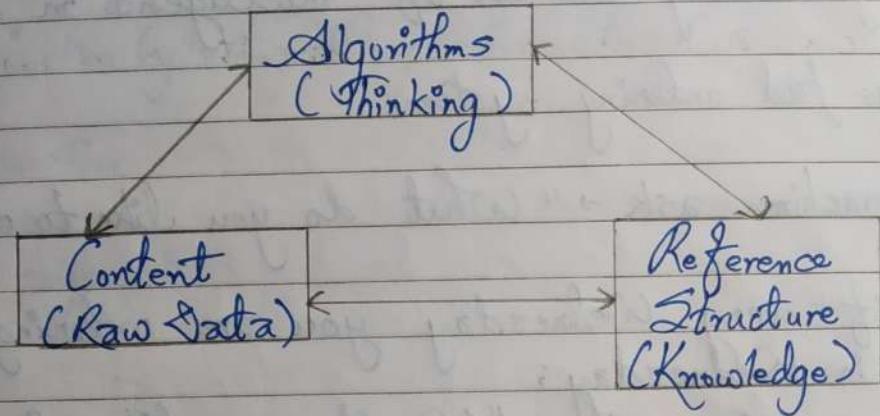


fig: Triangle of Intelligence.

### 1 Aggregated Content:

- dynamic rather than static
- large amount of data
- storage locations and origin could be geographically dispersed.

### 1 Reference Structure:

- provides semantic interpretation of the content.
- dictionaries, knowledge base, ontologies

### 1 Algorithm:

- modules that allows the application to extract the information.

Semantic  
Web

PAGE : 5  
DATE : 11/11/19  
5  
How to interface  
with users

- generalization, forecasting, prediction, improving interaction with users.

Q What applications can benefit from Intelligence?

- Social networking sites

- Job portal

- Aggregator (news)

- online games

Q How can I build intelligence in my own application?

1. Examine your functionality and your data:

- does it deal with free text → (natural language)

- does it deal with geographic data like map

- is fraud detection an issue to your application

- authentication

- does your application make automated decision making

2. Get more data from the web:

- Crawling

- Parsing

- Web Services API

## 8 Fallacies of Intelligent Application:

### 1. Your data is reliable:

- data might contain outlier/noise.
- missing values
- data may change
- different data types (numeric/non-numeric, discrete/continuous)

### 2. Inference happens instantaneously:

- computing a solution takes time.
- don't assume that an algorithm on all datasets run within the same response time.

### 3. The size of data doesn't matter:

- response time
- accuracy

### 4. Scalability of the solution isn't an issue:

- parallel computing
- Eg: not all clustering algorithm can run in parallel

5. Apply the same good library everywhere:

- when you are holding a hammer, everything looks like a nail.

Eg: Lucene Search Engine

6. The computation time is known:

- typically people expect that when we change the parameters of a problem, the problem can be solved consistently with respect to response time
- but it isn't true for all cases

Eg: distance between two locations.

7. Complicated Models are better

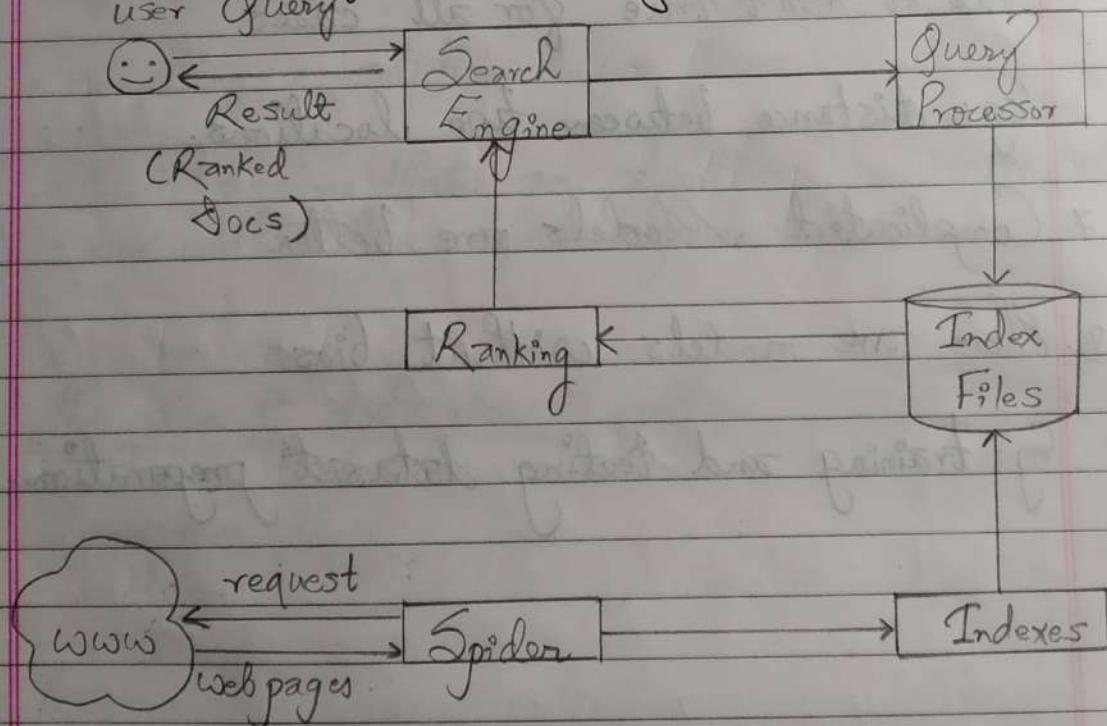
8. There are models without bias:

Eg: training and testing dataset preparation.

## 2 Search Engine:

- a program that searches for document for specified keyword and returns a list of the documents where the keywords are found.
- a search engine works by sending out the spider to fetch the documents as many as
- indexer then reads these documents and create the index.

## 3 Architecture of Search Engine:



## 4 How does Search Engine work?

- Web Crawling
- Indexing
- Searching
- Ranking

## 2. Web Crawling:

- is the process of gathering the pages from web.

## 2 Features of a Crawler:

### (i) Robustness:

- Detecting the trap.

### (ii) Politeness:

- follows the restrictions to spider.

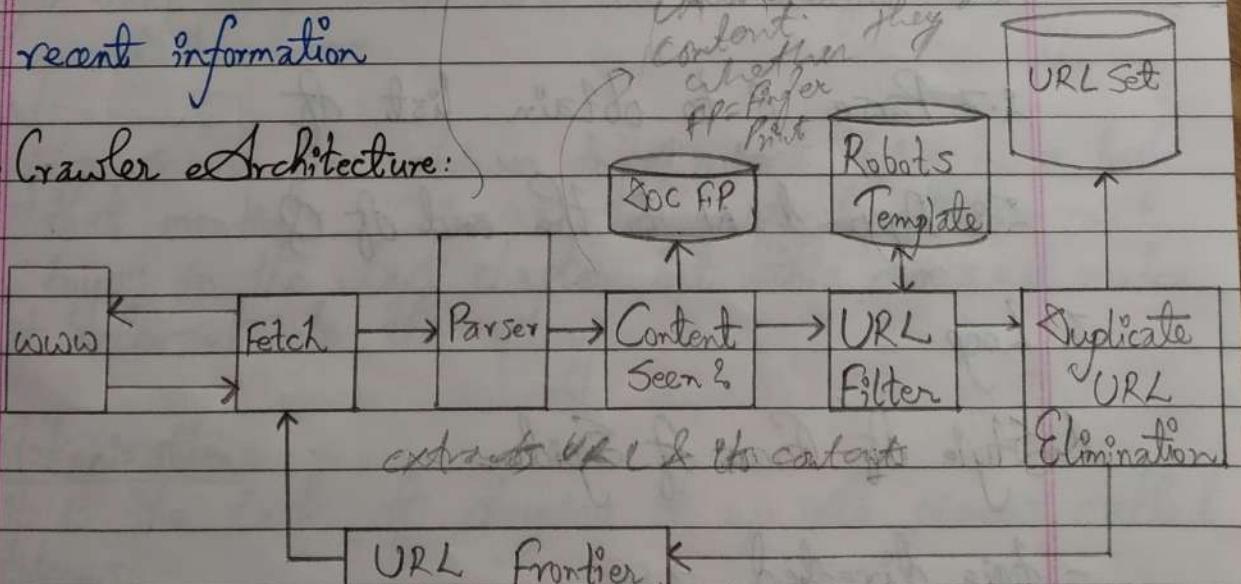
### (iii) Quality:

- Distinguishes fake

### (iv) Freshness:

- recent information

## 2 Crawler Architecture:



*if it gets to certain link it is URL Seed.*

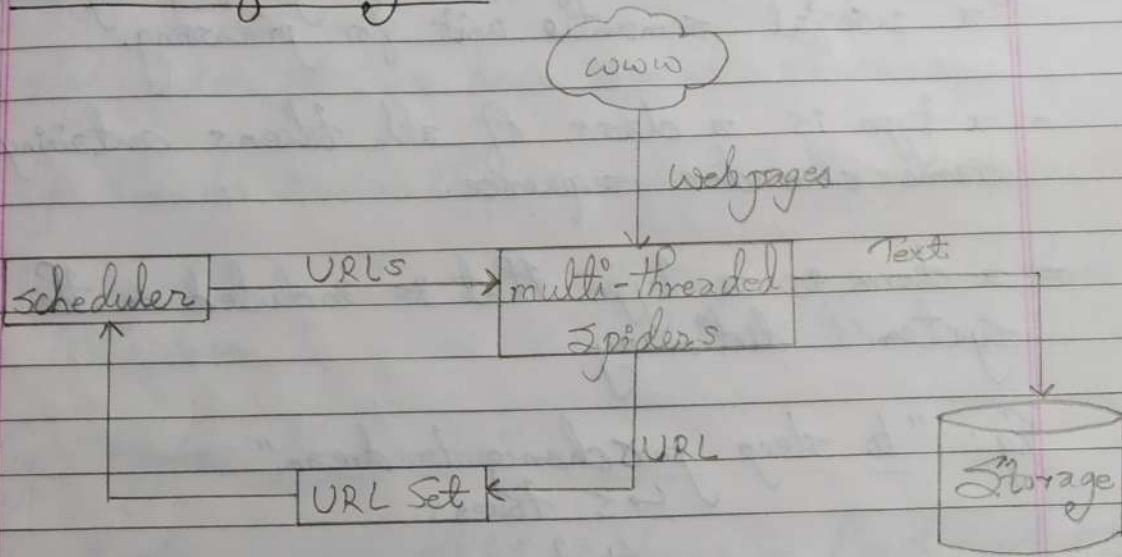
## 2. Crawling Algorithm / Spidering Algorithm:

1. Initialize queue ( $Q$ ) with initial set of known URLs.
2. Until  $Q$  is empty OR Page Limit exceeded OR time limit exceeded
  - 2.1 Pop URL,  $L$ , from  $Q$ .
  - 2.2 If  $L$  is not to an HTML page (.gif, .pdf) then continue loop.
  - 2.3 If  $L$  is already visited then continue loop
  - 2.4 Download page  $P$  for  $L$
  - 2.5 If cannot download  $P$  (e.g. 404 error, robots excluded), then continue loop.
  - 2.6 Index  $P$
  - 2.7 Parse  $P$  to obtain list of  $N$
  - 2.8 Append  $N$  to the end of  $Q$
3. Loop

## 2. Style of Focus of Spider

- topic directed  $\rightarrow$  OS
- link directed 

## 2 Structure of a Spider:



## 3 Indexing:

- choosing a document unit
- determine what the document unit for indexing is?
- for very long documents, the issue of indexing granularity arises.
- eg a search for "chinese toys" might bring up a book that mentions "China" in the first chapter and "toys" in the last chapter but this does not make it relevant to the query.

