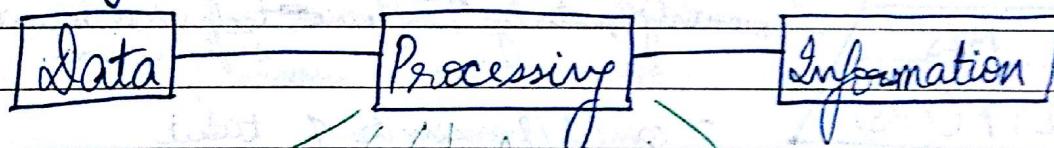


Information Retrieval



- 5 STEPS
- Collection
 - Classification - ① structural data (well defined attributes), ② semi-structured, ③ unstructured.
 - Coding - item code, bar code.
 - Sorting
 - Solving
 - Validation - not entering incorrect data. - restriction of pass.
 - Verification - password for login
 - Calculations
 - Storage - 2ndary storage
 - Retrieval

- 20 IR
- query suggestion - suggestion for 2nd word on typing one word
 - unstructured data
 - not well defined features
 - 'document' fetched from millions of sources
 - maintain in-cash & tries to give in those terms.
 - document ranking (Page rank)

25 User

Technical / Statistical

Non-technical.

30 OLT - Online Transaction System.

Top
level
Management

Middle



- suggestion in decision making & user achieves his/her goal
- assisting some contrib to decision taking - decision support
- monthly / quarterly / weekly decisions - technical decision based on info received.
- forecasting, targets, summaries, geoplots
- cancel / booking etc. of tickets

- launching of new product - static decision.

- weather forecasting - semi-structured data

↳ prev yr data + this yr data

Types of IR :

- ① Boolean Retrieval - True / False - Binary decision making.

- ② Vector Space IR model - frequency - basic unit of VSIR

e.g. Good post / kick

Given term

↳ suggesting based on frequency

- Web Crawling)

IR's goal = deliver most relevant info. - yet to be achieved.
 (objective) (exact)
 - query is not structured.
 - context dep. (T, semi) unstructured.
 (Search "apple")

- ① most q. are not 100% accurate: based on stats & probab.
- have to anticipate many things like context, why user writing this query,
- for same info., user can be diff. - diff. query answers (personalized)

* Uni. of Illinois - speech - summarization

* word WordNet

* Opinion Mining - analyse what ppl giving opinion want to say

* Sentiment Analysis

- blog via Twitter
- short vs long text analysis

What is IR?

- process of actively seeking out relevant info.

Document: web page, img, video clip, text → called doc in IR
 content

IR in practice

→ search - driven theoretical & experimental discipline.
 theories, proto, formulas, empirical research.

data

→ focus: context-dependent

→ query must be fast
 utilize space for text.

- 2) librarians - indexing by categorizing
 - e.g. one type → alphabetical indexing based on keywords
 - this where IR originated. - initially for books
- 3) cognitive sci. - sentimental analysis.
- 4) philosophy - relevance

#

computational linguistics

- grammar can be converted to CL
- 1. syntactical analysis
- 2. semantical analysis - 'not' has no significance.
 - find relation b/w keywords

#

- Wordcloud - diff colour, diff font, diff orientation



Info. visualization

- user feels comfortable to use
- retrieval is not boolean. → fuzzy

All IR are fuzzy in nature.

- Most relevant words are present.

How does it do this?

query - collection of words.

- even sequence of words (syntax & grammar) is not v. imp. (query suggestion)
- semantically imp keywords

Expanding a query → some related terms can be found and added in query. internally

- 1) • can generate a topic based thesaurus.
e.g. bomb entered by user., pick one/more of the
- 2) • generally query length \neq 25-30
• lexical/stat. analysis/content/concept - reads 100s of docs.
e.g. only goal of 8 after reading 100s of docs.
on football - now goal fast.
- 3) • relevance & feedback - in the beginning of Google, there was no text on search page, but now by seeing it / reading it, we know if it is relevant to us.
• Page ranking - dynamically changes.

Indexing and matching model 1 - all docs on thread, so how to make a mechanism to extract some documents not all.

- (Q) How to create index/sort as these are text docs?
- (A) - Keywords band hashing. (rest are common words)
 - extract keywords from docs.

Role of UI in IR

We use google search engine even though there are other search engines. (Q) Why? (A) Google provides excellent UI.

- not always able to define query well - prob. def.

Faceted organization

- # for some datasets we can do identify some features
- One of IR tasks is extracting features
- ① # medical database
 - On what features, the doc. should be extracted.
 - Doc. should be organized on the basis of some features
 - The doc. is composed of ... then ... features
 - ↳ composition (design term)

Structuring a doc. collection

- 1) Supervised - identify features & classify docs on features
- 2) Unsupervised (clustering) - features are not known.
(keywords ~~as~~ work as good features here)

Basic IR Models

1) Boolean models

(In DBMS - using AND, OR)

- either there is a result or nothing.

2) Vector space model

- each doc. is expressed in terms of keywords

$$d_1 = \{dt_1, dt_2, dt_3, \dots, dt_m\}$$

$$d_2 = \{dt_5, dt_6, \dots, dt_m\}$$

$$d_3 = \{ \dots \}$$

attribute value

Binary

present - 1
not - 0

Frequency

no. of times it's present

Weights

weight (imp) per word
(award one word & penalize another)

3) Cosine similarity

1 or 0 → doesn't belong.
↳ belongs to context

- semantically related.

