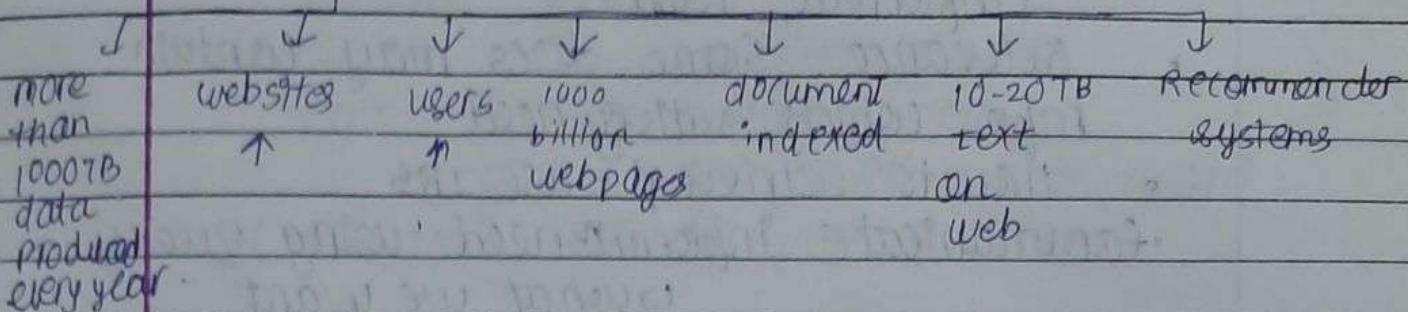


IRs→ use of IR

- more information of websites
- more choices call for relevant results
- we want personalized results
 - eg - content based recommender (movies - genre genre) - we check content of item & suggestions based on content similarity
- To deal with choices, we need IRS

→ HADOOP - distributed database clusters (HDFS)

- eg - facebook has billions of data stored in Hadoop.
- also billions of data added every day

→ key words in IRS

depend on system

- A large repo. of docs. stored on computers (corpus) (are distributed across the web).
- eg - static - file retrieved from own machine
- google - web based distribution
- dynamic → called corpus

- Topic for which we need information (information need)
- Relevance - Some docs may contain info. which satisfy need
- How to retrieve? - use IRS
- communicate informⁿ need using query
↳ what we want

→ Structures of data

↳ preferences

- IRS may extract implicit informⁿ
eg - if we buy a product, website
knows what we like - so implicit.
- explicit preference may/may not be provided.

- structured data (like DBMS)
 - clear, overt semantic structure

ID	name	-
-	-	-

↳ we can provide queries to it
for retrieval - also complex
queries (eg - complex DBMS query)

- But usually, we give simple query to unstructured data.
eg - google query

- Unstruc. data doesn't have clear, overt semantic structure.
 - It allows less expensive queries
 - eg - text on webpage etc.
 - give data with a keyword.
 - It allows less complex queries.
- ⇒ unstructured data req. IRS, as structured has DBMS queries.

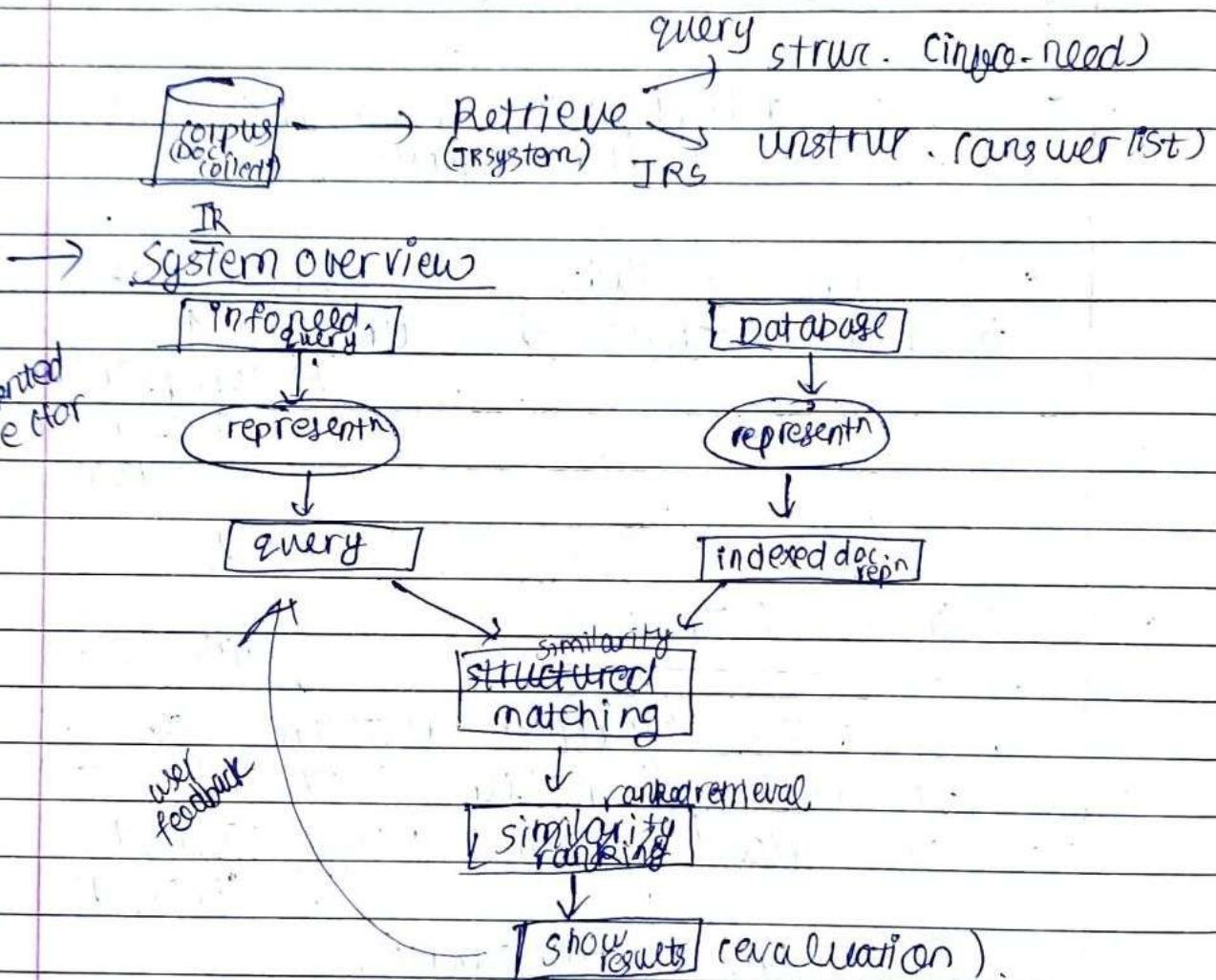
- Earlier, structured ↑, unstruc. ↓
but now, struc. ↓, unstruc. ↑
 - ↳ companies focus more here

→ Informⁿ Retrieval

- It is finding material of unstructured nature that satisfies an info. need from a large collectⁿ corpus
- Apart from web search,
eg ⇒ legal info. retrieval.
 - Laptop search
 - email search
 - corporate retrieval

- **Areal** - To find relevant informⁿ from large info-set / document set.

↳ show ranked results (for most relevance)



- ⇒ • IRS is built on corpus.
- B/w IRS & corpus we have a representation.

→ IRS vs Datamining

- • IR is ability to query a computer system & return relevant results.
eg - google web search engine
- . DM is ability to retrieve info. from one or more data sources in order to combine it, cluster it, visualize it & discover patterns in data.
↳ discover hidden pattern
- we can cluster all docs in one page so that query becomes easy.
- . BigData is the ability to manipulate huge amt. of data (eg - HDFS) in order to perform DM on the data.
↳ to reduce retrieval time.

→ Terminologies in IR

- websearch
- multimedia IR (eg - lyrics of song is query & upwarz song video)
- cross language IR (eg - query is India, we want info. related to India as well as "Bharat")
- Recommender system

how to
evaluate?

→ How to retrieve docs?

- Precision = A fraction of retrieved docs relevant to user info need.
eg- out of 100 docs, 60 retrieved, but user thinks only 40 are relevant.
- Fraction Recall = Fractⁿ of relevant docs in collectⁿ that are retrieved
 - If recall more, then more docs retrieved so precision less
 - So, we must maintain a balance b/w them both.

→ well known conferences: CA For Term paper - Only from 2020 onwards)

- SIGIR (special interest group on IR)
- ACM conference on multimedia retrieval
- KDD (knowledge discovery (ICMR) & DM)

→ eg- In 37th play docs. of shakespeare,
let our query - words having
BRUTUS & CAESER but not CALPURNIA?

- Possible answers
 - ⇒ manually search
 - ⇒ grep command on Linux
 - ↳ time consuming, one script does not work

→ Ranked retrieval - not supported using grep

Best solution - we preprocess the doc. corpus in advance in such a way we can process query faster.

→ we use - Term document incident matrices

	word1	word2	...	word37
all unique words	0	1	-	play37
word1
wordn				

• 0 - word not present

• 1 - word present

• every cell is filled in this way.

eg- doc1 = I am

doc2 → Hi I :

so, words = 3 unique

documents - now fill this

• Now for for Brutus & & , we extract row of Brutus.

Brutus → 110100

same for caesar → Bitwise

caesar → 110111 and

1100100

• Now we extract vector of calpurnia.

calpurnia → 010000

↳ we complement as we don't want not calpurnia.

so, NOT calpurnia = 10111

- Now we "and" this with Brutus & Caesar

$$\begin{array}{r} 110100 \\ 10111 \\ \hline 1100100 \end{array}$$

\rightarrow The 2 one's represent our docs which satisfy the query (play 1 & play 4)

\rightarrow Problems with term doc. incident matrix

- Frequency of word not considered
- Position of word not considered
 e.g. in research, if word present in keyword, then more relevant to query.

- e.g. $N = 1$ million docs., each with 1000 words
 - avg. 6 bytes/word
 $\therefore 6\text{GB}$ of docs.
 - if in hard disk, then everytime we need to bring it in main memory.
 - if in main memory then expensive
- Let 51akh unique words
 \therefore matrix is 51akh rows \times 1m⁹. cols.

- sparse - dispersed
- dense - collected

- also matrix is extremely sparse.
any column has max^m 1000 ones
(out of all 5 lakhs words only 1000)
- so sparsity = $\frac{1000 \times 100}{500000} = 0.2$
would be one in worst case
- (also, at max $1000 \times 1\text{mil} = 1\text{billion}$ ones in doc.)

→ Inverted index

→ better way.

- we only record those docs. in which term/word is present
- eg - Brutus only present in docs. 1, 2, 4, 11, 31,
 45, 73, 74
 they have 1
 in Brutus.

• similarly for Caesar & Calpurnia.

Brutus = (1, 2, 4, 11, 31, 45, 73, 74)

→ we can use array if static corpus

Calpurnia → (1, 5, 7) ← & Fixed array would waste
 only 3, so sparse
wanted sparse

- we can use Linked list when corpus is dynamic

- so array for time seeking time

• so, we must use optimized Data structure

- we can also use Hashing

- any row is posting list (1, 2, 3, ...)

dictionary of words:
 Brutus = (1, 2, ...)
 Caesar = (3, 4, ...)

→ postlist

(3, 4, 5, ...)

To find Brutus & Caesar using linked list, we can use merge operation of linked lists.
(Intersection)

⇒ merge → we need, atleast space of 2^2
(1 for doc.id, 1 for next pointer info.)

In merge - eg for Brutus & Caesar,

1) Locate Brutus in dictina dictionary

- retrieve its postings

2) Locate Caesar in the dictionary

- retrieve its posting.

Brutus - $2 \rightarrow 4 \rightarrow 8 \rightarrow 16 \rightarrow 32 \rightarrow 64 \rightarrow 128$

Caesar - $1 \rightarrow 2 \rightarrow 3 \rightarrow 5 \rightarrow 8 \rightarrow 13 \rightarrow 21 \rightarrow 34$

Now we apply the merge function by running through the 2 postings simultaneously
(like in merge sort) $O(n^2)$ worst case

intersect
function

if data at 2 pointers running
pointers match, then we store
it in answer & increment
both pointers

else, the pointer having in
a case of mismatch, pointer
having lesser data gets incremented

word

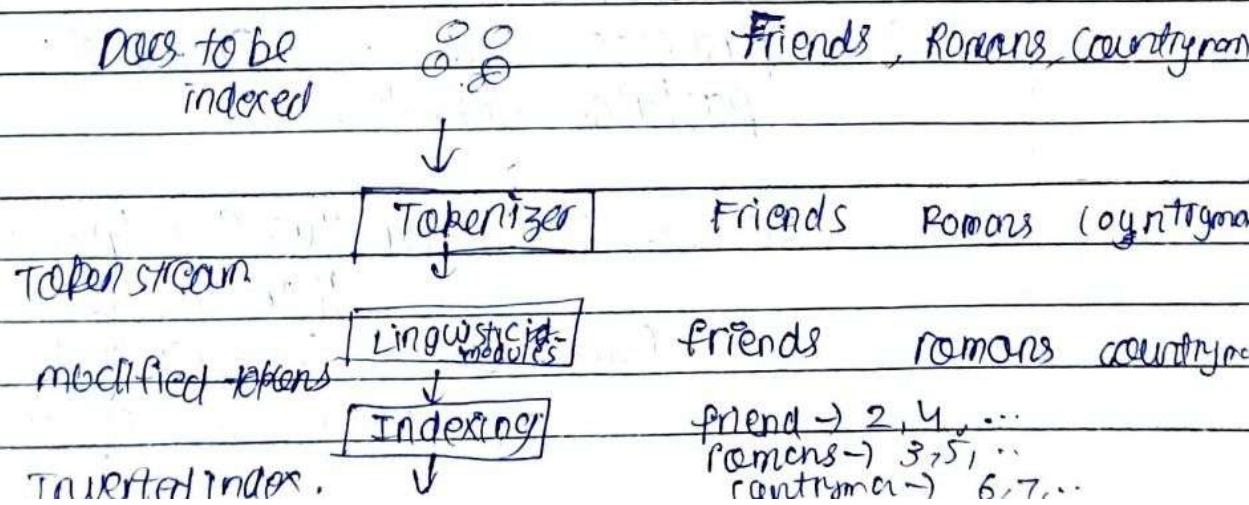
- To create a term from document, till now we just created the term.
- In every row, no repetition of term.
- If term present, put 1, else 0. By Boolean Retrieval.
- So, before the query, we need to present the document for constructing the matrix.
- Then we extract vectors based on our query.

$$\text{Sparsity} = \frac{\text{total no. of words in doc}}{\text{total unique words}}$$

Sparsity is biggest problem of Term doc. matrix

- Inverted index solves these problems
- we only store those docs. containing the words
- The words are dictionary & the docs. present (vector) is posting.
- To process query, we use merge operation.