## 4 Tokenization:

- it is the task of chopping it up into pieces called tokens.
- a token is an instance of characters in some

(12)

particular document that are grouped together as a useful semantic unit for processing.

- a type is a class of all tokens containing the same character sequence.
- a term is a type that is included in the IR system's dictionary.

Eg: "to sleep perchance to dream"  
    ↳ 5 tokens

    4 types  
    3 terms ("to" as stop word)

- issues of tokenization are language specific.

Eg: Chinese language has no space as word separation.

- email id
- URLs
- IP Address
- hyphen -

- but splitting on whitespaces can cause bad retrieval result sometimes.

Eg: Mr. Tribhuvan Shrestha never goes to University.

## 2 Dropping Common Terms:

- stop words are those having higher frequency.
- Eg: the, for, is, am, that, to, with, will etc.
- but a lot of time not indexing' stop word does little harm

Eg: "Flights to London"

- "As we may think"  $\Rightarrow$  Vannevar Bush's Article.
- "Let it be"  $\Rightarrow$  Song of Beatles
- So instead of dropping the stop words, system has to cope with them

## 3 Normalization:

- There are cases, where tokens are not quite the same, but we still want to match them.

Eg: U.S.A, US, USA

window, windows, Window, WINDOW

- - but sometimes it does more harm than good.

Eg: WHO  $\rightarrow$  who

(14)

~~saw's~~  
no meaning  
STAG

~~saw → see, gives meaning~~  
~~atm → atm → atm → atm~~

## 2 Stemming and Lemmatization:

- reduce the size of vocabulary
- stemming is defined as heuristic process that chops off the words.
- Lemmatization refer to doing things properly with the use of vocabulary and morphological analysis of word.

Eg: am, are, is → be

car, cars, car's, cars' → car

## 3 Porter Algorithm:

- most common algorithm for stemming English terms with a common stem will usually have similar meanings.

Eg: connect, connected, connection, connecting.

- a consonant in a word is a letter other than A,E,I,O,U and other than Y preceded by a consonant.

Eg: "JOY" → the consonants are J and Y.

"SYZYGY". → the consonants are S, Z, G.

- If a letter is not a consonant, it is a vowel.
- Let consonant be denoted by "c" and vowel by "v"
- any word has one of the following form:

c... c

c... v

v... c

v... v

- rules in Porter Stemmer can be written as

 $\langle \text{condition} \rangle \rightarrow S_1 \rightarrow S_2$ 

Eg:	Rules	Example
	1. SSESS $\rightarrow$ SS	caresses $\rightarrow$ caress
	2. IES $\rightarrow$ I	Ponies $\rightarrow$ Pon
	3. SS $\rightarrow$ S	caress $\rightarrow$ caress
	4. S $\rightarrow$	Cars $\rightarrow$ car
	5. ( $m \geq 1$ ) EMENT $\rightarrow$ excluding ement	replacement $\rightarrow$ replac but not cement $\rightarrow$ c $m=2$ (ement) <del>(atent)</del>

- here "m" means the measure of word in the form of "VC".

Eg:  $m=0 \Rightarrow \underline{\text{TR}}, \underline{\text{EE}}, \underline{\text{TREE}}$ , no pair of VC.

$m=1 \Rightarrow \underline{\text{TROUBLE}}, \underline{\text{OATS}}, \underline{\text{TREES}}$

$m=2 \Rightarrow \underline{\text{TROUBLES}}, \underline{\text{PRIVATE}}, \underline{\text{OATEN}}$

✓      ✓      ✓  
2 pairs    2 pairs    2 pairs

## 2 Building Inverted Index:

- also called postings.
- is a data structure storing a mapping from content

Eg: doc1 = new home sales

doc2 = home sales rise in July

doc3 = new home in July

home → doc1, doc2, doc3

July → doc2, doc3

new → doc3, doc3

rise → doc2, doc2

sales → doc1, doc2

1. distinct term

2. remove stop word

3. sort alphabetically

## 2 Biword Index:

- index every consecutive pair of terms in the text as a phrase.

Eg: black panther attack

black panther, panther attack

## 2 Positional Indexes:

- posting lists in a positional index in which each posting is a doc ID and a list of positions

Eg: Cat, 100  
 doc no. of occurrence      list of positions

<1, 6: <7, 8, 33, 72, 86, 231>;

2, 5: <1, 17, 74, 222, 255>;

4, 2: <8, 16>;

(17)

PAGE : / /  
DATE : / /

## 2 Document Processing:

1. Lexical Analysis
2. Elimination of stop words
3. Stemming
4. Selection of Index.
5. Thesaurus

### Lexical Analysis of Text:

- basically space is involved as word separator, but however the following cases also have to consider followings:
  1. Digits:
- Numbers are usually not good terms because without surrounding text, they are inherently vague
- Eg. query → "number of deaths between 1910 and 1984" might retrieve document which might contain the numbers only.
- but numbers like credit card number, special date might be index term.

2. Hyphens:

- state-of-the-art, co-education, ... B-99

3. Punctuation Marks:

- normally punctuation marks are removed entirely in the process of lexical analysis while some punctuation marks are integral part of the word.

Eg: Dr., B.C.

4. Case of Letters:

- Bank and bank, who & who.

Thesaurus

- is a collection of words with its synonyms and related words.

Eg: Doctor → doc, medicine, surgeon, physician, anesthetist, hospital.

5. Query Language:

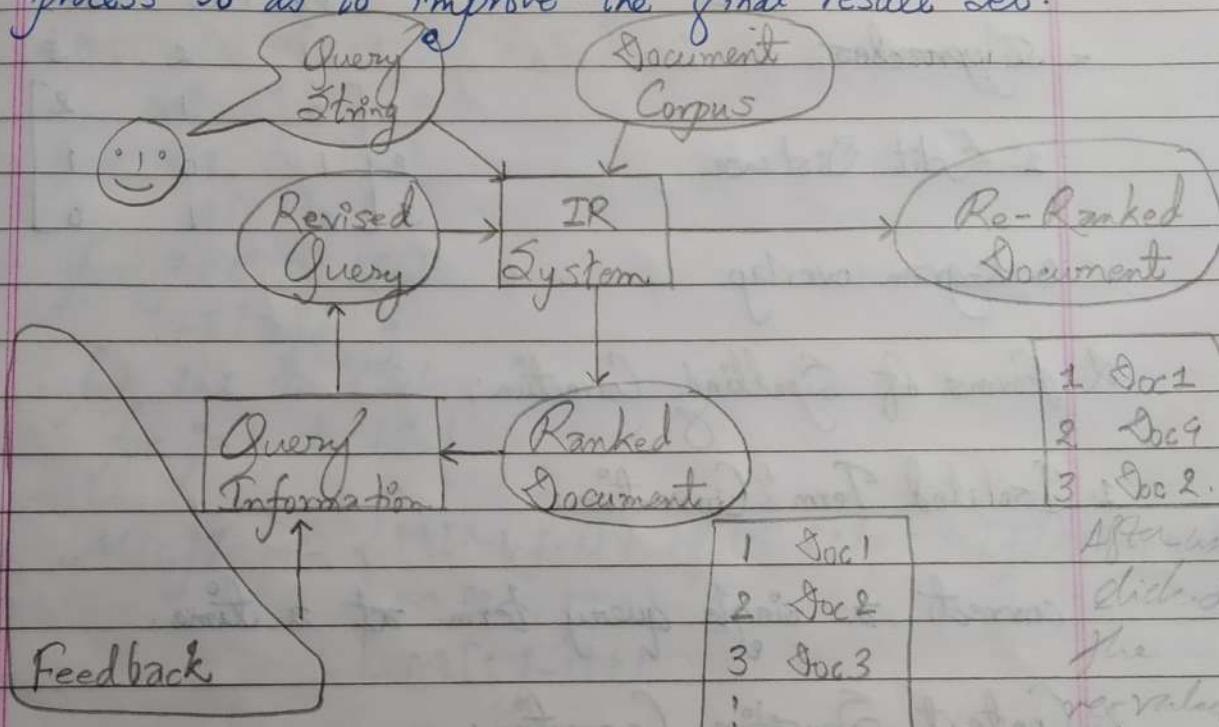
- a query is the formulation of a user information need.

- single word queries.
- boolean queries → cat or dog
- context queries → flights to London.

## (19)

## Relevance Feedback and Query Expansion:

- The idea is to involve the user in the retrieval process so as to improve the final result set.



1 Doc 1
2 Doc 2
3 Doc 3
:

After user  
clicked  
the  
re-ranked  
doc  
looks  
like this!

## 2. Spelling Correction:

- correcting spelling errors in queries       $s_1$        $s_2$   
Eg: cat      car

- Approaches:

1. Edit Distance

	1	2	3
1	0	1	2
2	1	0	1
3	2	1	0

2. k-gram overlap

## 2. Forms of Spelling Correction:

1. Isolated Term Correction:

- correct a single query term at a time.

2. Context Sensitive Correction:

- <sup>near</sup> ~~neig~~ the neighbour words.

Edit Distance:

- given two character strings  $s_1$  and  $s_2$ , the edit distance between them is the minimum number of edit operations required to transform  $s_1$  into  $s_2$ .

- Most commonly edit operations are:

1. Inserting a character

2. Deleting a character

### 3. Replacing a character

Algorithm:

EDIT DISTANCE ( $S_1, S_2$ )

```
int M[i, j] = 0
for i=1 to |S1|
    M[i, 0] = i;
for j=1 to |S2|
    M[0, j] = j
for i=1 to |S1|
    for j=1 to |S2|
```

$$M[i, j] = \min \begin{cases} M[i-1, j-1] + \text{if } (S_1[i] == S_2[j]) \text{ then } 0 \text{ else } 1, \\ M[i-1, j] + 1, \\ M[i, j-1] + 1 \end{cases}$$

K-gram indexes for spelling correction:

- a K-gram is a sequence of k characters

Eg: in term "castle"  $\Rightarrow$  "cas", "ast", "stl", "tle" are 3 gram

- use the k-gram index to retrieve vocabulary terms that have many k-grams in common with the query.

Eg: bo  $\rightarrow$  aboard, about, boardroom, border  
or  $\rightarrow$  border, lord, morbid, sordid  
rd  $\rightarrow$  aboard, ardent, boardroom, border

Matching 2-gram in the query "bord"

## 2. Weighting Mechanism for the term:

### 1. Term Frequency (TF)

- how often a term is found in a document
- denoted by  $tf_{t,d}$

### 2. IDF (Inverse Document Frequency)

- terms which appear very few in numbers
- may have higher probability of being relevant.

$$IDF_t = \log \frac{N}{df_t}$$

$$df \propto \frac{1}{\text{priority}}$$

Eg:

Terms	$df_t$	$IDF_t$
Computer	1059	0.152
Monitor	508	0.470
Keyboard	475	0.500
Device	1247	0.08
Optical	1500	0

Here,  $N$  = total no. of documents = 1500

### 3. TF-IDF

$$TF-IDF_{t,d} = TF_{t,d} \times IDF_t$$

$d_1 \rightarrow \text{cat, cal, dog}$

$TF_{cat,d_1} = 2$

$TF_{dog,d_1} = 1$

$TF_{cal,d_1} = 0$

## 2 Cosine Similarity:

- finding the similarity between the document.

$$\text{cossim}(\vec{A}, \vec{B}) = \frac{\vec{A} \cdot \vec{B}}{|\vec{A}| |\vec{B}|}$$

$$\text{cos} = 1$$

$$\text{cos} = 0$$

Ex:  
 doc1 = you say goodbye  
 doc2 = hello goodbye hello goodbye hello  
 doc3 = I say hello

query = hello goodbye

Find unique words: you, say, goodbye, hello, I

Alphabetical order: goodbye, hello, I, say, you

$$\text{doc1} < 1, 0, 0, 1, 1 >$$

$$\text{doc2} < 2, 3, 0, 0, 0 >$$

$$\text{doc3} < 0, 1, 1, 1, 0 >$$

$$\text{query} < 1, 1, 0, 0, 0 >$$

$$\text{cossim}(\text{query}, \text{doc1}) = \frac{1 \times 1 + 1 \times 0 + 0 \times 0 + 0 \times 1 + 0 \times 1}{\sqrt{1^2 + 1^2} \cdot \sqrt{1^2 + 1^2 + 1^2}} = \frac{1}{\sqrt{6}} = 0.41$$

$$\text{cossim}(\text{query}, \text{doc2}) = \frac{1 \times 2 + 1 \times 3}{\sqrt{2} \cdot \sqrt{2^2 + 3^2}} = \frac{5}{\sqrt{26}} = 0.98$$

$$\text{cossim}(\text{query}, \text{doc3}) = \frac{1 \times 1}{\sqrt{2} \cdot \sqrt{3}} = \frac{1}{\sqrt{6}} = 0.41$$

Ans:  
 ↗ doc2  
 ↗ doc1  
 ↗ doc3

Consider the following documents:

Doc 1 = Information Retrieval System

Doc 2 = Information Storage

Doc 3 = Digital Speech Synthesis System

Doc 4 = Speech Filtering Speech Retrieval

Rank the above document for the query  
"Speech System"

Solution: Unique words: Information, Retrieval, System,  
Storage, Digital, Speech,  
Synthesis, Filtering

Alphabetical order: Digital, Filtering, Information, Retrieval,  
Speech, Storage, Synthesis, System

$$\text{doc 1} = \langle 0, 0, 1, 1, 0, 0, 0, 1 \rangle$$

$$\text{doc 2} = \langle 0, 0, 1, 0, 0, 1, 0, 0 \rangle$$

$$\text{doc 3} = \langle 1, 0, 0, 0, 1, 0, 1, 1 \rangle$$

$$\text{doc 4} = \langle 0, 1, 0, 1, 2, 0, 0, 0 \rangle$$

$$\text{query} = \langle 0, 0, 0, 0, 1, 0, 0, 1 \rangle$$

$$\begin{aligned} \text{cosine}(\text{query}, \text{doc 1}) &= \frac{0 \times 0 + 0 \times 0 + 0 \times 1 + 0 \times 1 + 2 \times 0 + 0 \times 0 + 0 \times 0 + 0 \times 1}{\sqrt{1^2 + 1^2} \cdot \sqrt{1^2 + 1^2 + 1^2}} \\ &= \frac{1}{\sqrt{2} \cdot \sqrt{3}} \\ &= \frac{1}{\sqrt{6}} \end{aligned}$$

$$\text{cosim}(\text{query}, \text{doc2}) = \frac{0x0 + 0x0 + 0x1 + 0x0 + 1x0 + 0x1 + 0x0 + 1x0}{\sqrt{2} \cdot \sqrt{2}}$$

$$= \frac{0}{\sqrt{4}}$$

$$= 0$$

$$\text{cosim}(\text{query}, \text{doc3}) = \frac{0x1 + 0x0 + 0x0 + 0x0 + 1x1 + 0x0 + 0x1 + 1x1}{\sqrt{2} \cdot \sqrt{4}}$$

$$= \frac{2}{\sqrt{8}}$$

$$= \frac{2}{2\sqrt{2}}$$

$$= \frac{1}{\sqrt{2}}$$

$$\text{cosim}(\text{query}, \text{doc4}) = \frac{0x0 + 0x1 + 0x0 + 0x1 + 1x2 + 0x0 + 0x0 + 1x0}{\sqrt{2} \cdot \sqrt{2^2 + 1^2 + 2^2}}$$

$$= \frac{2}{\sqrt{2} \cdot \sqrt{6}} = \frac{2}{\sqrt{12}} = \frac{2}{2\sqrt{3}}$$

$$= \frac{1}{\sqrt{3}}$$

## 1. Link Analysis:

- use of hyperlinks for ranking web search results.
- link analysis is one of many factors considered by web search engines in computing a score for a web page on any given query.
- two methods :

### 1. Page Rank

### 2. HITS (Hyperlink Induced Topic Search)

#### 1. Page Rank:

- developed by Larry Page at SUN
- Link Analysis Algorithm.
- a hyperlink to a page counts as a vote of support
- a page that is linked to by many pages receives a high rank and if there is no links to a web page, there is no support for that page.
- a page rank of 0.5 means there is a 50% chance that a person clicking on a random link will be directed to the document with 0.5 page rank.

## Algorithm:

- Assume a small universe of four web pages A, B, C and D
- the initial approximation of page rank would be evenly divided between four documents i.e.  $\frac{1}{4} = 0.25$
- if pages B, C and D each only link to A, they would confer 0.25 page rank to A.

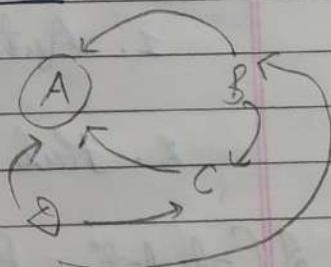
$$\text{i.e. } PR(A) = \frac{PR(B)}{1} + \frac{PR(C)}{1} + \frac{PR(D)}{1} = 0.75$$

- Suppose that page B has link to page C as well as to page A, which D has links to all three pages

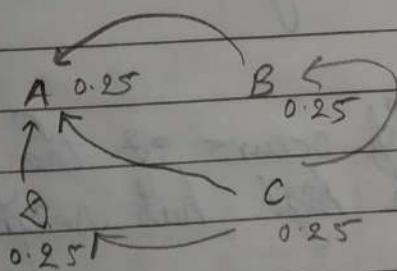
$$PR(A) = \frac{PR(B)}{2} + \frac{PR(C)}{1} + \frac{PR(D)}{3}$$

i.e. in general

$$PR(m) = \sum_{n \in B_m} \frac{PR(n)}{L(n)}$$



where,  $B_m \rightarrow$  set of all pages link to m  
 $L(n) \rightarrow$  outgoing links from n.



## 2 Authorities and Hubs:

- A page is called an authority for the query if it contains the valuable information on the subject
- eg. for a query car  $\Rightarrow$  "www.6mns.com"
- authoritative pages are truly relevant to the query.
- hubs contain useful links towards the authoritative pages.
- i.e. hubs point the search engine to the right decision.
- HITS algorithm is based on two parameters.
  1. Authority Weight
  2. Hub Weight

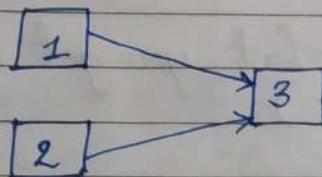
## 3 Calculation Process in HITS algorithm:

- authority weight and hub weight are calculated recursively.
- the higher authority weight occurs if the page is pointed to by pages with high hub weights.
- a higher hub weight occurs if the page points to many pages with high authority weights.

- for a web page  $p$ ,

$$h(p) = \sum_{p \rightarrow y} a(y) \quad | \quad \text{where } m \rightarrow n \text{ denotes the existence of hyperlink}$$
$$a(p) = \sum_{y \rightarrow p} h(y) \quad | \quad \text{from } m \text{ to } n.$$

Eg:



$$- \alpha = A^T h$$

$$- h = A \alpha$$

when  $A$  is adjacency matrix.

$$A_{ij} = \begin{cases} 1, & \text{if there is a hyperlink from } i \text{ to } j \\ 0, & \text{otherwise} \end{cases}$$

- in this case,

$$\begin{matrix} & 1 & 2 & 3 \\ 1 & 0 & 0 & 1 \\ 2 & 0 & 0 & 1 \\ 3 & 0 & 0 & 0 \end{matrix}$$

$$\alpha = A^T h$$

$$= \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 1 & 1 & 0 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} \xrightarrow{\text{initialization of } h \text{ vector}} = \begin{bmatrix} 0 \\ 0 \\ 2 \end{bmatrix}$$

$$h = A \alpha$$

$$= \begin{bmatrix} 0 & 0 & 1 \\ 0 & 0 & 1 \\ 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} 0 \\ 0 \\ 2 \end{bmatrix} = \begin{bmatrix} 2 \\ 2 \\ 0 \end{bmatrix}$$

2. Evaluation Criteria: (Is what you got what you want)

Precision:

- fraction of retrieved documents that are relevant.

$$\text{Precision} = \frac{\text{# (relevant item retrieved)}}{\text{# (retrieved item)}}$$

Recall:

- is the fraction of relevant documents that are retrieved.

$$\text{Recall} = \frac{\text{# (relevant items retrieved)}}{\text{# (relevant items)}}$$

F-Measure:

- harmonic mean of precision and recall.

$$F\text{-Measure} = \frac{2PR}{P+R}$$

Eg: An IR System returns 8 relevant documents and 10 non-relevant documents. There are 20 relevant documents in the collection. Calculate the precision, recall and F-measure.

$$8 \rightarrow R \\ 10 \rightarrow NR$$

$$R = 8$$

$$\text{Precision} = \frac{8}{18}$$

$$\text{Recall} = \frac{8}{20}$$

Recall  $\rightarrow$  100% retrieved.

$$F\text{-measure} = \frac{P+R}{2}$$

100%

250

In case of Ranked Documents:

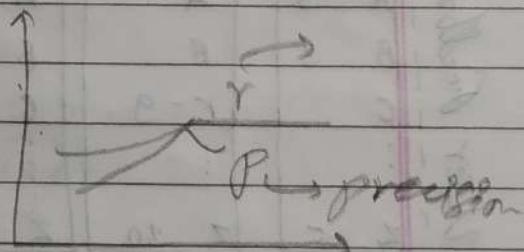
# of relevant documents = 6

Rank	docID	Relevance	Precision	Recall
1	100	✓	1/1	1/6
2	103	✓	2/2	2/6
3	104	✗	2/3	2/6
4	106	✓	3/4	3/6
5	105	✗	3/5	3/6

Assignment:

Click Analysis

- Scope and Applications
- Interest Mining
- Disambiguation
  - bat  $\rightarrow$  mammal
  - $\rightarrow$  cricket bat
- Rating (Relevance Feedback)



1. Recommender System:

- a specific type of information filtering system techniques that attempts to recommend information, items that are likely to be interest to the user.

2. Approaches:

1. Collaborative Filtering
2. Content Based Recommendation

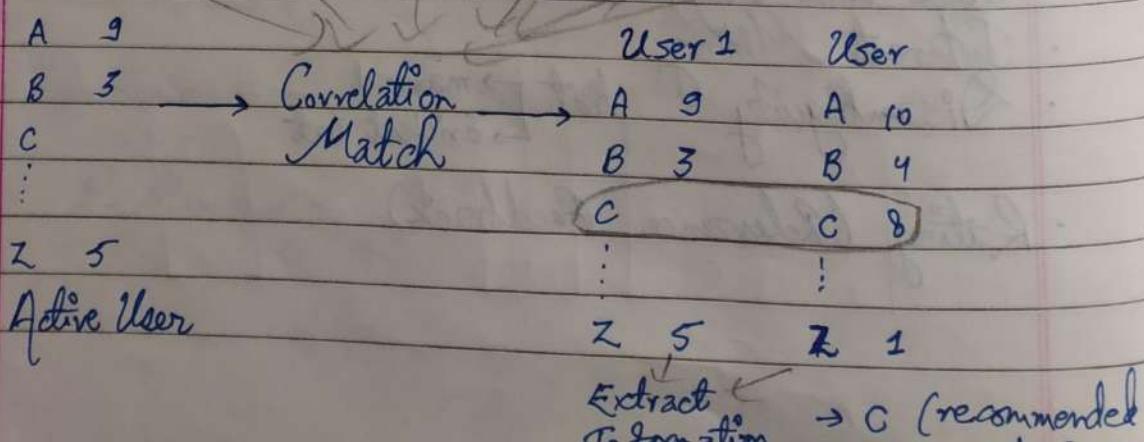
1. Collaborative Filtering:

- Tell me your friends tell I will tell you who you are.