→ How are inverted index prepared?



- steps of indexer
- To generate tokens we pass the docs & identify the terms.
 - Then we arrange the terms in alphabetical order.
 - Accordingly, ids are assigned.

words	id
caesar	1
caesar	2

- We convert this to document frequency.

doc. frequency.

caesar	1	2
--------	---	---

~~every optimizer~~ → Eg- query → Bratus & calpurnia & caesar.
 what is the best order for optimizn?

Eg - Bratus | 2 | 4 | 8 | 9 | 10 | 11 |

calpurnia → ---

caesar → ---

- We process the 2 smallest postings first.

so, order → Bratus & calpurnia
 then caesar.

(order of ↑ frequency).

→ process the smallest one first.

- Boolean model is the simplest model.
- It was primary commercial tool for 3 decades
eg- WestLaw uses it.

→ Before parsing a doc., we need to see :-

~~file format~~

- what is format of file (pdf, word, xml etc.)
- language of doc. (eg - `mail=french` attachment=german)
- Type of charset (CP1252, UTF-8 etc.)
- . Each of this is a classification problem.

- Docs. being indexed may be of diff. language - a single ~~single~~ index multiple lang
- Multiple lang in 1 doc. may exist
- Open source libraries may help

→ To what unit we consider a doc.?

- eg- 1 ppt is a doc. or 100 ppts. a doc.
- 5 attachments → 5 docs or 1 doc. (in emails).
- a group of files (eg- ppt of latex split over HTML pages)

→ some preprocessing techniques

① Tokenization → to generate tokens we parse the doc.

- tip = "friends, humans & country man"
- after each space/ punctuation we generate a token (usually) (here 4)

- Token is an instance of sequence of characters.
- each token generated is a candidate for index entry in dictionary.

Issues in Tokenizn

1) eg- Finland's capital → Finlands

~~only 3 equivalence class~~ → - Finland and s? } all 3 diff.
- Finlands
- Finland's Then how to compare.

eg - Hewlett - Packard

- Hewlett & packard as 2 tokens?

eg - state - of - the - art

- CO - education

- lowercase, lower-case, lowercase?

• It can be effective to get the user to put it in possible hyphens

eg - San Francisco - 1 or 2 tokens?

2) Date Format

- 3/20/91, Mar. 12, 1991 20/3/91

- Aadhar/License nos. ispace

- (891) 123 456

- (891) 123 - 234

3) → Language

- eg - some languages ~~do~~ do not have space. (German)
- in some languages we read from right to left (& nos. from left to right)
 - ↳ Japanese.
 - ↳ mix. of languages.

② ⇒
*intrinsic
modules*

Stop words

- with a stoplist, you exclude from the dictionary entirely the commonest words.
eg - a, an, the ...
 ↳ little semantic content
- They occupy 30% of doc. for top 30 words
- But in flights, A to B, King of Denmark
 ↳ to is imp. here.
 - so this stopword important.
- If space is an issue, we need to remove them.
- But if space not an issue, we can index them as terms.



(3) →

Normalizⁿ of terms

→ all mapped to same token
↳ equivalence class

- we need to normalize words in indexed as well as query words into the same form.
 - eg - to match U.S.A. & USA
- Result is terms - a term is a normalized word type, which is an entity in our IR system dictionary.
- so, we need to make system smart.

→ But,

French - résumé & resume

they become same
after
normalizⁿ

German - Tübingen & Tübingen

→

case folding

- convert everything to Lowercase.
- ↳ issues

↳ eg - General Motors (company)

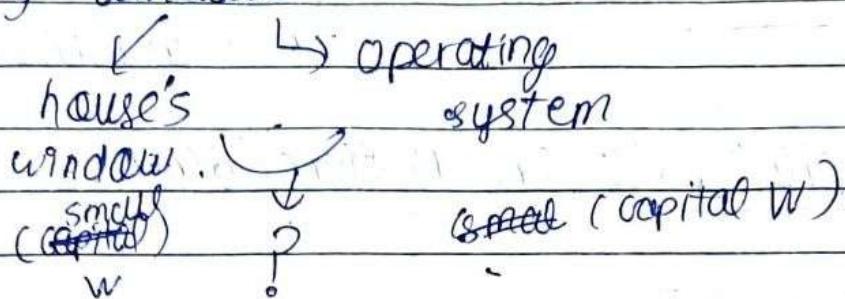
- fed or fed
- SAIL, sail

(6) →

→ exam

- query - C.A.T., we must not get cat.

eg- window



(4) → The Sauro & soundex

- Handling synonyms & homonyms -
- car = automobile. } Both should be retrieved
- color = colour }
- equivalence classes must be done →
- spelling mistakes can be solved by soundex using phonetic heuristic

(5) → Lemmatizⁿ → but does handle proper noun (e.g. Americans)

- converts word to root-form
- eg. car am, are, is → be
car, cars, car's, cars' → car.
- Reduce inflectional / variant forms to base forms.

(6) → Stemming → language dependent.

chops off the words.

eg- automate(s), automation, automation,
→ automat.

- Reduce terms to their roots before indexing.

→ Potter's ~~stammerer~~ → an algorithm for stemming.

~~weight
of word
sensitivity
rules~~ ← Typical rules in Potter's algo. →

- 1) - sses → ss
 - 2) - ies → i
 - 3) - ~~e~~ ational → ate
 - 4) - tional → ~~t~~ion
- } chopping rules

after chopping if $(m > 1)$,

- ENENT

eg- replacement → replace ($m > 1$) ✓
cement → ~~cem~~ (

X

↳ $m = 1$ not > 1 .

so, no chopping

: cement remains

eg- seperational → seperate (3rd rule)

⑦ part of speech tagging → Tags every word (POS)

• Required in recommender system.

• It is defined as the task of labeling each word in a sentence

with its appropriate part of speech

eg- whether word is noun, adjective etc



→ POS is good for -

- IRG
- Text to speech
- word sense disambiguation
- Recommender system is of 2 types
 - ⇒ non personalized

↳ average rating considered
(like most popular restaurant etc.)

and shown to every user.

⇒ personalized - based on personal user interest.

content based
eg - i've
watched several
movies

• It works based
on feature similarity
of items (e.g. - songs,
director, genre
etc.)

- problem :-
always only
similar item
shown, no
diversity

user doesn't
give review
perfectly
but on avg

collaborative
filtering

↳ similarity of
users

eg - A likes 5 movies
B likes 5 movies
Both give similar
ratings to movies.

we
cannot
understand
ratings
just based
on stars.

• so, A & B have
similar choices

• so, a new movie
comes. A watches
it. If A gives it
high rating,
it can be recommended
to B.

there
are many
features
(fine grains)

- Instead of rating, we can use sentiment analysis
- we can analyse user comments on features of item.
 - 1) If we analyze entire review, then its document analysis
 - 2) Other is feature / aspect level analysis.

eg- mobile is

good
as and I bought it...

- opinion = good

+ve

e.g- mobile's voice quality
is good

• Here feature = voice
quality,
opinion is good
+ve.

3) sentence level

- we analyse the sentence.

eg- The phone is bad.

- opinion = -ve.

→ For every feature, we extract the rating & calculate similarity based on user.

• How to extract feature & opinion from review?

→ If a word is a frequent noun, then its a feature.

- we can use POS tagging for this.

→ opinion could be an adjective / adverb

- we can see nearer words, if there are any adverb or adjectives, they can be opinions.

* Heuristic (doesn't give optimized answer but relatively good so in)

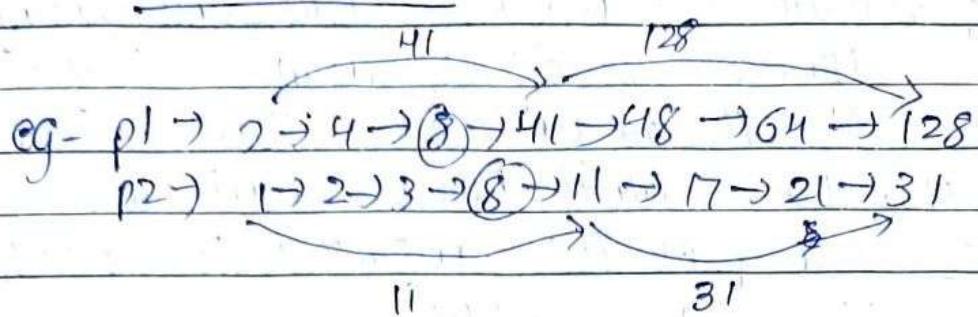
post1BH post1BT2
size size

Page No.

Date:

merge gives $O(m+n)$

→ How to improve posting list merges through skip pointer



• But where do we place skip pointer?

a heuristic approach \Rightarrow if length of list is n , at the gap of \sqrt{n} , place a skip pointer
(eg- if list = 16, skip pointer at $\sqrt{16} = 4$ intervals)

eg- if we have processed 8,
p1 reaches 41

- can be used, for root
be and "or"
- now p2 reaches 11.
 - is a skip pointer available at 11?
 - Yes there is.
 - also, $31 < 41$
 - so, we skip, and p2 reaches 31.

→ Trade off

more skips \rightarrow shorter skip spans \rightarrow more likely to skip

But lots of comparisons to skip pointers.

Fewer skips \rightarrow few pointer comparisons, but then long long skip spans \rightarrow few successful skips.

Best case $\Rightarrow O(\sqrt{n})$

I/O oper's are also imp.

→ Phrase queries

Eg - we want to retrieve "Stanford University"

- a phrase, we do not want "I am in Ahmedabad university", then it should this doc. should not be retrieved.

• So, our previous method is not good for phrase queries.

• There also, may not be space

⇒ "Stanforduniversity"

↳ User is free to enter

↳ 10% of queries have phrase terms.

heuristic) • we change the index ⇒ biword indexing
• simple approach

• If doc. contains - "Friends, Romans, countryman", we had 3 terms
• with biword indexing - 2 terms

↳ "Friends Romans" & "Romans countryman".

• If only "Romans," then none of the 2 terms match.

• But in "state-of-the-art", there is no guarantee for answer

• we can get false +ve

↳ model predicted +ve,
but actually not.

• In 4 to 5 words in the query,
biword indexing fails.

eg - Longer phrases can be processed by breaking them down.

- eg - Stanford university palo alto

↳ "stan. univ." & "palo alto".

↓
not carried

Problem - Index blowups due to more no. of words.

heuristic)

- Positional index help.

- we can ~~store~~ pos^n of word in the document.

- till now:-

$t_1 \rightarrow [doc] \rightarrow [doc]$

storing &
pos^n
in doc

} we store
doc. in which
term present
& pos^n where
term present.

eg - <bc: 993427;

bc appears in 1, 2, 4, 5

1:	7, 18, 33, 72, 86, 231	(bc is present in pos ⁿ 7, 18... in doc1)
2:	3, 149	
4:	17, 191, 291, 430, 434	
5:	363, 367, ... >	

→ doc. freq → in how many docs. is term present
 → count freq in doc1 10 times + in doc2 20 times
 $\therefore 10+20 = 30$

- For phrase queries, we use merge alg. recursively for at doc. level.

→ eg - we want \Rightarrow "to be or not to be"

\Rightarrow to :

\rightarrow doc. no. 2 \rightarrow posfr.
 \rightarrow 2: 1, 17, 79, 222, 551; 4: 8, 10, 190, 429, 433
; 7: 13, 23, 191 ...

\Rightarrow be :

\rightarrow 1: 17, 19 ; 4: 17, 191, 291, 430, 434;
5: 16, 19, 101 ...

\Rightarrow Here, 429, 430, 433 & 434 is possible
if explored.

• we can get "to be or not to be"
429 & 430 are consecutive.

→ Vector space / ranked retrieval / TFIDF model (RR)

retrieved results
have order;
so fastest
or failing
(\neq doc
or many
doc)
not a problem

- earlier in boolean retrieval - doc. either match or not
- Boolean is not good for user experience also
- "applied" processing is too much
- users wait too much in worst case
 - ↳ too many words then too much time.

- doc. gets retrieved, or entirely not retrieved.
 - many doc. may be retrieved, or 0 maybe retrieved.
-
- Rather than boolean model, in ranked retrieval system returns an ordering over the top docs. in called "for query"
 - Free text query - Rather than query or operators, user's query is in human language
 - Ranked ret. is used for free text processing.
 - Here, Large result set is not an issue.
 - we just show the top $K (\approx 10)$ results
 - we don't overwhelm the user.
 - Scoring PS basis for R.R
 - we assign score in $[0, 1]$
 - no ↪ perfect match
 - If query term doesn't occur in doc, score=0
 - Else if more frequent query for doc., score higher.