Collaborative Filtering Working Mechanism:

  
 User 1    User 2    User 3    User 4    User 5    User 6

User Feedback	A 9	A 3	A 6	A 10
B 3	B	B	B 4	B 4
C 9	G 9	C	C 8	C 8
:	:	:	:	:
Z 5	Z 10	Z 7	Z	Z 1



- Typically, Pearson correlation coefficient is used between active user and other user.

$$C_{a,u} = \frac{\text{covar}(r_a, r_u)}{\sigma_{r_a} \sigma_{r_u}}$$

where,

$$\text{covar}(r_a, r_u) = \frac{\sum_{i=1}^m (r_{a,i} - \bar{r}_a)(r_{u,i} - \bar{r}_u)}{m}$$

Covariance

$$\bar{r}_{xu} = \frac{\sum_{i=1}^m r_{x,i}}{m}$$

$$\sigma_{r_x} = \sqrt{\frac{\sum_{i=1}^m (r_{x,i} - \bar{r}_x)^2}{m}}$$

$m \rightarrow$  no. of commonly related items.

## 2 Significance Weighting:

$$w_{a,u} = S_{a,u} C_{a,u} \text{ where, } S_{a,u} = \begin{cases} 1 & \text{if } m > 50 \\ 0 & \text{otherwise} \end{cases}$$

## 2 Rating Prediction:

$$P_{a,i} = r_a + \frac{\sum_{u=1}^n w_{a,u} (r_{u,i} - \bar{r}_u)}{\sum_{u=1}^n w_{a,u}}$$

## 2 Limitations:

- Cold Start
- First Rater
- Popularity Bias (unique choice)

34

Q. Consider the following rating matrix:

User \ Item	I <sub>1</sub>	I <sub>2</sub>	I <sub>3</sub>	I <sub>4</sub>
User	4	?	5	5
U <sub>1</sub>	4	2	1	X
U <sub>2</sub>	3	X	2	4
U <sub>3</sub>	4	9	X	X
U <sub>4</sub>	2	1	3	5

Q. Calculate the prediction value of item I<sub>2</sub> for User U<sub>4</sub>.

Here,

$$P_{2,4} = \bar{r}_2 + \frac{\sum_{u \in U} (r_{u,2} - \bar{r}_u) w_{2,u}}{\sum_{u \in U} |w_{2,u}|}$$

where,

$$w_{2,u} = \frac{\sum_{i \in m} (r_{2,i} - \bar{r}_2) (r_{u,i} - \bar{r}_u)}{\sqrt{\sum_{i \in m} (r_{2,i} - \bar{r}_2)^2} \cdot \sqrt{\sum_{i \in m} (r_{u,i} - \bar{r}_u)^2}}$$

$i \in m = \{1, 3\}$   
↓  
for initial case

$m$  = other users who have ranked I<sub>2</sub>.

$m$  = commonly rated by 2 and u.

$w_{2,u}$  = similarity between two users 2 and u.

We have,  $r_1 = (4+5+5)/3 = 4.67$

$$r_2 = (4+2+1)/3 = 2.33$$

$$r_3 = (3+2+4)/3 = 3$$

$$r_4 = (4+4)/2 = 4$$

$$r_5 = (2+1+3+5)/4 = 2.75$$

Calculate correlated items:

users      items

$w_{1,2}$

$i, j$

$w_{1,3}$

$1, 3, 4$  ( $\because$  no need to compute as there is no rating on Item  $I_2$ )

$w_{1,4}$

$i$

$w_{1,5}$

$i, 3, 4$

$i = 1, 3$  (new)

$$w_{1,2} = \frac{(r_{1,1} - \bar{r}_1)(r_{2,1} - \bar{r}_2) + (r_{1,3} - \bar{r}_1)(r_{2,3} - \bar{r}_2)}{\sqrt{(r_{1,1} - \bar{r}_1)^2 + (r_{1,3} - \bar{r}_1)^2} \cdot \sqrt{(r_{2,1} - \bar{r}_2)^2 + (r_{2,3} - \bar{r}_2)^2}}$$

$w_{2,4}$

$$\checkmark$$

Here,

$$\bar{r}_1 = \frac{4+5}{2} = 4.5$$

$$\bar{r}_2 = \frac{5+1}{2} = 2.5$$

$$w_{1,2} = \frac{(4-4.5)(4-2.5) + (5-4.5)(1-2.5)}{\sqrt{(4-4.5)^2 + (5-4.5)^2} \cdot \sqrt{(4-2.5)^2 + (1-2.5)^2}}$$

$$= \frac{-1.5}{0.707 \times 2.1213}$$

$$= -1.0002.$$

$$w_{1,3} = \frac{(r_{1,1} - \bar{r}_1) \cdot (r_{3,1} - \bar{r}_3) + (r_{1,3} - \bar{r}_1) \cdot (r_{3,3} - \bar{r}_3) + (r_{1,4} - \bar{r}_1) \cdot (r_{3,4} - \bar{r}_3)}{\sqrt{(r_{1,1} - \bar{r}_1)^2 + (r_{1,3} - \bar{r}_1)^2 + (r_{1,4} - \bar{r}_1)^2} \cdot \sqrt{(r_{3,1} - \bar{r}_3)^2 + (r_{3,3} - \bar{r}_3)^2 + (r_{3,4} - \bar{r}_3)^2}}$$

$$\checkmark$$

$$\text{Here, } \bar{r}_1 = \frac{4+5+5}{3} = 4.67 \quad \bar{r}_3 = \frac{3+2+4}{3} = 3.$$

$$w_{1,3} = \frac{(4-4.67)(3-3) + (5-4.67)(2-3) + (5-4.67)(4-3)}{\sqrt{(4-4.67)^2 + (5-4.67)^2 + (5-4.67)^2} \cdot \sqrt{(3-3)^2 + (2-3)^2 + (4-3)^2}}$$

$$= 0$$

$$\omega_{1,4} = \frac{(r_{1,1} - \bar{r}_1) \cdot (r_{4,1} - \bar{r}_4)}{\sqrt{(r_{1,1} - \bar{r}_1)^2} \cdot \sqrt{(r_{4,1} - \bar{r}_4)^2}}$$

Here,

$$\bar{r}_1 = 4$$

$$\bar{r}_4 = 9$$

$$\omega_{1,4} = \frac{(4-4) \cdot (9-4)}{\sqrt{(4-4)^2} \cdot \sqrt{(9-4)^2}} \\ = 0$$

$$\omega_{1,5} = \frac{(r_{1,1} - \bar{r}_1) (r_{5,1} - \bar{r}_5) + (r_{1,3} - \bar{r}_1) (r_{5,3} - \bar{r}_5) + (r_{1,4} - \bar{r}_1) (r_{5,4} - \bar{r}_5)}{\sqrt{(r_{1,1} - \bar{r}_1)^2 + (r_{1,3} - \bar{r}_1)^2 + (r_{1,4} - \bar{r}_1)^2} \cdot \sqrt{(r_{5,1} - \bar{r}_5)^2 + (r_{5,3} - \bar{r}_5)^2 + (r_{5,4} - \bar{r}_5)^2}}$$

Here,

$$\bar{r}_1 = \frac{4+5+5}{3} = 4.67$$

$$\bar{r}_5 = \frac{2+3+5}{3} = 3.33$$

$$\omega_{1,5} = \frac{(4-4.67) \cdot (2-3.33) + (5-4.67) \cdot (3-3.33) + (5-4.67) \cdot (5-3.33)}{\sqrt{(4-4.67)^2 + (5-4.67)^2 + (5-4.67)^2} \cdot \sqrt{(2-3.33)^2 + (3-3.33)^2 + (5-3.33)^2}}$$

$$= \frac{1.333}{0.8167 \times 2.1603}$$

$$= 0.7555$$

$$\text{Now, } P_{2,i} = \bar{r}_2 + \frac{\sum_{u \in V} (r_{u,i} - \bar{r}_u) \cdot \omega_{2,u}}{\sum_{u \in V} |\omega_{2,u}|} \quad u \in V = \{2, 3, 4, 5\}$$

$$P_{1,2} = \bar{r}_1 + \frac{(r_{2,2} - \bar{r}_2) \cdot \omega_{1,2} + (r_{3,2} - \bar{r}_3) \cdot \omega_{1,3} + (r_{4,2} - \bar{r}_4) \cdot \omega_{1,4} (r_{5,2} - \bar{r}_5) \cdot \omega_{1,5}}{|\omega_{1,2}| + |\omega_{1,3}| + |\omega_{1,4}| + |\omega_{1,5}|}$$

Here,  $\bar{r}_2 = \frac{4+1}{2} = 2.5$      $\bar{r}_3 = \text{Not required since } V_3 \text{ has not rated } I_2.$

$$\bar{r}_4 = 4$$

$$\bar{r}_1 = 4.67 = \frac{4+5+5}{3}$$

$$\begin{aligned}
 P_{1,2} &= 4.67 + \frac{(2-2.5) \cdot (-1.0002) + (4-4) \cdot 0 + (1-3.33) \cdot 0.7555}{|-1.0002| + |0| + |0.7555|} \\
 &= 4.67 + \left( \frac{-1.2602}{1.7557} \right) \\
 &= 3.9522
 \end{aligned}$$

## 2 Content Based Recommendation System:

- This approach analyze a set of documents previously liked by a user, and build a model of user, and build a model of interest based on the features of the objects rated or liked by the user.
- Worked on three steps:
  1. Content Analyzer
  2. Profile Learner
  3. Filtering Component

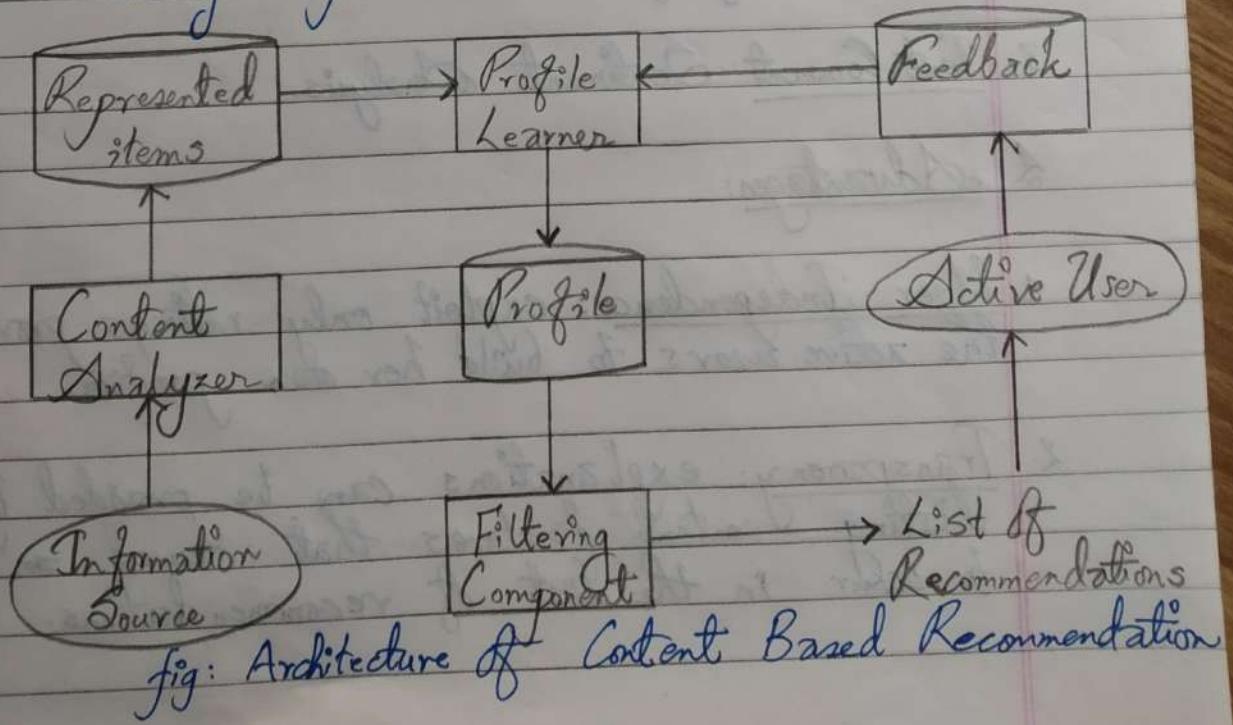


fig: Architecture of Content Based Recommendation

34

Content Analyzer: preprocessing to extract relevant information.

Profile Learner: collects data of user preferences and tries to generalize the data.

- like/dislike past

- +ve/-ve feedback

Filtering Component: suggest relevant items.

2. Approaches to get explicit relevance feedback.

1. Like/Dislike: items are classified as relevant or not relevant by adopting a simple binary rating scale. (1-5)

2. Ratings: a discrete numeric scale is usually adopted to judge items.

3. Text Comment: Sentiment Analysis.

2. Advantages:

1. User Independence: exploit only ratings provided by the active users to build her own profile.

2. Transparency: explanations can be provided by explicitly listing content features that caused an item to occur in the list of recommendations.

3. New Item: capable of recommending items not yet rated by user.

Limitations:

1. Limited Content Analysis:

- Domain knowledge is often needed. e.g. while movie recommendations, the system needs to know about director, actor etc.

2. Over Specialization:

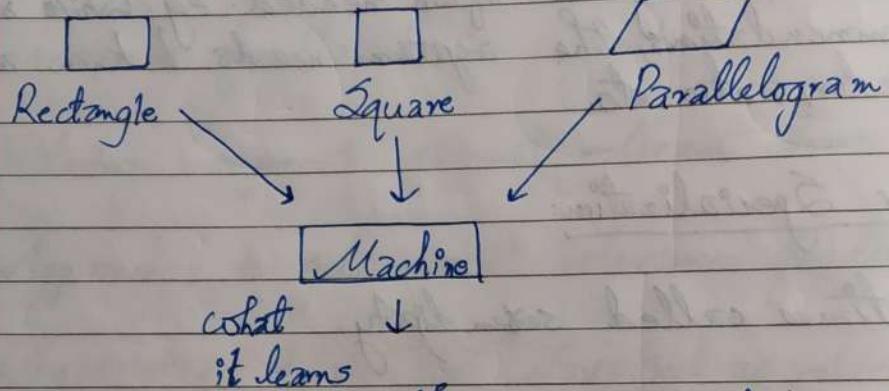
- Sometimes called serendipity.

3. New User:

- enough ratings have to be collected.

Clustering:

- grouping the things together.
- unsupervised learning.

Approach:

They are all bounded with 4 st. lines so they are rectangular.

Algorithm:

- K-means
- K-medoid
- Hierarchical Clustering  $\rightarrow$  Agglomerative / Divisive
- DBSCAN
- ROCK  $\rightarrow$  link based clustering algorithm.

41

4 K-means Algorithm:

Input: Datasize ( $D$ ), no. of clusters ( $k$ ).

Output:  $k$ -clusters

1. Select the  $k$  sample randomly as initial centroid.

2. Repeat

2.1 Find the distance of every data objects with the centroids.

2.2 Assign the object to the closest cluster.

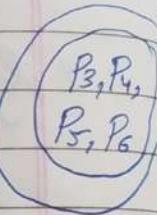
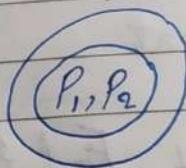
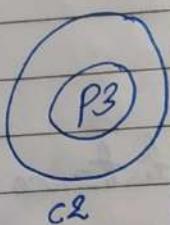
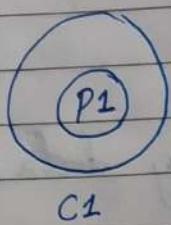
2.3. Update the centroid

3. Until

3.1 No change in centroid

3.2 No change in object of clusters.

Eg:  $P_1(x_1, y_1)$   
 $P_2(x_2, y_2)$   
 $P_3(x_3, y_3)$   
 $P_4(x_4, y_4)$   
 $P_5(x_5, y_5)$   
 $P_6(x_6, y_6)$

 $k=2$ 1st Iteration:

$$P_1 \leftarrow \begin{cases} C_1 & \min_{P_i \in C_1} \\ C_2 & \min_{d_{21}} \end{cases}$$

$$P_2 \leftarrow \begin{cases} C_2 & \min_{d_{12}} \\ C_1 & \min_{d_{22}} \end{cases}$$

 $P_3$ 

:

 $P_4$  $P_5$  $P_6$ 2nd Iteration:

$$P_1 \leftarrow \begin{cases} C_1 & \frac{-P_1.x + P_2.x}{2} \\ C_2 & \frac{P_1.y + P_2.y}{2} \end{cases}$$

:

:

$$P_6 \leftarrow \begin{cases} C_1 \\ C_2 \end{cases}$$

Q. Apply k-means ( $k=2$ ) to make clusters from following documents:

Doc 1 = go Longhorns go

Doc 2 = go Texas

Doc 3 = Texas Longhorns

Doc 4 = Longhorns Longhorns

- Assume Doc 1 and Doc 3 as initial centroid and perform the computation upto 2-iterations.

Solution: Unique words: go, Longhorns, Texas.

Sorting in Ascending Order we get: go, Longhorns, Texas

Now, let's construct the vector as:

$$\text{Doc 1} = \langle 2, 1, 0 \rangle$$

$$\text{Doc 2} = \langle 1, 0, 1 \rangle$$

$$\text{Doc 3} = \langle 0, 1, 1 \rangle$$

$$\text{Doc 4} = \langle 0, 2, 0 \rangle$$

Since Doc 1 and Doc 3 are considered as initial centroid so we classify them to cluster 1 and cluster 2 respectively and calculate the cosine distance to include the documents into the probable clusters.

(44)

We have vector as:

$$\text{doc } 1 = \langle 1, 0, 0 \rangle$$

$$\text{doc } 2 = \langle 1, 0, 1 \rangle$$

$$\text{doc } 3 = \langle 0, 1, 1 \rangle$$

$$\text{doc } 4 = \langle 0, 1, 0 \rangle$$

1st Iteration:

$$\begin{array}{l} \text{doc } 1 \left\{ \begin{array}{l} c_1 - \text{doc } 1 = 1 \\ c_2 - \text{doc } 3 = 0 \end{array} \right. \end{array}$$

$$c_1 = \text{doc } 1 \text{ and } c_2 = \text{doc } 3$$

$$\begin{array}{l} \text{doc } 2 \left\{ \begin{array}{l} c_1 - \text{doc } 1 = 0.75 \\ c_2 - \text{doc } 3 = 0.5 \end{array} \right. \end{array}$$

$$\begin{array}{l} \text{doc } 3 \left\{ \begin{array}{l} c_1 - \text{doc } 1 = 0 \\ c_2 - \text{doc } 3 = 1 \end{array} \right. \end{array}$$

$$\begin{array}{l} \text{doc } 4 \left\{ \begin{array}{l} c_1 - \text{doc } 1 \\ c_2 - \text{doc } 3 \end{array} \right. \end{array}$$

$$\text{Now, } \cosine(\text{doc } 1, c_1) = \frac{1 \times 1 + 0 \times 0 + 0 \times 0}{\sqrt{1^2} \cdot \sqrt{1^2}} \\ = 1$$

$$\cosine(\text{doc } 2, c_1) = \frac{1 \times 1 + 0 \times 0 + 0 \times 0}{\sqrt{1^2 + 0^2 + 1^2} \cdot \sqrt{2^2 + 1^2 + 0^2}}$$

$$(2a) = 0.63$$

$$\cosine(\text{doc } 3, c_1) = \frac{1 \times 0 + 0 \times 1 + 0 \times 0}{\sqrt{0^2 + 1^2 + 1^2} \cdot \sqrt{2^2 + 1^2 + 0^2}}$$

$$(3a) = 0.31$$

$$\cosine(\text{doc}_1, c_1) = \frac{0+2+0}{\sqrt{0^2+2^2+0^2} \cdot \sqrt{2^2+1^2+0^2}}$$

$$= 0.44$$

$$\cosine(\text{doc}_1, c_2) = \frac{0+1+0}{\sqrt{2^2+1^2+0^2} \cdot \sqrt{0^2+1^2+1^2}}$$

$$= 0.31$$

$$\cosine(\text{doc}_2, c_2) = \frac{0+0+1}{\sqrt{1^2+0^2+1^2} \cdot \sqrt{0^2+1^2+1^2}}$$

$$= 0.5$$

$$=$$

$$\cosine(\text{doc}_3, c_2) = 1$$

$$\cosine(\text{doc}_4, c_2) = \frac{0+2+0}{\sqrt{0^2+2^2+0^2} \cdot \sqrt{0^2+1^2+1^2}} = 0.70$$

$$c_1 = \text{doc}_1, \text{doc}_2.$$

$$c_2 = \text{doc}_3, \text{doc}_4.$$

Iteration 2:

$$c_1 = \left( \frac{2+1}{2}, \frac{1+0}{2}, \frac{0+1}{2} \right) = (1.5, 0.5, 0.5)$$

$$c_2 = \left( \frac{0+0}{2}, \frac{1+2}{2}, \frac{1+0}{2} \right) = (0, 1.5, 0.5)$$

$$\cosine(\text{doc}_1, c_1) = \frac{3+0.5+0}{\sqrt{2^2+1^2+0^2} \cdot \sqrt{(1.5)^2+(0.5)^2+(0.5)^2}}$$

$$= 0.94$$

$$\cosine(\text{doc2}, \text{c1}) = \frac{1.5 + 0 + 0.5}{\sqrt{1^2 + 0^2 + 1^2} \cdot \sqrt{(1.5)^2 + (0.5)^2 + (0.5)^2}}$$

$$\begin{aligned}\cosine(\text{doc3}, \text{c1}) &= \frac{0 + 0.5 + 0.5}{\sqrt{0^2 + 1^2 + 1^2} \cdot \sqrt{(1.5)^2 + (0.5)^2 + (0.5)^2}} \\ &= 0.42\end{aligned}$$

$$\begin{aligned}\cosine(\text{doc4}, \text{c1}) &= \frac{0 + 1 + 0}{\sqrt{0^2 + 2^2 + 0^2} \cdot \sqrt{(1.5)^2 + (0.5)^2 + (0.5)^2}} \\ &= 0.30\end{aligned}$$

$$\begin{aligned}\cosine(\text{doc1}, \text{c2}) &= \frac{0 + 1.5 + 0}{\sqrt{2^2 + 1^2 + 0^2} \cdot \sqrt{0^2 + (1.5)^2 + (0.5)^2}} \\ &= 0.42\end{aligned}$$

$$\begin{aligned}\cosine(\text{doc2}, \text{c2}) &= \frac{0 + 0 + 0.5}{\sqrt{1^2 + 0^2 + 1^2} \cdot \sqrt{0^2 + (1.5)^2 + (0.5)^2}} \\ &= 0.44\end{aligned}$$

$$\begin{aligned}\cosine(\text{doc3}, \text{c2}) &= \frac{0 + 1.5 + 0.5}{\sqrt{0^2 + 1^2 + 1^2} \cdot \sqrt{(0)^2 + (1.5)^2 + (0.5)^2}} \\ &= 0.89\end{aligned}$$

$$\begin{aligned}\cosine(\text{doc4}, \text{c2}) &= \frac{0 + 3 + 0}{\sqrt{0^2 + 2^2 + 0^2} \cdot \sqrt{0^2 + (1.5)^2 + (0.5)^2}} \\ &= 0.99\end{aligned}$$

Hence,  $\text{c1} = \text{doc2}, \text{doc2}$  and  $\text{c2} = \text{doc3}, \text{doc4}$ .