→

→ Jaccard similarity (b/w A & B)

$$\Rightarrow \text{Jac}(A, B) = |A \cap B| / |A \cup B|$$

$$\Rightarrow \text{Jac}(A, A) = 1$$

$$\Rightarrow \text{Jac}(A, B) = 0 \text{ if } A \cap B = 0$$

- A & B don't have to be same size

- Jac(A, B) assigns Value in range 0 to 1.

eg - doc. 1 = person died in march

doc 2 = the long march

query = idea of march.

→
Term
doc.

count
means

$$\text{Jac}(q, d_1) = \frac{1}{6}$$

$$\text{Jac}(q, d_2) = \frac{1}{5}$$

- Here as d_2 is shorter, it is more relevant to query

- ↳ if is not necessarily true

- ↳ this is a problem

- with Jac. sim.

- If $d_2 = \text{long idea of march}$

$$\text{Jac}(q, d_2) = \frac{2}{5}$$

Problems
of Jaccard

- ↳ Jac. doesn't consider term freq
- ↳ Rare terms are more imp. than frequent terms

- we need normalization to make vector magnitude same.
- To prevent shorter-longer doc problem, we normalize them.

For normalization, $\text{Term A B I A N B I} / \sqrt{\text{A B I}}$

\downarrow
to normalize

\rightarrow Now lets consider freq. of words in doc.
 Term doc. count means \rightarrow \hookrightarrow so each vector is a count vector document.

e.g. - Antony

doc A (55) Antony appears 55 times in A.

\rightarrow Considering bag of words model for RR, we ignore position & index for RR

\therefore A B C & C B A are same vectors
 \hookrightarrow bag of words model (no ordering)

The term freq. (tf) = no. of times word occurs in doc.

A doc. with 10 occurrences of a term is more relevant than 1 occurrence of a term.

\Rightarrow * (But not 10 times more relevant!)

\therefore Relevance doesn't increase proportionally with tf

: log freq. weighting:-

- Log freq. of term t in doc d

$$\Rightarrow w_{t,d} = \begin{cases} 1 + \log_{10} t f_{t,d} & \text{if } t f_{t,d} > 0 \\ 0 & \text{otherwise} \end{cases}$$

becomes

$$0 \rightarrow 0 \quad \text{after } w_{t,d}$$

$$1 \rightarrow 1$$

$$2 \rightarrow 1.3$$

$$10 \rightarrow 2$$

$1000 \rightarrow 4$ } these are weights of relevance.

\therefore Not linearly increases.

$$\therefore \text{Score} = \sum_{t \in \text{doc}} (1 + \log_{10} t f_{t,d})$$

every document has common terms.

The score is 0 if every term not present in the doc.

- Rare terms are more informative than frequent terms.
 - like Stopwords.

* log change magnitude,

not direct
(so, normalize)

e.g. "arachnacritic" is a rare word.

↳ give more weightage.

so, a doc query with this word, then the doc. containing the word is highly relevant.

TO give weightage to the terms
(based on how many docs. is it present)

df_t is term doc. frequency of t .

df_t is an inverse measure of informativeness of t .

$$df_t \leq N$$

we define

$$idf_t = \log_{10}(N / df_t)$$

e.g. for 1 mil. docs,

if California appears only in 1 doc.

$$\therefore i.e. df_t = 1,$$

$$idf_t = \log_{10} \left(\frac{10^6}{1} \right) = 16$$

if "the" appears in all docs,

$$\therefore i.e. df_t = 10^6$$

$$idf_t = \log_{10} \left(\frac{10^6}{10^6} \right) = 0$$

↳ no info.

→ collect frequency of t is total

occurrences of t in entire corpus/all docs.

so, it doesn't help us in knowing importance of that word.

↳ it may appear in only 1 doc, but sometimes, $idf_t = 10000$, is bad

→ T-f-idf indexing, ~~do not normalize if given then do not consider~~

$$W_{i,d} = \left(\log \left(\frac{1 + f_{i,d}}{1 + f_{i,d}} \right) \right) \times \log \left(\frac{N}{d_f} \right)$$

↑
occurrences
are relative
frequent
rarity

→ We represent our documents as vector.

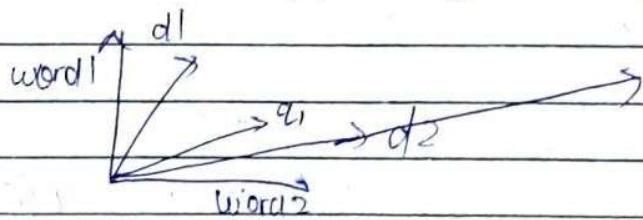
- Terms are axes
- Doc. are points (vectors)
- Very high dimensional
- Very sparse matrix

→ Query also must be vector.

- Represent queries as vectors.
- Rank docs. according to prox. similarity of query

similarity \propto 1 / distance.

- Euclidean distance is not a good measure. It is large for vectors of different lengths.



- even though q_1 , q_2 , d_1 , d_2 are similar distribution, Euclidean may give large distance

Cosine
similarity
is a good
measure

Instead of distance, we could use similarity measure.
eg - cosine similarity.

It gives b/w 0 & 1

as cos is b/w 0 & 1

✓
90° → fully match.
↳ vectors overlap
do not match.

near to 1 → both match

near to 0 → do not match.

$$\text{cosine } (v_1, v_2) = \frac{\vec{v}_1 \cdot \vec{v}_2}{|\vec{v}_1| |\vec{v}_2|}$$

→ normalized vector length cosine sim :-

$$\cos(\vec{q}, \vec{d}) = \vec{q} \cdot \vec{d}$$

eg - query → "gold silver truck"

d_1 = "shipment of gold damaged in a fire"

d_2 = "delivery of silver arrived in a silver truck"

d_3 = "shipment of gold arrived in a truck"

→ unique terms = shipment, of, gold, damaged, in, a, fire, delivery, silver, arrived, truck

ordered → a, damaged, delivery, fire, gold, in, of, arrived, silver, shipment, truck

	d	d ₁	d ₂	d ₃	df
a	0	0	0	1	3
damar arrived	0	0	1	1	2
damaged	0	1	0	0	1
delivery	0	0	1	0	1
fire	0	1	0	0	1
gold	1	1	0	1	2
in	0	1	1	1	3
oil	0	1	1	1	3
silver	1	0	2	0	8
shipment	0	1	0	1	2
truck	1	0	1	1	3

- Here, corpus is small, hence we didn't normalize.
- But if yes, then we would have put

$$\text{eg. } 1 + \log_{10} 2$$

df	idf $\Rightarrow C \log \frac{N}{df+1}$	weights	d ₁	d ₂	d ₃	$\Rightarrow 1.3$
3	0.47		0	0	0	
2	0.17		0	0.17	0.17	$\rightarrow idfxtf_{d_3}$
1	0.47					
1	0.47					
1	0.47					
2	0.17					
3	0					
3	0					
1	0.47					
1	0.47					
2	0.27					

\rightarrow this d₁ is d₁ also

df

* 3

2

1

now cosine similarity with $d_1, d_2 \text{ & } d_3$
is calculated with Q.
most similar is the answer.

Q-
 d_1 = supervised unsupervised artificial learning artificial
 d_2 = intelligence machine learning display
 d_3 = graphics processor processor
 d_4 = processor intelligence artificial
 d_5 = raster display processor

const. normalized term freq. - inverse
doc. freq. index per the doc. collect.

unique words = supervised, unsupervised, artificial,
learning, intelligence, machine, display,
graphics, processor, raster

	d_1	d_2	d_3	d_4	d_5	d_6	idf	$tf \times idf$
artificial	2	0	0	1	0	2	0.397	1.477
supervised	1	0	0	0	0	1	0.698	0.698
unsupervised	1	0	0	0	0	1	0.698	0.698
learning	1	1	0	0	0	2	0.397	1.301
intelligence	0	1	0	1	0	2	0.397	0.397
machine	0	1	0	0	0	1	0.698	0.698
display	0	1	0	0	1	2	0.397	1.301
graphics	0	0	1	0	0	1	0.698	0.698
processor	0	0	2	1	1	3	0.221	1.602
raster	0	0	0	0	1	1	0.698	0.698

TF-IDF Numericals

→ Resolve 2 quest'n's (practice)

(Q1) query = "gold silver truck"

d1 = "shipment of gold damaged in a fire"

d2 = "delivery of silver arrived in a silver truck"

d3 = "shipment of gold arrived in a truck"

→ Unique words = gold, silver, truck, shipment, of, damaged, in, a, fire, delivery, arrived,

$\frac{\log(\frac{1}{df})}{\log(\frac{N}{df})}$ normalized to 1.0

query-idf	word → doc →	query	d1	d2	d3	df	idf	H _{avg} f _p	1/H _{avg} f _p	1 + log(1/f _p)	1 + log(1/f _p) / idf	1 - ap / tp
0	a	0	1	1	1	3	0	1	1	1	1	0 0 0
0	arrived	0	0	1	1	2	0.176	0	1	1	1	0 0.176 0.176
0	damaged	0	1	0	0	1	0.477	1	0	0	0	0.477 0 0
0	delivery	0	0	1	0	1	0.477	0	1	0	0	0.477 0 0
0	fire	0	1	0	0	1	0.477	1	0	0	0	0.477 0 0
0.176	gold	1	1	0	1	2	0.176	1	0	1	0.176	0 0.176
1.0	in	0	1	1	1	3	0	1	1	1	0	0 0 0
1.0	of	0	1	1	1	3	0	1	1	1	0	0 0 0
0	shipment	0	1	0	1	2	0.176	1	0	1	0.176	0 0.176
0.477	silver	1	0	2	0	1	0.477	0	1.3	0	0	0.620
0.176	truck	1	1	0	1	2	0.176	0	1	1	0	0.176 0.176

$$\text{Spm. of query \& d1} = 0 + 0 + 0 + 0 + 0.176 \times 0.176 + 0 + 0 + 0 + \\ 0.477 \times 0 + 0$$

$$\sqrt{(0.176)^2 + (0.477)^2 + (0.176)^2}$$

$$\times \sqrt{(0.477)^2}$$

$$+ (0.477)^2$$

$$\sqrt{(0.176)^2 + (0.176)^2}$$

$$= \frac{0.03}{0.387} = 0.077$$

Sim. of query & d₂

$$\begin{aligned}
 &= \frac{0.176 \times 0 + 0.477 \times 0.62 + 0.176 \times 0.176}{\sqrt{(0.176)^2 + (0.477)^2} \times \sqrt{(0.176)^2 + (0.477)^2}} \\
 &\quad + (0.62)^2 + (0.176)^2 \\
 &= \frac{0.326}{0.441} = 0.739
 \end{aligned}$$

sim. of query & d₃

$$\begin{aligned}
 &= \frac{0.176 \times 0.176 + 0.477 \times 0 + 0.176 \times 0.178}{\sqrt{(0.176)^2 + (0.477)^2} \times \sqrt{(0.176)^2 + (0.178)^2}} \\
 &= \frac{0.06}{0.189} = 0.317
 \end{aligned}$$

So, similarity ranking:-

d₂ > d₃ > d₁ for query



(Q2) d1 = "supervised unsupervised artificial learning artificial"

d2 = "intelligence machine learning display"

d3 = "graphics processor processor"

d4 = "processor intelligence artificial"

d5 = "raster display processor".

⇒ construct tf-idf matrix.

unique words = supervised, unsupervised, artificial, learning, intelligence, machine, display, graphics, processor, i.e. raster

words →	d1	d2	d3	d4	d5	tf	idf	1	2	3	4	5	$\frac{\log(N/d)}{1 + \log(d)}$	$tf \cdot idf (w)$						
artificial	2	0	0	1	0	2	0.397	1.3	0	0	1	0	0.397	0	0	0.397	0			
display	0	1	0	0	1	2	0.397	0	1	0	0	1	0	0.397	0	0	0.397	0		
graphics	0	0	1	0	0	1	0.699	0	0	1	0	0	0	0.699	0	0	0.699	0		
intelligence	0	1	0	1	0	2	0.397	0	1	0	1	0	0	0.397	0	0.397	0	0.397	0	
learning	1	1	0	0	0	2	0.397	1	1	0	0	0	0.397	0.397	0.397	0	0	0	0	
machine	0	1	0	0	0	1	0.699	0	1	0	0	0	0	0.699	0	0	0	0	0	
processor	0	0	2	1	1	3	0.221	0	0	1.3	1	1	0	0	0.221	0.221	0.221	0.221	0.221	0.221
raster	0	0	0	0	1	1	0.699	0	0	0	0	1	0	0	0	0	0	0	0.699	0.699
supervised	1	0	0	0	0	1	0.699	1	0	0	0	0	0	0.699	0	0	0	0	0	0
unsupervised	1	0	0	0	0	1	0.699	1	0	0	0	0	0	0.699	0	0	0	0	0	0

→ Document classification

e.g. Naive Bayes, Bernoulli

→ Role of doc. classifier in IRS

↳ entertainment

↳ doc1

↳ politics

↳ doc3

} mapping of doc.

- we have training & testing.

- Repo. of doc. with labels available
(e.g. doc1 :- sports etc.)

- when new doc. comes, we can automatically classify entertainment or politics (others)

- By creating clusters, we can reduce our search space.
- more relevant queries
- reduction of search time. (optimiz'g of retrieval func)
- e.g. in email :- spam or not spam
- more precision & recall.
- if image, audio, videos etc. already classified, then helpful for specific query like - cat image

- Naive Bayes most popular for text classif'n.

→ Supervised Learning

- A doc. d

$$\text{classes} = C = \{C_1, C_2, \dots, C_j\}$$

- A training set D of doc. with a label from C.

Bernoulli trial \Rightarrow variable, 2 values \Rightarrow over
multivariate Bernoulli trial \Rightarrow multiple variables with over
for each test.

$$P(C|d) = P(C) \times P(d|C)$$

$P(d)$

Page No. _____ Date _____

For Naive Bayes

Determine

- A learning classifier to classify class y .
- we assign class y to d
 $y(d) \in C$
- we train a model to classify on test object/dan.