## Lab No. 2

### News Aggregator:

1. Design a spider to extract the news(nepali) from web.  
Seed url → ekantipur, onlinekhabar, ratopati.
2. Implement the parser and filter the topics.
3. Use k-means algorithm to make the clusters of the news ( $k=5$ ).

## Lab No. 2

### News Aggregator:

1. Design a spider to extract the news(nepali) from web.  
seed url → ekantipur, onlinekhabar, ratopati.
2. Implement the parser and filter the topics.
3. Use k-means algorithm to make the clusters of the news ( $k=5$ ).

### ↳ K-medoids method:

- pick actual object to represent clusters instead of mean values.
- each remaining object is clustered with the representative object (medoid) to which it is the most similar.
- the algorithm minimizes the sum of dissimilarities between each object and its corresponding reference point.

$$E = \sum_{i=1}^k \sum_{p \in C_i} |P - O_i|$$

Idea:

- initial representation are chosen randomly

for each representative object  $O_1$

→ for each non representative object  $R$ , swap  $O_1$  and  $R$

→ choose the configuration with the lowest cost

Eg:

$O_1(2,6), O_2(3,4), O_3(3,8), O_4(4,7), O_5(6,2), O_6(6,4), O_7(7,8)$   
 $O_8(7,4), O_9(8,5), O_{10}(7,6)$

$K=2$

$O_2(3,4)$

$C_1$

$O_8(7,4)$

$C_2$

$C_1 \rightarrow \{O_1, O_2, O_3, O_4\}$

$C_2 \rightarrow \{O_5, O_6, O_7, O_8, O_9, O_{10}\}$

$$\begin{aligned} Err(O_2, O_8) &= (O_2 - O_1) + (O_2 - O_2) + (O_2 - O_3) + (O_2 - O_4) + \\ &\quad (O_8 - O_5) + (O_8 - O_6) + (O_8 - O_7) + (O_8 - O_9) + (O_8 - O_{10}) \\ &= (3+4+4) + (3+1+1+2+2) \\ &= 20 \end{aligned}$$

$$Err(O_2, O_7) = (2+4+4) + (2+2+1+3+3) = 22.$$

∴ It is bad idea to swap  $O_8$  and  $O_7$ .

## L Hierarchical Clustering:

- Start with one cluster, individual item in its own cluster and iteratively merge clusters until all the items belongs to one cluster.
- Dendograms are pictorially used to represent hierarchical clustering.

## L Approaches:

1. Single Linkage: distance between the closest members of two clusters.
2. Complete Linkage: distance between the numbers that are farthest apart.
3. Average Linkage: involves looking at the distance between all pairs and average of these distances.

Eg:  $P_1 \quad P_2 \quad P_3 \quad P_4 \quad P_5$

$P_1 \quad 0$

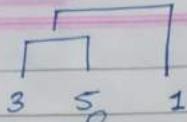
$P_2 \quad 9 \quad 0$

$P_3 \quad 6 \quad 7 \quad 0$

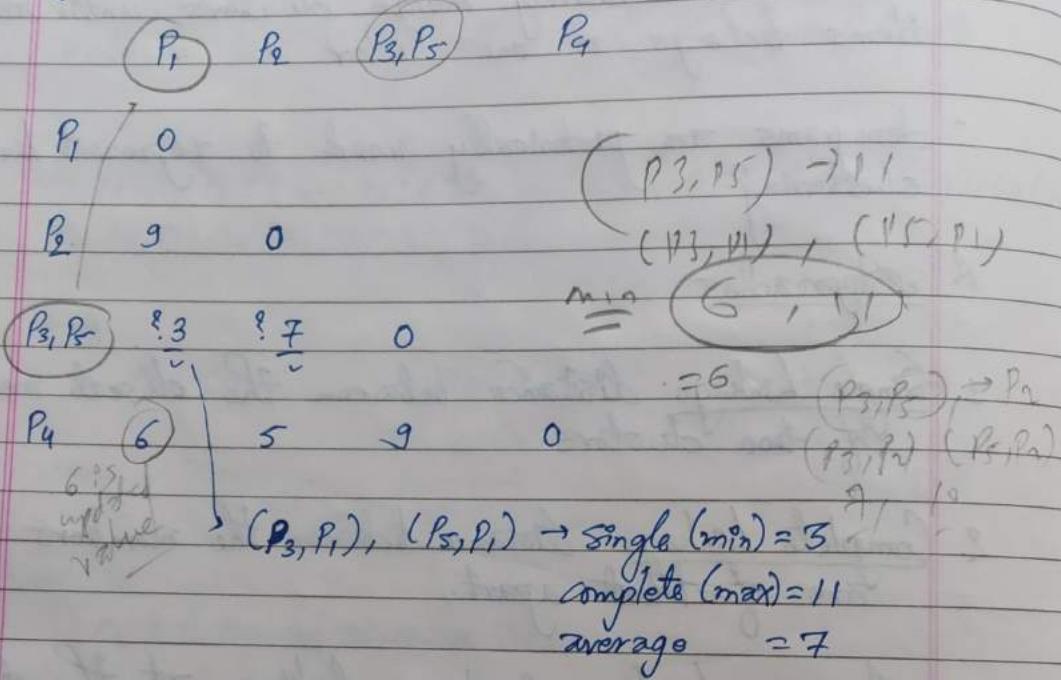
$P_4 \quad 3 \quad 5 \quad 9 \quad 0$

$P_5 \quad 11 \quad 10 \quad 2 \quad 8 \quad 0$

↑  
min



update the matrix



Q Consider the following dataset:

- $P_1 (0.40, 0.53)$
- $P_2 (0.22, 0.38)$
- $P_3 (0.35, 0.32)$
- $P_4 (0.26, 0.19)$
- $P_5 (0.08, 0.41)$
- $P_6 (0.45, 0.30)$

Apply the hierarchical clustering approach to cluster these points using complete linkage approach. Use Euclidean distance to find the similarity.

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

	$P_1$	$P_2$	$P_3$	$P_4$	$P_5$	$P_6$
$P_1$	0					
$P_2$	0.23	0				
$P_3$	0.22	0.14	0			
$P_4$	0.37	0.19	0.16	0		
$P_5$	0.34	0.14	0.28	0.63	0	
$P_6$	(0.06)	0.84	0.10	0.22	0.39	0

Merge  $P_1, P_6$ .

	$P_1, P_6$	$P_2$	$P_3$	$P_4$	$P_5$
$P_1, P_6$	0				
$P_2$	0.24	0			
$P_3$	0.22	(0.14)	0		
$P_4$	0.37	0.19	0.16	0	
$P_5$	0.39	(0.14)	0.28	0.63	0

$$d(P_1, P_2) = 0.23$$

$$d(P_6, P_2) = 0.24 \text{ choosing maximum}$$

$$d(P_1, P_3) = 0.22$$

$$d(P_6, P_3) = 0.10$$

$$d(P_1, P_4) = 0.37$$

$$d(P_6, P_4) = 0.22$$

$$d(P_1, P_5) = 0.34$$

$$d(P_6, P_5) = 0.39$$

Merge  $P_2, P_3, P_5$

	$P_1, P_6$	$P_2, P_3, P_5$	$P_4$
$P_1, P_6$	0		
$P_2, P_3, P_5$	0.39	0	
$P_4$	(0.37)	0.63	0

$$d(P_1, P_2) = 0.23$$

$$d(P_1, P_3) = 0.22$$

$$d(P_1, P_5) = 0.34 \quad \text{max}$$

$$d(P_6, P_2) = 0.24$$

$$d(P_6, P_3) = 0.10$$

$$d(P_6, P_5) = 0.39$$

$$d(P_2, P_4) = 0.19$$

$$d(P_3, P_4) = 0.16$$

$$d(P_5, P_4) = 0.63$$

Merge  $P_1, P_6, P_4$

$P_1, P_6, P_4$

$P_2, P_3, P_5$

$P_1, P_6, P_4$

0

$P_2, P_3, P_5$

0

$$d(P_1, P_2) = 0.23$$

$$d(P_1, P_3) = 0.22$$

$$d(P_1, P_5) = 0.34$$

$$d(P_6, P_2) = 0.24$$

$$d(P_6, P_3) = 0.10$$

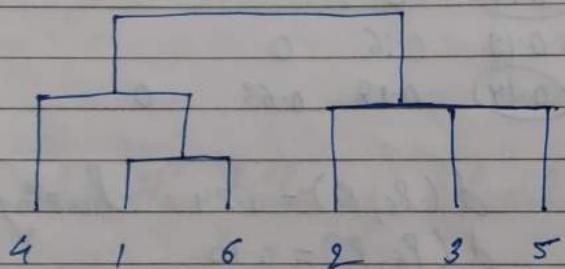
$$d(P_6, P_5) = 0.39$$

$$d(P_4, P_2) = 0.19$$

$$d(P_4, P_3) = 0.16$$

$$d(P_4, P_5) = 0.63$$

Complete linkage Dendrogram.

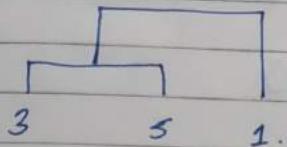


Previous Example:

	$P_1$	$P_2$	$P_3$	$P_4$	$P_5$
$P_1$	0				
$P_2$	9	0			
$P_3$	3	7	0		
$P_4$	6	5	9	0	
$P_5$	11	10	2	8	0

↑  
min

	$P_1$	$P_2$	$P_3, P_5$	$P_4$
$P_1$	0			
$P_2$	9	0		
$P_3, P_5$	3	7	0	
$P_4$	6	5	8	0



$(P_3, P_1)$  and  $(P_5, P_1)$   
single  $(\min) = 3$   
complex  $(\max) = 11$   
average = 7.

Merge  $P_1, P_3, P_5$

	$P_1, P_3, P_5$	$P_2$	$P_4$
$P_1, P_3, P_5$	0		
$P_2$	7	0	0
$P_4$	6	5	

$$d((P_1, (P_3, P_5)), P_2) = \min((P_1, P_2), ((P_3, P_5), P_2))$$

$$= \min(9, 7) = 7.$$

$$\begin{aligned} d(P_1, (P_3, P_5), P_4) &= \min(d(P_1, P_4), d((P_3, P_5), P_4)) \\ &= \min(6, 8) = 6 \end{aligned}$$

Merge  $P_2, P_4$ .

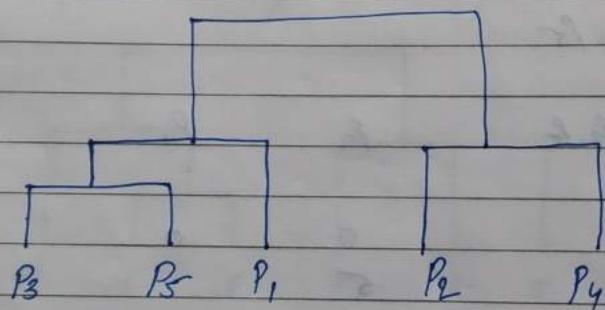
	$(P_1, (P_3, P_5))$	$(P_2, P_4)$
$(P_1, (P_3, P_5))$	0	
$(P_2, P_4)$	6	0

$$d((P_1, (P_3, P_5)), (P_2, P_4)) = \min \{ d((P_1, (P_3, P_5)), P_2), d((P_1, (P_3, P_5)), P_4) \}$$

$$= \min \{ 7, 6 \} \\ = 6$$

Merge:  $(P_1, (P_3, P_5))$  and  $(P_2, P_4)$

The single linkage dendrogram is:



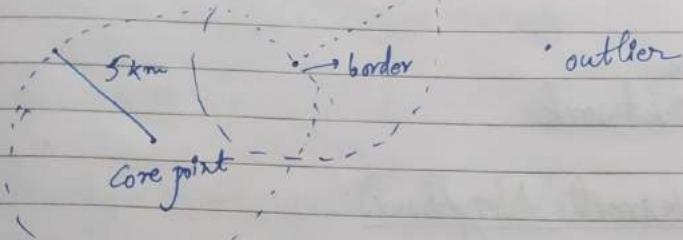
### DBSCAN:

- Density Based Spatial Clustering Application.
- Based on two parameters
  - 1.  $\epsilon$ -neighborhood  $\rightarrow$  threshold for determining the range.
  - 2. Min pts  $\rightarrow$  minimum number of points to form a cluster.

- core points, border points, outlier/noise.

$$n_{\text{pts}} = 10$$

$$\epsilon = 5 \text{ km}$$



Unit 5

## Classification

- placing the things where they belong?
- based on two

1. e Model Construction (Training Phase)

2. e Model Usage (Testing Phase)

Fr  
Paradigm:

Training Data Set:

Name	Post	Years	Pension	Machine	YES
A	Prof.	1	Yes	→ [Machine]	→ YES
B	Asst. Prof.	2	No	↓ what it learns	
C	Asst. Prof.	5	Yes	IF POST == "Prof" OR Years	
D	Prof	5	Yes		$> 5$
E	Asst. Prof	6	Yes	THEN Pension = Yes.	

Algorithm:

- KNN
- Bayesian
- Decision Tree
- Neural Network

KNN (K-Nearest Neighbour):

Algorithm:

1. For  $x$  the object  $x$
2. Compute the feature vector of  $x = d$
3. for each  $y, c(y) > \epsilon \emptyset$

$$S_y = \text{assim}(d, y)$$

4. Sort  $S_y$  in decreasing order.
5. Let  $N$  be the first  $k$ .
6. Return the majority class

Eg:

Food = turkey stuffing =  $d_1$ Food = buffalo wings =  $d_2$ Beverage = Cream Soda. =  $d_3$ Beverage = Orange Soda. =  $d_4$ Iteration:

turkey, stuffing, buffalo, wings, cream, soda, orange

 $d_1 \langle 1, 1, 0, 0, 0, 0, 0 \rangle$  $s_1$  $d_2 \langle 0, 0, 1, 1, 0, 0, 0 \rangle$  $s_2$  $d \langle 1, 0, 0, 0, 0, 1, 0 \rangle$  $d_3 \langle 0, 0, 0, 0, 1, 1, 0 \rangle$  $s_3$  $d_4 \langle 0, 0, 0, 0, 0, 1, 1 \rangle$  $s_4$ 

Classify the document,

 $\langle \text{turkey}, \text{soda} \rangle$  with  $k=3$ 

$$s_1 = \frac{1}{\sqrt{2} \cdot \sqrt{2}} = \frac{1}{2} = 0.5$$

 $\checkmark_{\text{mt}}$  $0.5 \rightarrow s_1 (d_1) \rightarrow B$  $0.5 \rightarrow s_3 (d_3) \rightarrow F \quad \checkmark_{\text{food}}$  $0.5 \rightarrow s_4 (d_4) \rightarrow F$  $0 \rightarrow s_2 (d_2) \rightarrow$ 

$$s_2 = 0$$

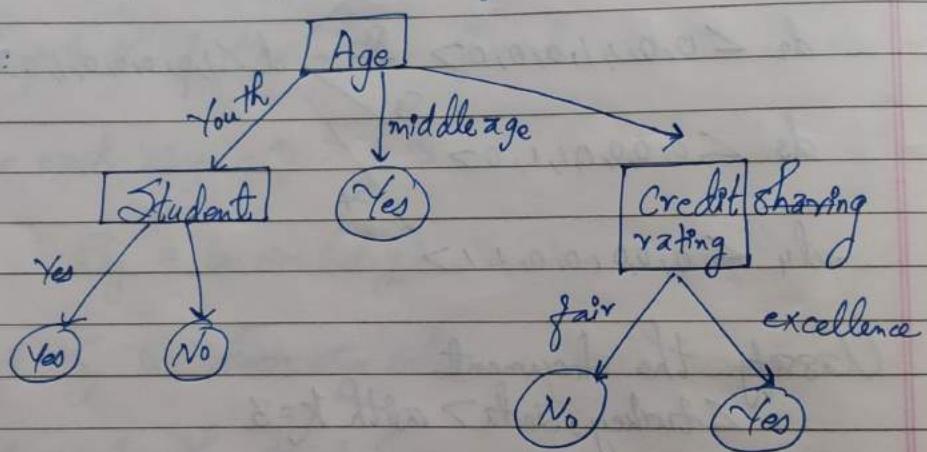
$$s_3 = \frac{1}{\sqrt{2} \cdot \sqrt{2}} = \frac{1}{2} = 0.5$$

$$s_4 = \frac{1}{\sqrt{2} \cdot \frac{1}{\sqrt{2}}} = 0.5$$

## A Decision Trees:

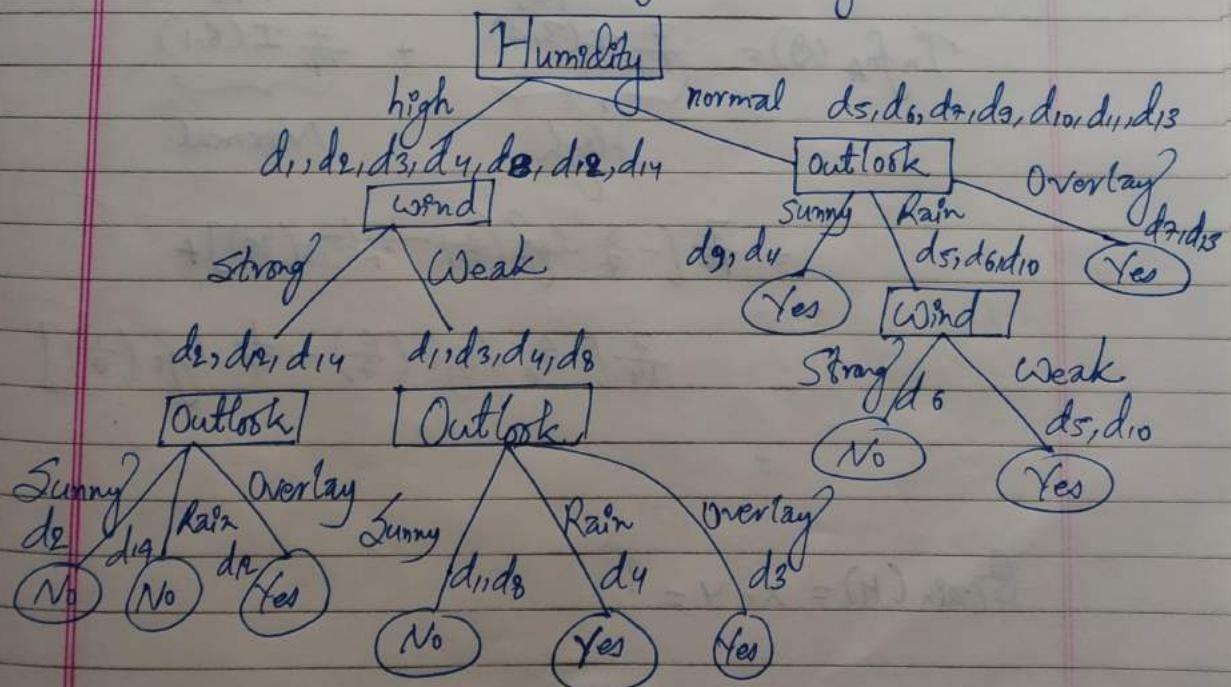
- is a flowchart like structure where each internal node represent the attribute and leaf represents the label of the class.
- The values on arc represent the possible outcomes of an attribute
- each internal node is represented by rectangle.
- each leaf is represented by oval.

Eg:



Training Dataset:

Day	Outlook	Temperature	Humidity	Wind	Play Tennis	Class
d <sub>1</sub>	Sunny	Hot	High	Weak	No	
d <sub>2</sub>	Sunny	Hot	High	Strong	No	
d <sub>3</sub>	Overcast	Hot	High	Weak	Yes	
d <sub>4</sub>	Rain	Mild	High	Weak	Yes	
d <sub>5</sub>	Rain	Cool	Normal	Weak	Yes	
d <sub>6</sub>	Rain	Cool	Normal	Strong	No	
d <sub>7</sub>	Overcast	Cool	Normal	Strong	Yes	
d <sub>8</sub>	Sunny	Mild	High	Weak	No	
d <sub>9</sub>	Sunny	Cool	Normal	Weak	Yes	
d <sub>10</sub>	Rain	Mild	Normal	Weak	Yes	
d <sub>11</sub>	Sunny	Mild	Normal	Strong	Yes	
d <sub>12</sub>	Overcast	Mild	High	Strong	Yes	
d <sub>13</sub>	Overcast	Hot	Normal	Weak	Yes	
d <sub>14</sub>	Rain	Mild	High	Strong	No	



Testing Data  $\rightarrow \langle \text{Sunny}, \text{Hot}, \text{Normal}, \text{Weak} \rangle \Rightarrow ? \quad \text{Yes}$

Attribute Selection Algorithm (ID3):

Select the information with the highest information gain.

$$\text{Gain}(A) = \text{Info}(\emptyset) - \text{Info}_A(\emptyset) \text{ where,}$$

$$\text{Info}(\emptyset) = -\sum_{i=1}^m p_i \log_2(p_i)$$

$$\text{Info}_A(\emptyset) = \sum_{i=1}^m \frac{|S_i|}{|\emptyset|} * I(S_i)$$

Eg:  $A = \text{Humidity}$

$$\text{Info}(\emptyset) = -\frac{9}{14} \log_2\left(\frac{9}{14}\right) - \frac{5}{14} \log_2\left(\frac{5}{14}\right)$$

$$\text{Info}_A(\emptyset) = \underbrace{\frac{7}{14} I(3,4)}_{\substack{\text{High} \\ \text{Yes}}} + \underbrace{\frac{7}{14} I(6,1)}_{\substack{\text{Normal} \\ \text{No}}}$$

$$= \frac{7}{14} \left[ -\frac{3}{7} \log_2\left(\frac{3}{7}\right) - \frac{4}{7} \log_2\left(\frac{4}{7}\right) \right] +$$

$$\frac{7}{14} \left[ -\frac{6}{7} \log_2\left(\frac{6}{7}\right) - \frac{1}{7} \log_2\left(\frac{1}{7}\right) \right]$$