\Rightarrow we use Naive Bayes

$$P(A|B) = P_B A$$

$$P(C|d) = P_C C \prod_{k=1}^n P(t_k | c)$$

prob. of
 d belonging
to class C

t_k is 1nd

multiplication of term
presenting
for all terms

division term
ignored assume
for all $\frac{3}{5} \geq \frac{2}{5}$

If class c is popular, then new doc. is highly probable for new test class d .

(eg 90 of 100 are C , so C is probable for test d).

$n_d = \text{length of doc. (no. of tokens)}$

including t_k . $P(t_k | c)$ is conditional prob. of term t_k occurring in doc. of class C

$P(t_k | c)$ is a measure of how much evidence t_k contributes that c is correct class

(\oplus) we check for each term t , how many times it appears in all docs. of c .

- Our goal in naive Bayes is to find best class.
- Best class is most likely or max^m a posterior (MAP) class. Cmap is calculated use argmax.

$$C_{map} \quad P(C) = \frac{N_c}{N} \xrightarrow{\text{docs of class } c} \text{Prior prob.}$$

N total docs in corpus

$$\text{conditional prob.} = P(+|c) = \frac{T_{ct}}{T_{t+}}$$

eg. + appears
 3, 2, 4 | + in my
 total 23,
 $\therefore T_{t+} = 3 + 1$
 ≈ 6

• we make naive bayes posit' al independence assumption

$$P(t_k, | c) = P(t_k | c)$$

1st place or any other posit'

\hookrightarrow same prob.

• we need to add 0's conditional probability

due to 1 term which may not be in training phase, prob. becomes 0
 \hookrightarrow so need for Laplacian correction.

eg. class a, b, c

$$P(C|W/a) = \begin{cases} 0 & \text{if } W \text{ not in } a \\ 1 & \text{otherwise} \end{cases}$$

So, we add +1 to numerator, & no. of classes in denominator.

~~calculated~~ ∵ $P(c \rightarrow c) = \text{P}(c)$ - Test: chinese chinese chinese Tokyo
Japan.

3 classes - c

1 class - j

$$P(c) = \frac{3}{4}$$

$$P(j) = \frac{1}{4}$$

$$P(\text{chinese} | c) = \frac{5+1}{8+6} = \frac{6}{14} = \frac{3}{7}$$

↑ Chinese in c
↑ Abundance correct
↓ Total Chinese in total of C
↓ Total words in c

$$P(\text{c} \rightarrow \text{Tokyo} | c) = \frac{0+1}{8+6} = \frac{1}{14}$$

$$P(\text{c} \rightarrow \text{Japan} | c) = \frac{0+1}{8+6} = \frac{1}{14}$$

$$P(c \rightarrow c | d_s) \propto \frac{3}{4} \times \left(\frac{3}{7}\right)^3 \times \frac{1}{4} \times \frac{1}{4}$$

→ Chinese appears 3 times.

- Similarly for $P(\text{chinese} | j)$, $P(\text{Tokyo} | j)$
 $\& P(\text{Japan} | j)$
 multiple. $\Rightarrow P(g | d_s)$.
 whichever greater is class.

$$P(w|c) = \frac{\text{count}(w, c) + 1}{\text{count}(c) + IV \text{ to vocabulary}}$$

Bernoulli

- Doesn't see term freq. only presence or absence.
- $P(t|c)$ estimated as prob' of doc. of class containing term t .
- Probability of non occurrence is also calculated.
- Binary values put - 0, 1.

$$\text{e.g. } P(\text{chinese}|c) = \frac{3+1}{3+2}$$

we calculate probability of each unique term in corpus, to check if term may not appear.

$$\Rightarrow P(\text{chinese}|c) = \frac{\text{in how many it appears}}{\text{total no. of terms}} = \frac{3+1}{3+2} \rightarrow \text{apart from}$$

total
docs
of

total classes
(Eg)

$$P(\text{Japan}|c) = P(\text{Tokyo}|c) = \frac{0+1}{3+2} = \frac{1}{5}$$

$$P(\text{Chinese}|c) = \frac{1+1}{3+2}$$

similarly others.

$$\therefore P(\text{C}|c) = P(c) \times P(\text{Chinese}|c) \times P(\text{Japan}|c) \\ \times P(\text{Tokyo}|c) \\ \times (-P(\text{Macau}|c)) \times P(\text{Shanghai}|c)$$

multinomial

Naike Bayda question

training set	doc words	in c=chinese?
1	Chinese Beijing Chinese	yes
2	Chinese Chinese Shanghai	yes
3	Chinese manab	yes
4	Tokyo Japan Chinese	no
TEST	d 5 Chinese Chinese Tokyo Chinese Tokyo Japan	P

$$\rightarrow P(\text{Chinese}/c) = \frac{5+1}{8+6} = \frac{6}{14} = \frac{3}{7}$$

$$P(\text{Tokyo}/c) = \frac{0+1}{8+6} = \frac{1}{14}$$

$$P(c) = \frac{3}{4}$$

$$\therefore P(\text{Chinese}/\bar{c}) \propto \frac{3}{4} \times \left(\frac{3}{7}\right)^3 \times \frac{1}{14} \times \frac{1}{14} \approx 0.0003$$

$$\rightarrow P(\text{chinese}/\bar{c}) = \frac{1+1}{3+6} = \frac{2}{9}$$

$$P(\text{Tokyo}/\bar{c}) = \frac{1+1}{3+6} = \frac{2}{9}$$

$$P(\text{Japan}/\bar{c}) = \frac{1+1}{3+6} = \frac{2}{9}, P(\bar{c}) = \frac{1}{4}$$

$$P(\text{Chinese}/\bar{c}) \propto \frac{1}{4} \times \left(\frac{2}{9}\right)^3 \times \frac{2}{9} \times \frac{2}{9} \approx 0.0001 \therefore \bar{c} \text{ is more probable}$$

Laplacean

correction is for
smoothing.

Page No. _____

Date: _____

steps

From training corpus extract vocabulary

calculate req. $P(c_j)$ & $P(x_k | c_j)$ terms.

For each class c_j in C_{cls} :

• $\rightarrow C_{cls}, \leftarrow$ subset of C_{cls} for which the target class is c_j

$$P(c_j) = \frac{|C_{cls}|}{\text{total} |C_{cls}|} \quad \begin{matrix} \text{1 class of class } c_j \\ \text{total } |C_{cls}| \end{matrix}$$

$$P(x_k | c_j) = \frac{n_{jk}}{n_j + \alpha / \text{vocab}} \quad \begin{matrix} \text{in how many docs. } \\ \text{doc } k \text{ is in } \\ \text{for Laplacean } \\ \text{correct} \end{matrix}$$

where

$K = 1$ to all unique words in query

n_j = no of words in C_{cls}

n_{jk} = no of words in corpus.

positions = all word positions in current doc. which contain tokens found in vocab.

we calculate $\Rightarrow C_{NB} = \operatorname{argmax}_{c_j \in C} (P(c_j)) \times \prod_{i \in \text{positions}} P(x_i | c_j)$

($c_j \in C$)

ie positions.

for

class

having

max^m value.



Bernoulli Naive Bayes classifier

check absence for
all remaining words in vocabulary

- considers presence & absence of words (0 or 1)
- 2 possible values for a word \Rightarrow 0 or 1.

$$P(x_i | c_j) = \frac{N(c_j | x_i = 1, c_j) + 1}{N(c_j) + k}$$

↗ no. of time word appears
 in documents containing label C
 ↘ Laplace correction

↗ no. of values for award
 ↘ here 2 - 0 or 1

	class IP	words	in c = China?
training set	1	Chinese Beijing Chinese	yes
	2	Chinese Chinese Shanghai	yes
	3	Chinese Macao	yes
	4	Tokyo Japan Chinese	no.
test set	5	Chinese Chinese Chinese Tokyo Japan	?

$$P(\text{Chinese} | c) = \frac{\text{3+1}}{\text{3+2}} = \frac{4}{5}$$

↗ Chinese appears in 3 documents
 ↗ total documents (0,1)

$$P(\text{Tokyo} | c) = \frac{0+1}{3+2} = \frac{1}{5}$$

$$P(\text{Japan} | c) = \frac{0+1}{3+2} = \frac{1}{5}$$

$$P(\text{Beijing} | c) = \frac{1+1}{3+2} = \frac{2}{5}$$

$$P(\text{Macao} | c) = \frac{1+1}{3+2} = \frac{2}{5} = P(\text{Shanghai} | c)$$

$$P(c) = \frac{3}{5} \quad \therefore P(\text{China}) = \frac{3}{5} \times \frac{4}{5} \times \frac{1}{5} \times \left(1 - \frac{2}{5}\right) \times \left(1 - \frac{2}{5}\right) \times \left(1 - \frac{2}{5}\right) \approx 0.005$$

$$P(c \text{ Chinese}/\bar{c}) = \frac{1+1}{1+2} = \frac{2}{3}$$

$$P(c \text{ Tokyo}/\bar{c}) = \frac{1+1}{1+2} = \frac{2}{3} = P(c \text{ Japan}/\bar{c})$$

$$\begin{aligned} P(c \text{ Beijing}/\bar{c}) &= P(c \text{ Macao}/\bar{c}) = P(c \text{ Shanghai}/\bar{c}) \\ &= \frac{0+1}{1+2} = \frac{1}{3} \end{aligned}$$

$$P(\bar{c}) = \frac{1}{9}$$

$$\therefore P(\bar{c}/d_5) = \frac{1}{4} \times \frac{2}{3} \times \frac{2}{3} \times \frac{2}{3} \times \left(1 - \frac{1}{3}\right) \times \left(1 - \frac{1}{3}\right)$$

$$\approx 0.022$$

$\therefore \boxed{\bar{c}}$ is predicted for d_5

→ How to improve relevance through feedback

- Objective is to provide relevant docs to user
- We want to improve relevance of docs \Rightarrow goal
- Purchase implies high rating
 - if we click on retrieved result, like us
 - click 4th doc of out of 5 retrieved, then 1st 3 are not relevant, but 4 was.
- This is a kind of feedback
- We do query expansion.
- We want to provide better relevance to users.
- Improved relevance model, can be provided to same or other users, when same query was asked.
 - ↳ all this is query expansion.

eg - We ask "plane"

- results:- (synonyms as well as wrongly spelled ones are retrieved)

= flight

= aircraft

= by air

= aeroplane

= by air

= fly

= fight figt

= arcrft.

→ How to ensure good results

automatic

query
refinement

local

- we
use query
or refined
results
to reformulate
& improve
query

to for relevance

e.g. pseudorelevance

- indirect

relevance

feedback

- relevance
feedback

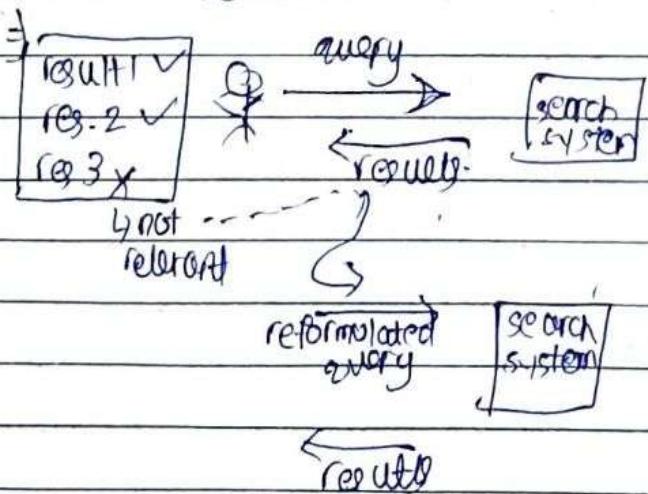
global

- we do not
use query
or results
to improve
& reformulate
query

for relevance

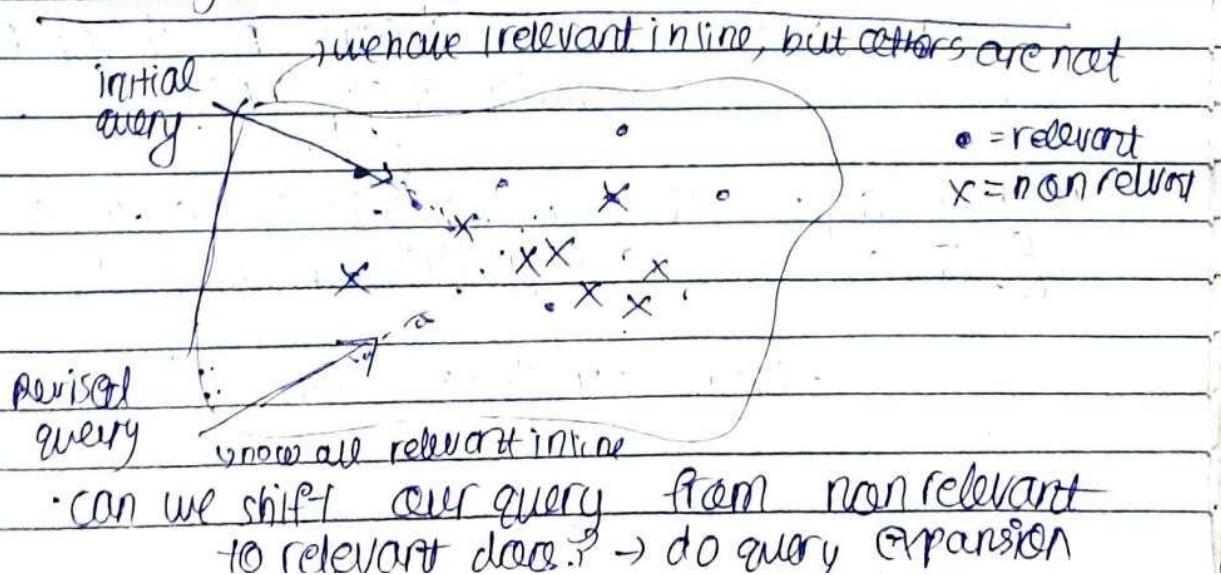
e.g. -thesaurus

- spelling
correction



- Refinement takes explicit feedback from user.
 - It asks if result is relevant or not.
 - But it's bad.
 - User doesn't have time & patience to provide explicit feedback.
 - But, still it is one of the ways.
 - * Indexed content is unknown to user
 - We don't know if other relevant docs. are present but not shown.
 - "Inform" need changes after looking at the results.
- eg - we search for a specific song. But looking at song results, we change to our mind & search other song.
- So need changes & it's necessary to do query expansion.

→ Rocchio algorithm (Relevance feedback (local))



- In practice, we only know few relevant & non-relevant.

\vec{q}_0 = initial query vector

D_r = set of relevant docs

$\beta \uparrow$ its weight

D_{nr} = set of non rel. docs

$\beta \downarrow$ its weight

$$\vec{q}_m = \alpha \vec{q}_0 + \beta \frac{1}{|D_r|} \sum_{d_j \in D_r} \vec{d}_j - \gamma \frac{1}{|D_{nr}|} \sum_{d_j \in D_{nr}} \vec{d}_j$$

α, β, γ are weights.

- We represent query as weighted vector.
- Use Recchio formula to get new query vector.
- Now calculate cosine similarities.
- Recchio model increases recall.
- +ve feedback is much more valuable than -ve feedback.

eg - initial query = " cheap CDs cheap DVDs extremely cheap CDs"

relevant d_1 = " CDs cheap Software cheap CDs"

not relevant d_2 = " cheap drives DVDS"

- Now construct term freq. vector.

	Cheap	CDS	+ -
q_0	3	-	
d_1	2	-	
d_2	1	-	

eg - to move $(2, 2)$ to $(5, 2)$, add 3 to x.

$$q_{\text{new}} = 1 \times q_0 + 0.75 \times d_1 - 0.25 \times d_2$$

do this for each term.

$$\begin{aligned} \text{eg. } q_{\text{new}} &= 1 \times 3 + 0.75 \times 2 - 0.25 \times 1 \\ &\quad (\downarrow) \\ &= 4.25 \end{aligned}$$

$$q_{\text{new}} = \boxed{4.25} \quad | \quad + \quad + \quad + \quad -$$

we can't have $q_{\text{new}} = -0.5$, this is wrong, so -ve values are put as 0.
 $\therefore q_{\text{new}} = 0$

unable
feedback.

Pseudo Relevance (IPnd) Relevance Feedback

- No user judgement
- we completely trust our web search engine
- we assume top k given by model are relevant to query
- now we can use ranking algorithm to change the query.

->

problem

- research : " Dhoni's elites from 1st class cricket".
- now, southafrica tests, etc. are also relevant hence retrieved.
- so matrix includes southafrica, newzealand etc.
- The vector query to vector will thus drift.
so, query drift problem occurs, as more weightage to terms which are not relevant.

→ Indirect (implicit) Relevance Feedback

- we do not ask explicit feedback.
 - we consider clickstream mining
userclicks on nth doc,
then top n-1 not
relevant,
if n is relevant
 - query expansion would ~~be~~ be done
based on implicit user preference (clicks)
- In global approaches, we do not consider either query or retrieved results.



→ Global (User) result - independent

automatic thesaurus gener'n.

eg - fast = rapid

fair = boundary \Rightarrow cricket

sound = noise

tall = height

- how to improve relevance?

- slangs?

\hookrightarrow brother \approx bro

. we can do co-occurrence analysis for the doc.

eg - mainframes
(doc)

mainframe
(doc?)

primarily, we
are ...

usually, we
refer

primarily & usually
are used like
synonyms.

eg - term incident matrx :-

words(book)	BOOK 1	B2	B3	B4	BOOK 1 & 3, BOOK 2 & 4 are similar.
w1	100	1	110	2	
w2	50	9	47	8	
w3	25	3	24	4	
w4	75	6	73	7	

- So, we can do co-occurrence analysis
- 2 terms are similar if term vectors are similar.
- Context is important
 - ↳ same video may not be relevant on weekend with but relevant on week.

eg - term-doc matrix

term-doc matrix • term doc matrix (Transpose)
= term-term matrix

from term-term matrix,
we can see results :-

	t_1	t_2	t_3	t_4	t_5	t_6
t_1	3	-	1	2	1	2
t_2	1		-			1
t_3	2					2
t_4	1					1
t_5	2					2
t_6	3	1	2	1	2	3

similar

so, $t_1 \approx t_6$

→ Precision & Recall

- Precision - fraction of retrieved docs. are relevant

$$P(\text{relevant}) / \text{retrieved}$$

- Recall - fraction of relevant docs that are retrieved

$$P(\text{retrieved}) / \text{relevant}$$

		true	
		relevant	non-relevant
retrieved	true	tp	fp
	false	fn	tn
not retrieved			

↳ false +ve ↳ false -ve ↳ confusion matrix

• Precision $P = \frac{tp}{tp+fp} \rightarrow$ correct prediction

($tp+fp$) → all retrieved docs.

• Recall = $\frac{tp}{tp+fn} \rightarrow$ correctly predicted

($tp+fn$) → total relevant instances

- There has to be balance b/w precision & recall

- If we try to increase one, other may decrease

eg - we retrieve more docs for more recall, but precision may decrease, & vice versa

- Diff. b/w recall & precision must not be more than 5%.

e.g - if precision = 95%, recall = 10%, then diff. $> 5\%$

- For recommender systems,

	Item 1	i2	i3	... in
User 1	-	-	-	-
User 2	-	-	-	-
User 3	-	-	-	-
⋮	-	-	-	-
User n	-	-	-	-

→ we try to predict each cell rating.

Then we calculate error by $\sqrt{\frac{1}{n} \sum (actual\ rating - predicted\ rating)^2}$ (Root mean square error)

This is our accuracy parameters.

eg - IRS retrieved 20 docs.

there are 100 relevant in corpus.
8 are correct & relevant



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

(similarly
20 cells)

$$\text{Recall} = \frac{8}{100}$$



Retrieved docs → RRNNNNNNN RNRNNNR
NNNNR

→ here
all the
Rare

correct
recommend.

$$\text{precision} = \frac{6}{20}$$

let 8 relevant docs (total)

$$\therefore \text{recall} = \frac{6}{8}$$

we use HM as we want both precision & recall to group

F OFF F measure to go up

→ F measure (F1 score)

- Takes both precision recall into account

$$F = \frac{2 \cdot P \cdot R}{P + R} = \frac{2}{\frac{1}{P} + \frac{1}{R}}$$

→ more conservative

Q - Why F measures takes harmonic mean?
(not arith. or geometric)?

A - arithmetic - biased towards high values.
(tail)

the geometric mean - ? - gives consistent result
for normalized values.

→ eg -

R	R	R	R
R	R	R	R

 case 1

↳ here ranking is not good

R	R	R	R
R	R	R	R

 case 2

↳ both have same precision &
recall, but case 2 is better,
as 1 to 8 are ranked,
and ^{case 2} correctly
retrieves
ranking wise.

→ solution - precision at kth index.



retrieved
1000

Page No. _____
Date: _____

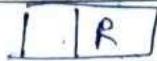
similar
if total
relevant docs
recall@1 = 100
recall@2 = 100
recall@3 = 100

$$P@1 = \frac{0}{1} = 0$$



← At 1st index / location

$$P@2 = \frac{1}{2}$$



← At 2nd index

$$P@3 = \frac{2}{3}$$



$$P@4 = \frac{2}{4}$$



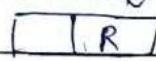
adding more & more
docs only relevant
values

disadv. - If only 4 total relevant docs,
will retrieve 10 docs.
max precn only 0.4

→ Interpolated precision

cut off at kth relevance level.

$$P@1 = \frac{1}{2}$$



↳ eg- if total
5 relevant,

1st locn where R appears

$$P@2 = \frac{2}{3}$$

2nd R

then 5
levels.

1st locn
= where
first R
appears,

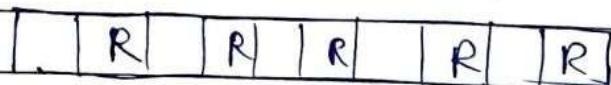
simil only
2nd locn
where 2nd
R appears

↳ if likewise.

$$\text{interpol. Avg. precn} = \frac{0.5 + 0.6}{2}$$

(if we only care
only interested in 2
levels of precn)

* case 1 → To calculate avg. precn



$$\text{avg. precn} = \frac{\frac{1}{2} + \frac{1}{2} + \frac{1}{2} + \frac{1}{2} + \frac{1}{2}}{5}$$

eg - AP at level 4:  \rightarrow level 3

Page No.
Date :

ex/lication

\therefore T- φ 10 does,

Index

$$\text{avg} = \frac{1+1}{2} + \dots + \frac{10+1}{2}$$

case 2 → for Avg.pressⁿ (AP) at level 4,

$$\frac{P@1 + P@2 + P@3 + P@4}{4}$$

If multiple queries are given we can calculate MAP (mean avg. precision) is important - calculates precision for all queries for all relevant levels

$$\therefore \underline{P@1 + P@2 + P@3 + P@4 + P@5} \quad \text{Qf reduce}$$

$$= \frac{1}{1} + \frac{2}{3} + \frac{3}{6} + \frac{4}{9} + \frac{5}{10}$$

5

$$= \frac{1 + 0.66 + 0.5 + 0.5 + 0.44}{5} = 0.62$$

query-2) 门田田田田田田田

$$\underline{P@1 + P@2 + P@3}$$

$$\geq \frac{\frac{1}{2} + \frac{2}{5} + \frac{3}{7}}{3} = \frac{0.5 + 0.4 + 0.428}{3} = 0.428$$

$$\therefore \text{final MAP} = \frac{0.62 + 0.44}{2} = [0.53]$$

(3) Case 1: P@1

(3) 2 ways (\Rightarrow)

\rightarrow Also,

$$\text{Precision} = \frac{|\{\text{relevant docs}\} \cap \{\text{retrieved docs}\}|}{|\{\text{retrieved docs}\}|}$$

$$\text{recall} = \frac{|\{\text{relevant docs}\} \cap \{\text{retrieved docs}\}|}{|\{\text{relevant docs}\}|}$$

(MAP)

\rightarrow Mean avg. precⁿ is the average of average precⁿ of all queries.

e.g. q1 = [R | R | R | R | R]

q2 = [| | R | | R | | R]

q3 = [| | R | | R | | R]

$$\text{AP for q1} = \frac{1 + \underbrace{\frac{1}{2} + \frac{1}{2}}_{\text{level 1}} + \underbrace{\frac{1}{2} + \frac{1}{2}}_{\text{level 2}} + \underbrace{\frac{1}{2} + \frac{1}{2}}_{\text{level 3}}}{5} = \frac{1}{2}$$

AP for q2: $P@1 + P@2 + P@3$

$$= \frac{1}{3} + \frac{2}{6} + \frac{3}{9} = \frac{1}{3}$$

$$\text{NP for } q_3 = \frac{1+2+3}{3} = \frac{6}{3} = 1$$

$\therefore \text{mNP} = \text{AP of } q_1 + \text{AP of } q_2 + \text{AP of } q_3$

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

$$\Rightarrow \therefore \text{MAP}(Q) = \frac{1}{|Q|} \sum_{j=1}^{|Q|} \frac{m_j}{m_j} \sum_{k=1}^{m_j} \text{Precision}(R_{jk})$$

* Problems - we cannot compute MAP if
 we do not know the total no. of
 relevant results for any query
 This also happens in webpage

- For recommended RMSE
- For webpage - & MAP (meaning. presn)



Crawling

⇒ Role of crawling / crawler

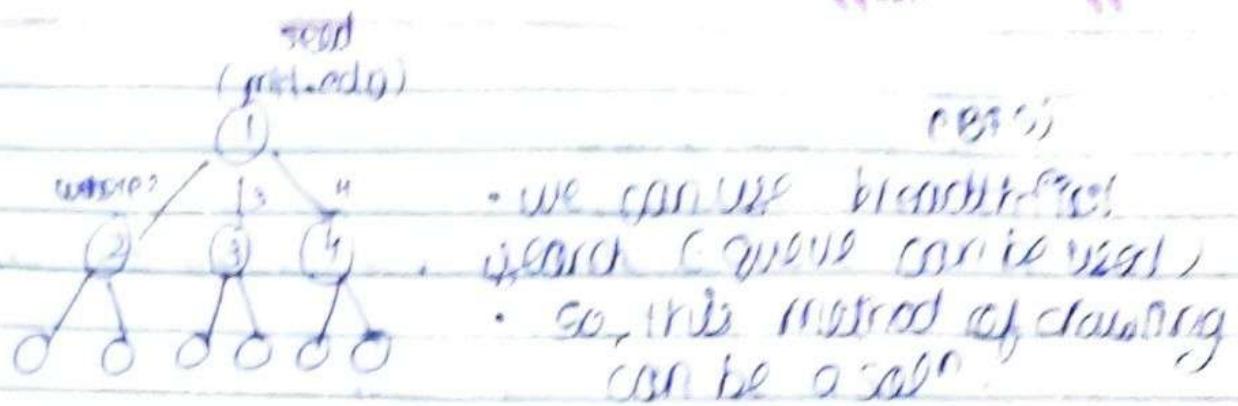
- we discuss IRS where corpus is available.
- But, we need a proper corpus
- So, crawler is responsible to create proper corpus.
- Every second, multiple web pages are added on internet
- So, corpus must be ~~up to~~ update.
- It is also called "spider" / "robot"

Basic crawler opern

- Begin with known "seed" URLs
 eg- cisco.com ← needed ↗ initial known pages
 MIT.edu ← authoritative to fetch other docs. ↗ some authoritative pages
 ↗ some websites we trust.
- Fetch & parse them
 - extract URL they point to
 - place the extract URLs on a queue
- Fetch each URL on the queue & repeat.