=

$$\text{Gain}(H) = X - Y = ?$$

2. Rocchio Algorithm:

Training Phase:

- For each class  $C_i$ , calculate

$$P_i = P_0 + d \text{ where } d \in C_i$$

Testing Document:

- Let  $y$  be the featured vector of document to be tested.
- For each class find cosine similarity

$$S_i = \cos(\theta_{i,y})$$

return  $C_i$  with maximum  $S_i$

Eg: Food = turkey stuffing  $d_1$

Food = buffalo wings  $d_2$

Beverage = cream soda  $d_3$

Beverage = Orange soda  $d_4$

Testing  $\rightarrow$  turkey soda  $d$

turkey, stuffing, buffalo, wings, cream, soda, orange

$$d_1 \langle 1, 0, 0, 0, 0, 0, 0 \rangle \quad ] \quad C_1 \text{ (Food), } P_1 = d_1 + d_2 = \langle 1, 1, 1, 0, 0, 0 \rangle$$

$$d_2 \langle 0, 1, 0, 1, 1, 0, 0 \rangle$$

$$d_3 \langle 0, 1, 0, 0, 1, 1, 0 \rangle \quad ] \quad C_2 \text{ (Beverage), } P_2 = d_3 + d_4 = \langle 0, 0, 1, 0, 1, 1, 1 \rangle$$

$$d_4 \langle 0, 0, 0, 0, 0, 1, 1 \rangle$$

$$d \langle 1, 0, 0, 0, 0, 1, 0 \rangle$$

$$S_1 = \cos \theta (d, P_1) = \frac{1}{\sqrt{7} \cdot \sqrt{2}} = \frac{1}{\sqrt{14}} = \frac{1}{2\sqrt{7}}$$

$$S_2 = \cos \theta (d, P_2) = \frac{2}{\sqrt{6} \cdot \sqrt{2}} = \frac{2}{\sqrt{12}} = \frac{2}{2\sqrt{3}} = \frac{1}{\sqrt{3}}$$

Classification using Backpropagation:

1. Initialize the input and weights.
2. Calculate the output

2.1 For each input layer

$$O_j^o = I_j^o$$

2.2 For each hidden layer and output layer

$$I_j^o = \sum_k w_{kj} O_k^o$$

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

3. Calculate the error:

3.1 For output layer  $\xrightarrow{\text{Target}}$

$$E_{rrj} = O_j^o (1 - O_j^o) (T - O_j^o)$$

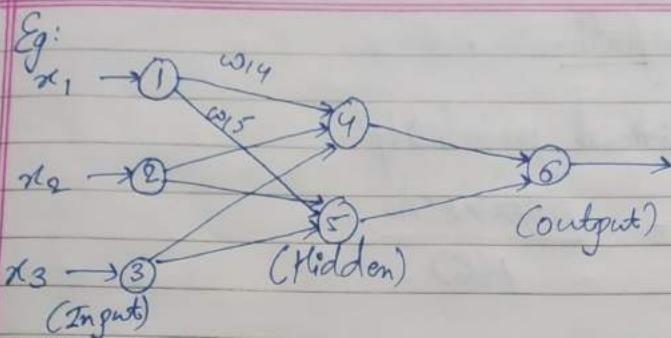
3.2 For Hidden layer

$$E_{rrj} = O_j^o (1 - O_j^o) \sum_k E_{rrk} w_{jk}^o$$

4. Update the weight

$$w_{ij}^o (\text{new}) = w_{ij}^o (\text{old}) + l \cdot E_{rrj}$$

↓ learning rate.



Initialization:

$$\begin{matrix} x_1 & x_2 & x_3 & w_{14} & w_{15} & w_{24} & w_{25} & w_{34} & w_{35} & w_{46} & w_{56} \\ 0 & 1 & 0 & 0.1 & 0.2 & 0.3 & 0.8 & 0.1 & 0.2 & 0.3 & 0.4 & 0.5 \end{matrix}$$

Input:  $O_1 = x_1 = 0$ ;  $O_2 = x_2 = 1$ ;  $O_3 = x_3 = 0$

Hidden Layer:

for  $O_4$

$$I_4 = O_1 w_{14} + O_2 w_{24} + O_3 w_{34}$$

$$Err_4 = O_4 (1 - O_4) Err_6 \cdot w_{46}$$

$$Err_5 = O_5 (1 - O_5) Err_6 w_{56}$$

$$O_4 = \frac{1}{1 + e^{-I_4}}$$

Update the weight

for  $O_5$

$$O_5 = O_1 w_{15} + O_2 w_{25} + O_3 w_{35}$$

$$w_{15}(\text{new}) = w_{15}(\text{old}) + 0.9 \frac{Err_4}{Err_5}$$

$$O_5 = \frac{1}{1 + e^{-I_5}}$$

$$w_{25}(\text{new}) = w_{25}(\text{old}) + 0.9 \frac{Err_4}{Err_5}$$

for  $O_6$

$$I_6 = O_4 w_{46} + O_5 w_{56}$$

$$w_{46}(\text{new}) = w_{46}(\text{old}) + 0.9 \frac{Err_6}{Err_6}$$

$$O_6 = \frac{1}{1 + e^{-I_6}}$$

Calculate Error:

$$Err_6 = O_6 (1 - O_6) (1 - O_6)$$

If  $I_6 Err_6 \approx 0$  Terminate

## • Bayesian Classification:

- based on conditional probability?

$$P(A|B) = \frac{P(B|A) \times P(A)}{P(B)}$$

Eg.	Age	Income	Credit Rating	Buys computer
	20-30	>10k	Fair	No
	20-30	>20k	Excellence	Yes
	30-40	>10k	Fair	Yes
	30-40	>20k	Excellence	Yes
	20-30	>20k	Excellence	Yes
	(20-30, >10k, Fair)			?

Are your result credible?  
Actual Result

Predictive Result	Yes (True Positive)	No (False Positive)
Actual Result	No (False Negative)	Yes (True Negative)

$$\text{Precision} = \frac{TP}{TP + FP}$$

$$\text{Recall} = \frac{TP}{TP + FN}$$

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

$$X = \left( \begin{array}{c} \text{Age} \\ \text{Income} \end{array}, \begin{array}{c} \geq 10k \\ \geq 20k \end{array}, \begin{array}{c} \text{Fair} \\ \text{Excellence} \end{array} \right) = ?$$

$$P(\text{Yes}/x) = \underline{P(x/\text{Yes})} * P(\text{Yes})$$

$$= \frac{2}{4} * \frac{2}{4} * \frac{1}{4} * \frac{4}{5}$$

$$P(\text{No}/x) = P(x/\text{No}) * P(\text{No})$$

(as:

- Q. Compute the entropy value at each level in previous training set using ID3.

(6x)

- significantly reduce the risk of using a classifier that will be inadequate on unseen data.

## 2 McNeuman's Test:

- comparing two classifier on the same data.

Algorithm A failed      Algorithm A succeeded

Algorithm B failed

$N_{sf}$

$N_{sf}$

$\hookrightarrow$  success/fail.

Algorithm B succeeded

$N_{fs}$

$N_{ss}$

$$z = \frac{(|N_{sf} - N_{fs}|) - 1}{\sqrt{N_{sf} + N_{fs}}}$$

$\hookrightarrow$  testing set

- when  $z=0$ , the two algorithms are said to show similar performance

## Cochran's Q Test:

- Comparing more than two classifiers

- Test the hypothesis, that there is no difference.

$$H_0: P_1 = P_2 = \dots = P_L$$

$$Q = (L-1) \frac{\sum_{i=1}^L G_i^2 - T^2}{LT - \sum_{j=1}^{N_{fs}} (L_j)^2}$$

$\hookrightarrow$  no. of classifier  
degree of freedom

$$\text{where, } T = \sum_{i=1}^L G_i$$

$G_i \rightarrow$  no. of object out of  $N_i$ s, correctly classified by  
 $D_i, i=1, 2, \dots, L$

$L_j \rightarrow$  is the no. of classifier out of  $L$  that correctly  
classified object.

### Ensemble Method:

- uses a combination of models to increase accuracy

- methods:

1. - Bagging

2. Boosting

1. Bagging:

- averaging the predictions over a collection of classifiers.

### Training:

- a classifier model  $M_i$  is learned from each training set  $D_i$  (with replacement)

### Testing:

- classify the unknown sample  $x$
- each classifier  $M_i$  returns its class prediction.
- return the majority class.

## Boosting:

- is an iterative technique which adjust the weight of an observation based on the last classification.

## Steps:

- draw a random subset of training samples  $d_1$  without replacement from the training set  $\mathcal{D}$  to train a weak learner  $C_1$ .
- draw a second random training subset  $d_2$  without replacement from the training set and add 50% of the samples that were previously falsely classified to weak learner  $C_2$ .
- find the training samples  $d_3$  in the training set  $\mathcal{D}$  on which  $C_1$  and  $C_2$  disagree to train  $\mathcal{D}_3$  a third weak learner  $C_3$ .

2 Romantic Web:

- provides a common framework that allows data to be shared with understanding of its semantic meaning.

- XML

3 RDF (Resource Description Framework)

- provides the technology for expressing the meaning of terms and concepts in a form that machine can process.

e.g: <Album>

<Singer> Naryan Gopal </Singer>

<Price> — </Price>

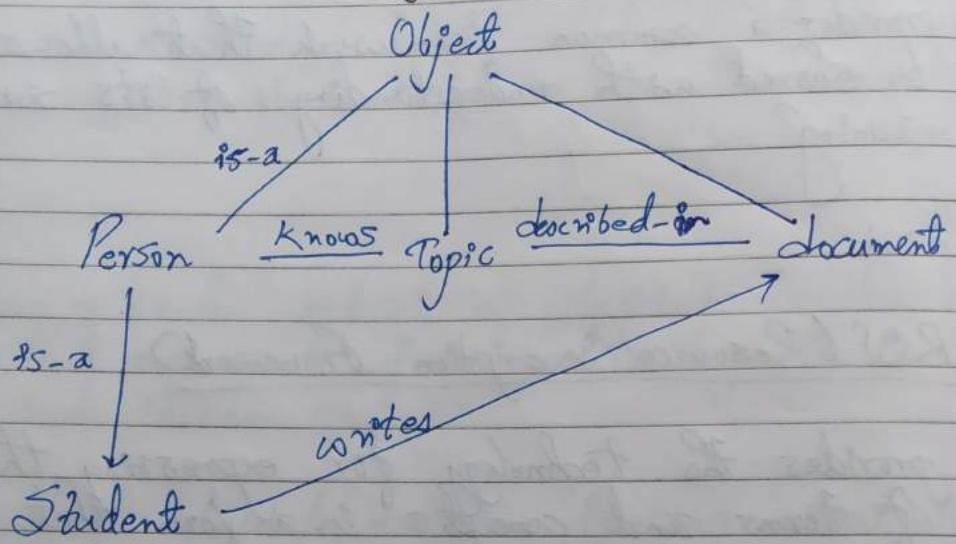
<Edition> — </Edition>

</Album>

4 Ontologies:

- defines the relation between concepts and specifies logical rules for reasoning about them.

#### 4 Sample of Ontology:



#### 4 Reasoning Over Ontologies:

- inferencing capabilities

Eg:  $X \text{ is author of } Y \Rightarrow Y \text{ is written by } X$ .