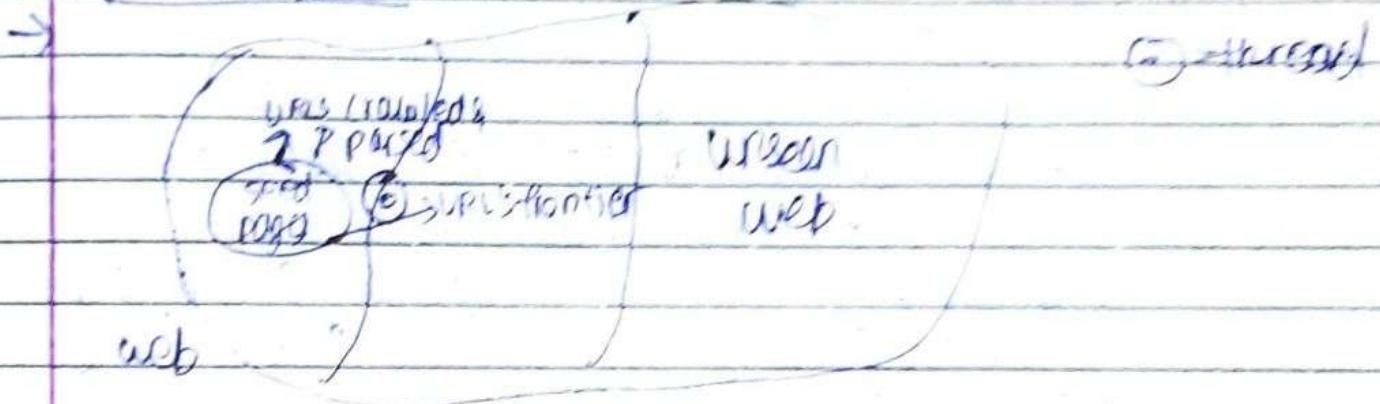
eg- we have MIT.edu. → seed URL

- If we have this URL, then crawler extracts links to other external websites (this is parsing).
- From those webpages, again parse & retrieve docs.
- We can use queue.



- So, crawler may or may not use depending on TPS
- But its req. for websites.

Crawler picture



URL frontier - set of URLs, not yet parsed

e.g. queue: 7773456

URL frontier

crawled
parsed

But not pollute, as 7 didn't

still 3 not parsed

use response
hitting PE

problematic aspects

- Not possible with 1 machine.
- distributed architecture needed
- ↳ more processing.

• malicious pages

- spam page

↳ title of page & content different

- spider trap trap

↳ start with seed, then next page
then again new page then

again bad page \rightarrow spam

↳ efforts of crawlers are wasted

• so, we need to deal with them.

• Even non malicious pages pose challenge -
↳ latency / bandwidth to remote servers vary
↳ eg - busy server.

↳ manager of website

↳ gives quick reply,
other busy

- webmaster stipul's

↳ constraints for how deep can
crawler crawl the content
of website,

↳ rules in robot.txt.

* - crawler needs to follow this
protocol ↳ rules (restrictions).

↳ crawler gets limited access.

↳ eg - only educational crawler can
crawl mit.edu.

↳ so, only some
pages allowed
to be crawled.

e.g - no robot should visit any URL starting with ". /yoursite /temp" except robot called "search engine"

• user-agent :- *

• disallow = "yoursite /temp"

• allow = "search engine"

- other problems:-

- site mirrors & duplicate pages

↳ visiting same/similar pages.

- politeness - crawler must not bombard the server, don't hit a server too often.

eg- if nirmaluni.ac.in doesn't give reply, wait for sometime then send.

→ what crawler must do

- Be polite, even if explicit politeness (in robot.txt) not mentioned - implicit politeness

- only crawl allowed pages.

- respect robot.txt.

- Be robust & handle spider traps.

- Be capable of distributed open, or multi-threaded environment.

- Be scalable, to ↑ crawl rate by adding more machines

all threads & machines
should be busy.

- Perf. / efficiency - permit full use
of available processing & network
resources

scalable -

webpage → more breadth,
of webpage,
even then crawler must
must not ↓ speed.
crawl rate

- Fetch higher quality pages first.
- continuous oper. - continue
fetching fresh copy of previously
fetched page
 - eg - news page should
be crawled more
often compared
to nirmauni.ac.in
(priority
of fetch rate)
- Extensive - adapt to new data formats /
protocols.

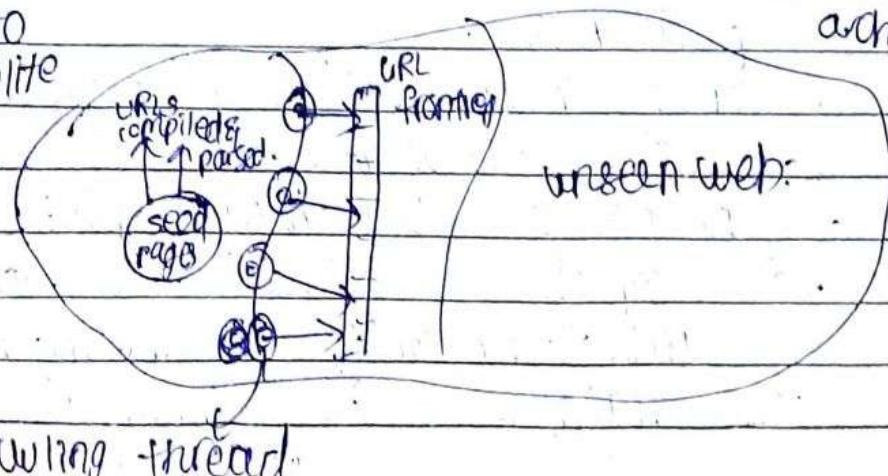
→ updated crawling picture

we have seen
discusses, so we
need multi-threaded
architecture

- We want to
prevent inactivity
ness

• all threads

must be busy



crawling thread

→ URL frontier

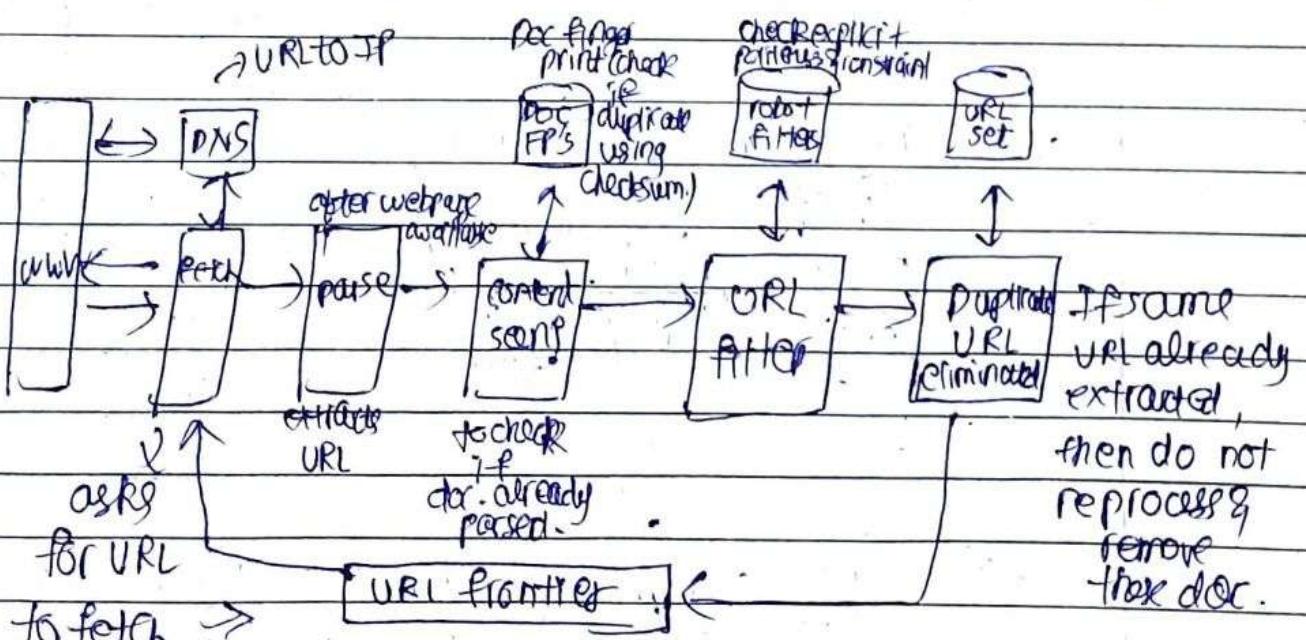
- set of URLs yet to be parsed
- can include multiple pages from same host.
- must avoid trying to fetch them at same time.

→ Implicit & explicit politeness

- Explicit → use robot.txt.
- Implicit → understand on own not to hit server too much.

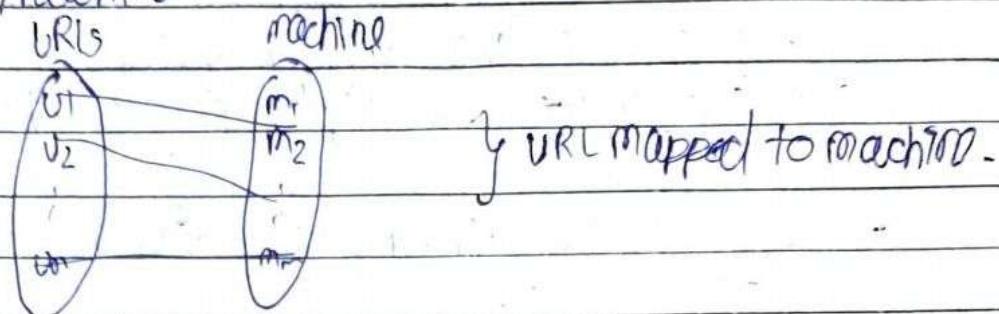
→ Basic crawler architecture

- If we have 1 machine, then all in 1 machine.
- If 20 machines, then architecture distributed & runs simult. on all machines.



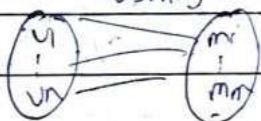
• Deals with explicit but not implicit politeness.

- we can use concept of hashing - ~~URL~~ URL-to-machine.



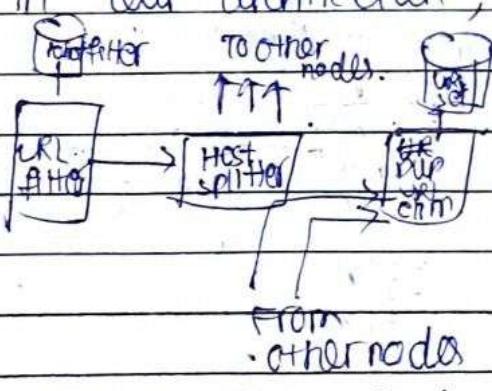
- Suppose the came ~~is~~^m comes across U_2 , but it's not responsible for it, then it forwards the request to the responsible machine.

- To achieve distributed archi., server must be divided so, parsing done by particular machine we use hashing to map URL's to machine hashing.



- Nodes are our machines which are geographically distributed

so, in our architecture,



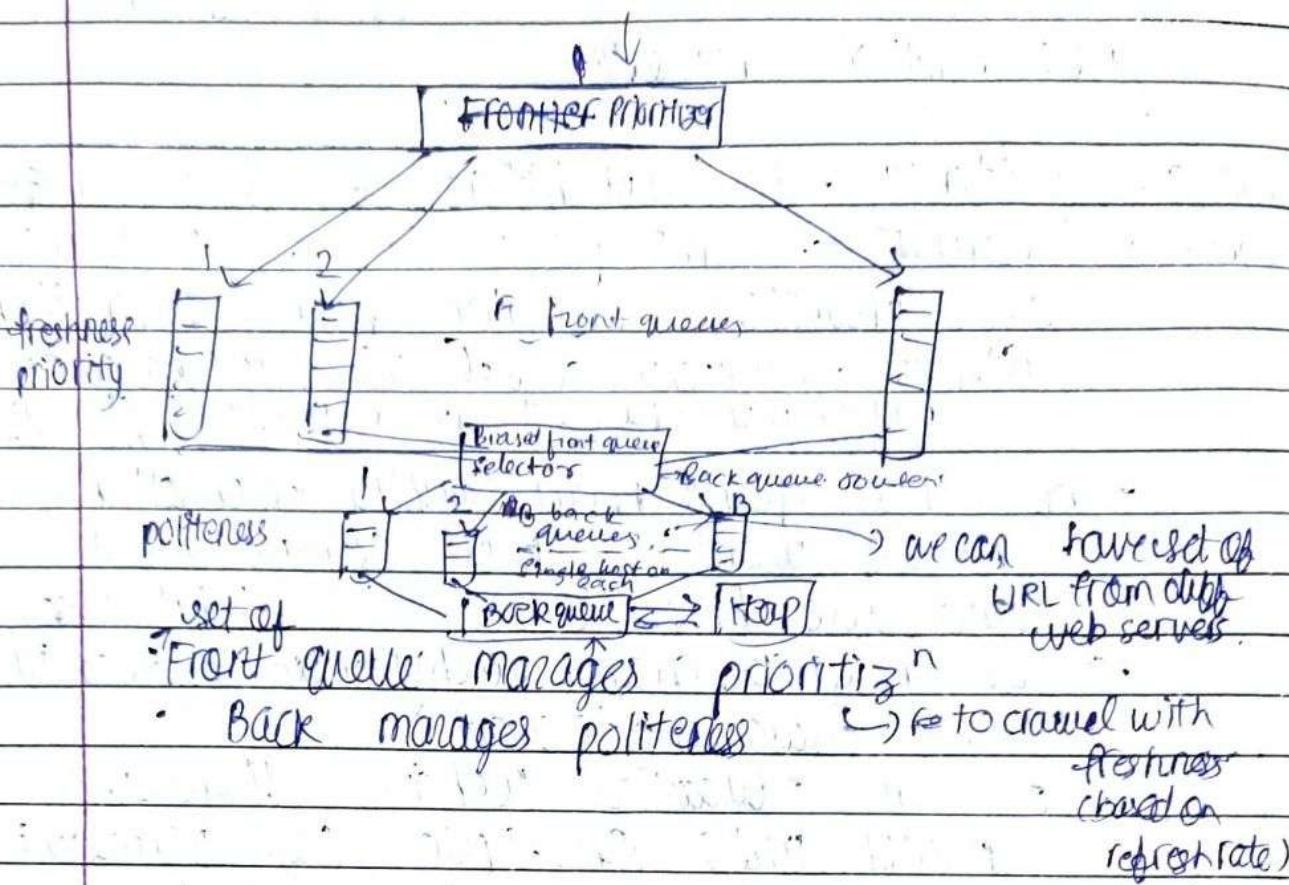
- We have a URL frontier for each node, which sends it to URL clim. of appropriate node

- If M1 comes to know U1 is not its responsibility, if forwards U1 to M2;

simple priority queue fails as - many links out of a page go to its own site creating a burst of access to that site
eg cisco.com may be links going to same 1110 server

URL frontier considerⁿs

- 2 min^m f. measured
- politeness - do not hit a server too frequently
 - freshness - some pages like news, news, are crawler crawled more often than others. ↳ we can crawl this page based on refresh rate of page. If p. refresh rate ↑, then crawl more often.
 - But these goals may conflict each other.
 - As we want to keep threads busy, and also be polite. ↳ hit servers.
 - so we need to change architecture of frontier.
 - (Simple queue won't work)
 - eg - on educational websites, there are many links on 1 page (eg. Nirmauni law, Engg, Pharm etc.)
↳ so crawler hits Nirmauni too often.
 - even if we restrict only 1 thread to fetch from a host, can hit it repeatedly.
 - heuristic - if server replies after 10ms, multiply it by 10, ∵ after 10ms send request again.
 - insert time gap bw successive req. from to a host \gg time for server for most recent reply time / fetch time
 - ↳ usually, we multiply by 10 times



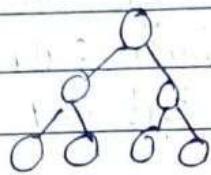
→ Mercator URL frontier

- URL flow from top to frontier.
 - prioritizer assign URL integer priority b/w 1 to k. (all queues are FIFO).
 - Then it assign priority to 'carries' queue
 - Heuristic to assign priority - refresh rate can be checked from previous "crawls" (e.g. news have high refresh rate)
 - i) we can select prioritized queue directly using its priority assigned or
 - ii) we give probability 0.5 0.75 generate random no. & choose
- [1 2 ...] → priority for URL to one of

Now for back queues:-

- For each URL we have a set of timestamp
- we use minheap for selecting minimum timestamp

e.g. - lets say timestamp is 10s, 30, we can send after that it is 20s
so, root node gets changed.



we req at 10:25, we get response at 10:35.
so, we delete 10:25 by some other (lets say 10:26), & the 10:35 is inserted in heap.

so, back queues help in politeness.

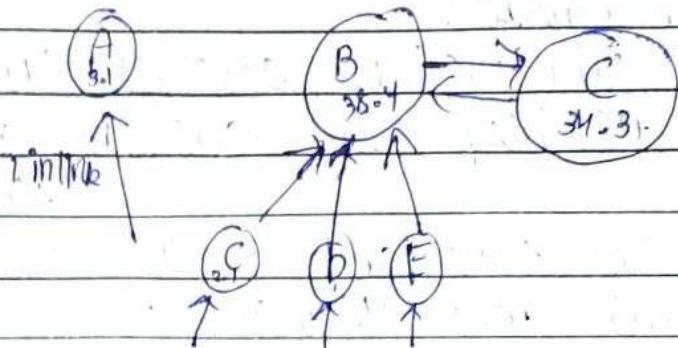
→ when a back queue req. a URL, ask front queue to give URLs.

→ How to check if page is authoritative?

use page rank - assigning rank to page.
we have - outlink - how many websites on our page
- inlink - how many websites reach our page.

more reliable
as it's
difficult to fake

eg -

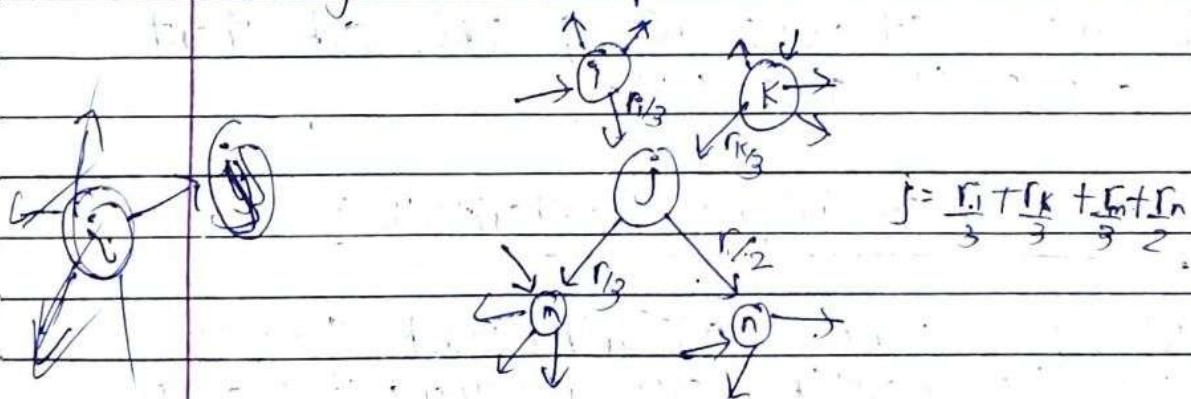


- So, B has many inlinks, so more also
- Also, C has high score, as B is authoritative, and gives inlink on C.

→ each link's vote is α to importance of its source page.

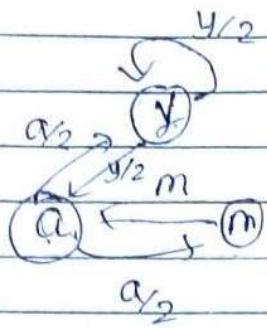
page
• If j has importance r_j , has n outlinks, each link gets $\frac{r_j}{n}$ rates

• j 's own importance is sum of its inlinks



• If page is vote from page is more important if it is pointed to by other pages.

rank for j : $r_j = \sum_{i \in d_j} \frac{r_i}{|d_i|}$ for each inlink



8/6/17

importance of $r_y = 2$ (2 inlinks)

$$\therefore r_y = \frac{r_u}{2} + \frac{r_a}{2} + \frac{r_m}{2}$$

only 5 outlinks
for inlinks. for y

+ no reln
from a
blw m

$$v_{ya} = \frac{v_u + v_m}{2}$$

$$v_m = \frac{v_a}{2}$$

no reln with y.

3 eqns, 3 unknowns. (no unique solⁿ, as no const.)
~~* additional constraint $\Rightarrow r_y + r_a + r_m = 1$.~~ scale factors
~~now unique const.~~

$$\text{sol}^n \Rightarrow r_y = \frac{2}{5}, \quad r_a = \frac{2}{5}, \quad r_m = \frac{1}{5}$$

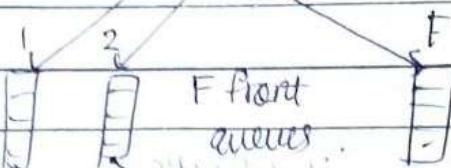
we can use Gaussian elimⁿ to solve,
 but for bigger values, we can
 use matrix formulⁿ.

Self

→ URL frontier

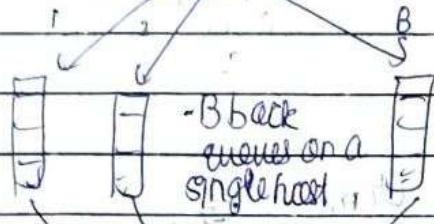
↓ i/p URL

prioritizer



pulled out of front queue!

B) biased front queue selector



pulled out of back queue

back queue selector

Heap

Sent to
crawler thread
then fetch

URL's

↓
prioritizer

↓↓↓↓↓
K front queues

↓↓↓↓↓
Biased front queue selector. Back queue ready

↓↓↓↓↓
B back queues
single host on each

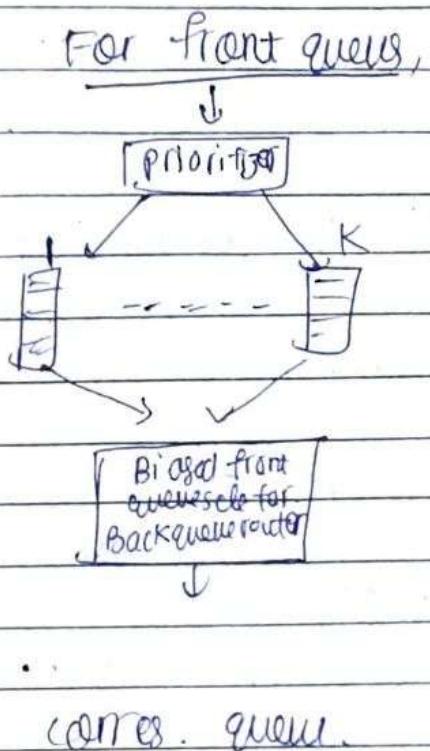
↓↓↓↓↓
Back queue selector

URL frontier

mercator store

crawl/thread req. URL - then fetch

- All this is mercator architecture
 - ↳ old classical crawler.
- Front queue prioritize based on refresh rate
 - (eg - news crawled more)



- There are K front queues, and URL incoming is assigned priority from 1 to K .
- Priority 1 need more freq. crawling than others (eg - news?)
- Prio. K need less freq. crawling
- So, we append URL to corres. queue.

- To assign priority, heuristic: - refresh rate sampled from previous crawls

For back queues

- How to pull URL's from front queues?
- We may pull in a round robin way (first queue F1, F2, .. FK). But it's wrong, as we need to prioritize lower no. of queues
- So, we need more URLs from queue 1 than queue K .



There are 2 ways to pull from front queues:

- 1) \rightarrow First pull an item from queue 1.
Then pull next time, pull an item from queue 2,
Then 1, 2, 3

1

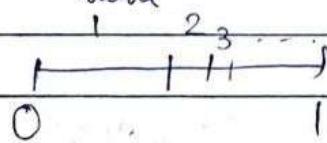
2

3

- we can see than 1's URL will be pulled more often, than 2, then 3 ...
so, it's kind of round robin.

- 2) \rightarrow we assign prob. to each queue

- eg - if random no. bw 0 to 0.5,
pick URL from queue 1
If 0.5 to 0.75, pick from queue 2, so on



so, 1 has larger prob. of queue being pulled & thus algo. is biased for more priority queue.

- \rightarrow So we have solved freshness problem by front queues

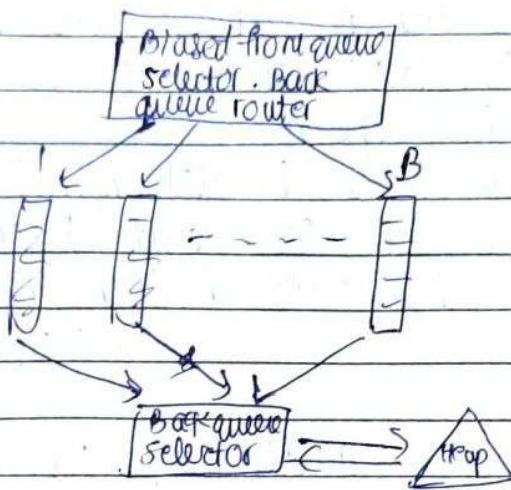
- URL going in are in any order
but coming out based on priority.

\rightarrow

- Now, for the back queue,
- URL pulled out. will go in a back queue from 1 to B.
- In a front queue, it is possible for
 - to have URL's of diff servers

e.g. CNN.com & BBC.com

have same priority, & so their links may be on same front queue.



- each back queue has links only from 1 server
- e.g. queue 1 has only links from CISCO server
- queue 2 has only from Stanford server, etc.

- we maintain a min-heap.
- heap has B vertices, each vertex corresponds to each back queue (or each server)
- The value of the key: threshold time
 - ↳ time before which you can't hit server.

vertex

e.g -

(10:15 PM
cisco)

→ we can't send query to cisco before 10:15 PM.

- so, root has least value of timestamp
- so, first we send query to root URL.

- So, we extract ~~HTTP~~ root URL
- After that, suppose removed URL was of queue 2, we go to back queue give it to fetcher.
- So, in back queue
 - one entry for each back queue
 - ~~The~~ The entry at earliest time to at which the host comes to back queue can be hit again.
 - The earliest time is determined from \Rightarrow last access to that host.
 \Rightarrow any other time buffer heuristic we choose.
- * Suppose we pulled all URL's from back queue, & it is empty (server of cisco)
 - We ask front queue to give a URL:
 - If it is stanford website, we put it in Stanford queue.
 - We again ask for URL, till we get atleast 1 URL of cisco, so that atleast 1 URL in back queue.
- Now there are millions of servers, so we require millions of back queues, which is not feasible.
- So, once a queue gets empty, we relabel it from cisco to mit.edu.
 - ↳ some other server queue.

- Back queue processing:-
 - A crawl thread seeking to crawl
 - Extract root of the heap
 - Fetches URL at head of corrs. back queue
- []
- ↳ This URL
- Checks if queue q is now empty, if so, pulls a URL v from front queue
 - if there's already a back queue for v 's host, append v to q & pull another URL from front queue, repeat
 - else add v to q
 - When q is nonempty create a heap entry for it
 - so when root removed we again take new URL from the corrs. back queue, calculate new timestamp & reinsert it in heap (now it could be anywhere in the heap)
 - we maintain a table from hosts to back queues

Host name	Back queue
cisco	3
mit	2
?	1

- each back queue is kept non empty while crawl is in progress.
- each back queue has URL only from 1 host.

number of back queues B :-

- Keep all threads busy while respecting politeness.
- mercator heuristic - 3 times as many back queues as crawler threads.

class :-

For bigger graphs, we can represent using adjacency matrix.

e.g. if $i \rightarrow j$ then $M_{ji} = \frac{1}{d_i}$ else $M_{ji} = 0$
 \downarrow out degree

matrix representation.

M is column stochastic

\hookrightarrow sum of all values in col. = 1.

Rank vector $r \rightarrow$ vector with an entry per page
 $\sum r_i = 1$ (r_i is importance of i th page).
 $r = M \cdot r$.

\rightarrow now, $r = M \cdot r$ can be mapped like eigen vector.

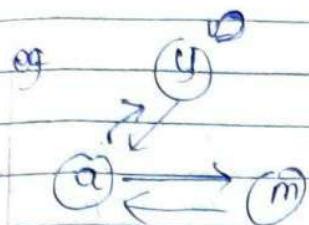
Suppose.

$$I \cdot r = M \cdot r$$



similar to $Ax = \lambda x$

$\therefore r$ is eigenvector of Matrix M .



$$\begin{array}{c|ccc}
y & y & a & m \\
\hline
y & \frac{1}{2} & \frac{1}{2} & 0 \\
a & \frac{1}{2} & 0 & 1 \\
m & 0 & \frac{1}{2} & 0
\end{array} \rightarrow \text{matrix } M$$

(for link $i \rightarrow j$, we fill M_{ji})

↳ basically
we fill invertedly

• we need to assume initial r vector.

$$\text{set } r_i = \frac{1}{N} \quad \begin{matrix} \text{no. of nodes} \\ (\text{here } N=3) \end{matrix}$$

→ power iterⁿ method

initial
set $\Rightarrow r = \frac{1}{N}$

$$\therefore r = \left[\begin{array}{c} \frac{1}{3} \\ \frac{1}{3} \\ \frac{1}{3} \end{array} \right] \quad \begin{matrix} \text{3 nodes so } (3 \times 1) \\ \text{dim} \end{matrix}$$

next iterⁿ $\rightarrow r = M \cdot r$

• keep doing this this r and r of next iterⁿ become same.
(convergence).

$$\therefore \text{eg - } r = \begin{bmatrix} \frac{1}{2} & \frac{1}{2} & 0 \\ \frac{1}{2} & 0 & 1 \\ 0 & \frac{1}{2} & 0 \end{bmatrix} \cdot \begin{bmatrix} \frac{1}{3} \\ \frac{1}{3} \\ \frac{1}{3} \end{bmatrix}$$

$$= \begin{bmatrix} \frac{1}{3} \\ \frac{1}{2} \\ \frac{1}{6} \end{bmatrix} \quad \text{new } r.$$

- now again

$$r = M \cdot r$$

- Keep doing this till new value of r & previous value of r almost same (or exactly within some specified threshold)

let final $r = \begin{bmatrix} 6/15 \\ 6/15 \\ 3/15 \end{bmatrix}$ score of page Y
 after convergence. $\begin{bmatrix} 6/15 \\ 6/15 \\ 3/15 \end{bmatrix}$ page a page m

These r give our score/rank

so we solved a big problem using power iter. without gaussian elimin. or flow model.

using Hubs &
authority

- we assumed that only 1 score is assigned to a page.
- we use hubs & authority algo. for this.

eg - we have some newspapers
 \downarrow
 hubpage

- these papers have ~~at~~ refer other good papers \hookrightarrow authoritative page.
- hubs give outlinks to auth. pages.
- auth. page have inlinks., & useful inform

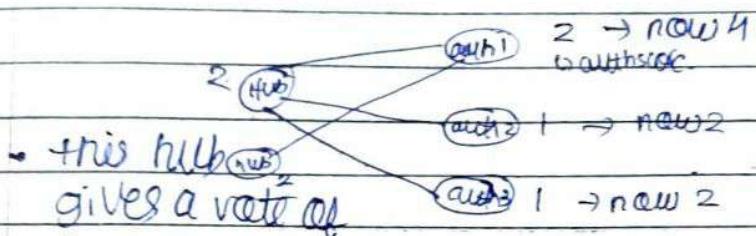
2 ways to check quality

- quality as an expert \rightarrow Hub
 - quality as a content provider \rightarrow authority
- To know good is content.

- ⇒ Auth. page have useful info. (usually).
- Hubs to have list of all auth. pages
(eg - newspapers)

- To count q score,

- If there is a link from hub to auth. page, then we count it as a unit.

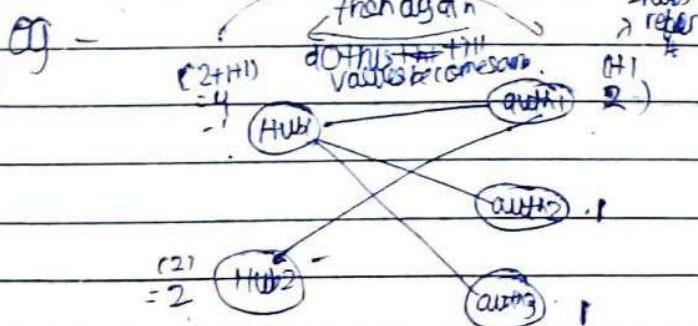


- This hub gives a vote of 1 to each auth. page.
- auth. page 1 has score of 2. (2 hublinks).
- Now using score of auth. page we can calculate hubscore.
- We add the sum of score of all auth. pages being referenced by a hub, this gives the hubscore.
- We keep calculating scores (hub & auth.) till convergence.

⇒ This HITS algorithm

Hyper

now again change it



IHS (contd.)

- Every page has 2 score \rightarrow hub score & authority score
- A good hub links to many good auth. & good auth. linked to good hubs.
- We have self reinforcing recursive fund.
- Represented as h & a

$\begin{matrix} h \\ a \end{matrix}$ $\begin{matrix} \text{hub} \\ \text{score} \end{matrix}$ $\begin{matrix} \text{auth} \\ \text{score} \end{matrix}$

\rightarrow HITS - hypertext induced Topic Selectn

$$\text{Initialize } \Rightarrow q_j^{(0)} = \frac{1}{\sqrt{n}}, b_j^{(0)} = \frac{1}{\sqrt{n}}$$

• keep iterating till convergence

$$\forall i \text{ auth} = q_i^{(t+1)} = \sum_{j \ni i} h_j^{(t)}$$

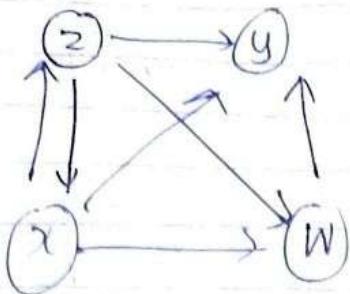
$$\forall i \text{ hub} = h_i^{(t+1)} = \sum_{i \ni j} q_j^{(t)}$$

$\forall i$: normalize:-

$$\sum_i (q_i^{(t+1)})^2 = 1 \Leftrightarrow (\sum_i h_i^{(t+1)})^2 = 1$$



eg -



$$H(z) = H(y) = H(x) = H(w) = 1$$

$$\alpha(z) = H(x) = 1$$

$$\alpha(y) = H(z) + H(x) + H(w) = 3$$

$$\alpha(x) = H(z) = 1$$

$$\alpha(w) = H(x) + H(z) = 2$$

$$\text{total } \alpha = (3+1+2) = 7$$

$$\therefore \text{normalized: } \alpha(z) = \frac{1}{7}, \alpha(y) = \frac{3}{7},$$

$$\alpha(x) = \frac{1}{7}, \alpha(w) = \frac{2}{7}.$$

$$\text{now, } H(z) = \alpha(x) + \alpha(w) + \alpha(y) \quad (3 \text{ outlinks})$$

$$= \frac{1}{7} + \frac{2}{7} + \frac{3}{7} = \frac{6}{7} = 0.858$$

$$H(y) = 0 \rightarrow \text{no outlinks}$$

$$H(x) = \alpha(z) + \alpha(y) + \alpha(w)$$

$$= 0.858$$

$$H(w) = \alpha(y) = \frac{3}{7} = 0.429$$

$$\text{Normalized total } H = H(z) + H(y) + H(x) + H(w) \\ = 2.145$$

$$\therefore \text{normalized: } H(z) = \frac{0.858}{2.145}$$

$$H(y) = 0$$

$$H(x) = \frac{0.858}{2.145}$$

$$H(w) = \frac{0.429}{2.145}$$

- \Rightarrow we repeat this, till values become nearly same. (convergence)
- This is to make good hub & auth. pages.
 - This is HTTS algorithm.

→ Need of Recommender system

- If user knows what he wants, then he directly searches.
- He can also look & choose
- But what if more choices? \rightarrow more confusion.
- so, to deal with this, we need recommender system

\hookrightarrow type of IRS

- Provides personalized service to user.

\Rightarrow Earlier we had retail store
 \hookrightarrow limited shelf space.

also, TV \rightarrow only 1 channel at a time
 movie \rightarrow one at a time.

\Rightarrow so, we have limited choices

⇒ new web enables near 0 cost dissemination of info about products
from scarcity to abundance.

⇒ so, many choices require better filter
recommendation engine.

e.g. one book of author became famous due to other book which became famous.

as book A becomes famous, book B of same author gets recommended,
so, B become more famous even though it was published earlier, but got famous now.

→ Non personalized recommender system

- editorial & hand curated
 - list of favourites
 - list of essential items

↳ e.g. recomm.
restaurants based on avg. rating

- simple aggregates

- top 10, most popular, recent buys ↗ take avg. ratings.

- Tailored to indiv. users
 - Amazon, netflix etc.

Formal model

let C = set of customers

let I = set of items.

VB: $C \times S \rightarrow R$

utility matrix $R = \text{set of ratings, } R \text{ is totally ordered set}$
 $\text{eg - 0 to 5 stars, 0 means bad, 5 is best.}$ Ordered

eg - utility matrix \rightarrow user v/s item

		Moviel	2	3	4
		1			
person 1	1		0.2		
2		0.5		0.3	
3	0.2		1		
4				0.4	

- A person may or may not rate a movie
- so matrix can be sparse.
- we predict the ratings which are missing, then we recommend.

key problems

- 1) Dataset collection \leftarrow gathering known ratings.
 - 2) Extrapolate Unknown from known ratings.
 \hookrightarrow we are interested here
 - 3) evaluating extrapolation method \hookrightarrow how to measure success / failing
- * - we are not interested in what we don't like
only interested in what we like - may or may not be true.

1) gathering ratings

- explicit

- ask people to rate

- many users don't like this

- implicit

- if user purchases, then high rating

- learn from user actions.

- what about low rating? → issue.

2) extrapolating utilities

- key problem :-

- most people don't rate

- so sparse matrix (U matrix)

- cold start problem

- ↳ new items have no rating

- ↳ new users have no history.

⇒ 3 approaches to recommender system → to find value:-

↪ 1) content based

2) collaborative

3) Latent Factor based.

↳ similar to SVD, matrix factoriz.

wrks
based on
similar
attribut.

content based

→ we watch a movie, our profile is prepared based on attributes of movie. If new movie's feature match with our matched movie, recommend new movie.
 - problem: - overspecializ?
 only similar movies recommended.

collaborative filtering

- suppose user x
- set of N other users have similar choices to x
- we check ratings of ~~N users~~ ~~x user~~ N users to predict based on pred ratings of ~~N other users~~ to predict rating of x user.

eg- set $N = \{3, 4\}$ for x , if
 N_1, N_2 rating can be $= \frac{3+4}{2}$

$= 3.5$,
 given x is similar to
 N users. (here $N=2$)

- to find if x is similar, we use different techniques.



	movie1	m2	m3	m4	m5	m6
A	4		0	5	1	
B	5	5	4			
C			2	4	5	
D	3			5		3

- consider users X & Y with rating vectors r_X & r_Y
- we see, that A & B are not similar initially.

eq. similarity $(A, B) = |r_A \cap r_B| / |r_A \cup r_B|$

Jaccard. \rightarrow common $\rightarrow m_1$

$$= \frac{1}{5} \quad \text{sim}(A, C) = \frac{2}{4} \rightarrow \text{common } \rightarrow m_1, m_2$$

m_1, m_2, m_3, m_4, m_5 \rightarrow total $4 + m_1, m_3, m_4, m_5$

\downarrow total in $A \& B$.

- Jaccard only uses common items.
- so, it's a problem. \rightarrow

\Rightarrow so, we can use cosine similarity.

$$\text{sim}(A, B) = \cos(r_A, r_B) \rightarrow (5, 4, 0, 0, 0)$$

\rightarrow considering missing values as 0.

$$\text{sim}(A, B) = 0.38, \quad \text{sim}(A, C) = 0.32$$

$\text{sim}(A, B) > \text{sim}(A, C)$, but not by much.

• Now, 0.38 & 0.32 are nearly same problem - this happened as we considered missing values as 0 \rightarrow bad rating - just because we don't rate doesn't mean we didn't like it.

- Instead, we can replace missing ratings by average rating.
- This is called centered cosine similarity.