- 4) Tfidf score • based on score, arranged/displayed decreasing order.
(most to least rel.)
• set a threshold. (only above score of '0.5')

10

15

20

25

30

Tokenization & Indexing

Quality measures

A- How to measure relevance?

A- Recall & Precision

$$\begin{aligned}
 &\text{# no. of doc} = 1000 \\
 &\text{# relevant doc} = 350 \\
 &\text{# retrieved } n = 500 \\
 &\text{# rel. } n = 200
 \end{aligned}$$

$$\text{Precision} = \frac{\text{no. of rel. ret. doc.}}{\text{no. of ret. doc.}} = \frac{200}{500}$$

$$\text{Recall} = \frac{\text{no. of rel. ret. doc.}}{\text{no. of rel. doc.}} = \frac{200}{350}$$

Goal High precision + High recall

$$\begin{array}{lcl}
 \text{F-Measure} & = & \frac{2PR}{P+R} \\
 (\text{NM of PR}) & &
 \end{array}$$

P↑ R↓ & R↑ P↓



Need for tuning of para.

Every doc is set of tokens.

convert
to
HTML
version
2.0

convert words to their roots

Lemmatization

e.g. good, better, best → good

• 5 Stemming (← Lemmatization)

e.g. play, plays, played, playing → play
- reduce terms to their roots before indexing.

• 10 Example from book — Shakespeare.

Tokenization

- 15 separate out the words
- Not all words are taken as tokens. keywords

Normalization

Match U.S.A. & USA

- 20 Normalize terms in indexed text

Case Folding

- 25 reduce all letters to lowercase.

Stop words

- 30 Most common words -

latest ~~not~~ Eng. stopword → 720

- if not performing syntactical analysis can remove them
- Obj - to remove reduce dimension.
- e.g. 'King of Denmark' L_{Needed} $\text{S}_{\text{--}}$

Indexing

- **Inverted index** = word \rightarrow doc containing 'word'
- flexible & fast

Best rep. of concept :-

- 1 word - coverage good, not precise
- 2 phrase - poor coverage, more precise
- 3 concept - set of rel. words | good coverage, precise

For each term T , we must store a list of all doc. that contain T . (also keep freq.)

• Organized in sorted order.

② • Post list - updated everytime

• Term & Doc. Freq.

~~A~~
Assignment - 1



17 Aug 2018
Deadline

docs = 500

5



find unique terms in all docs



10

index terms

freq. of each term

(Hash Map)



index terms

15

"Cryp's Law"



20

Come up with ideas for indexing

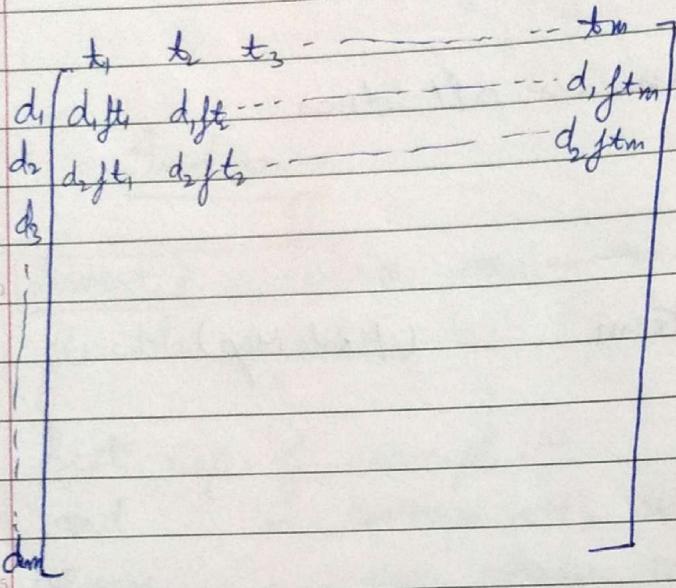
(TRAIN)

25

30

Similarity Measure

(Metrics)



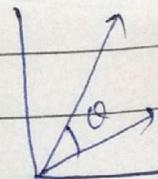
1) Euclidean Distance

k-means → Uses Euclidean distance

$$\sqrt{\sum_{i=1}^m (d_2 f t_i - d_1 f t_i)^2}$$

↳ sum w/w doc1 & doc2

2)



3) Jaccard Coefficient

indep. of freq. — if occurs! 1
else: 0.

$$J(d_1, d_2) = \frac{|d_1 \cap d_2|}{|d_1| + |d_2| - |d_1 \cap d_2|} \quad dm_1 \quad dm_2$$

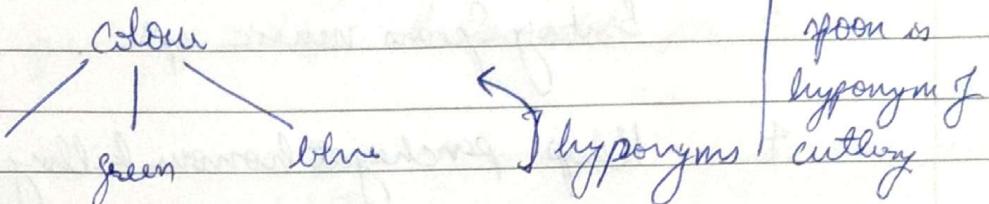
$$\left(\frac{|d_1 \cap d_2|}{|d_1| \cup |d_2|} \right)$$

Semantic Relatedness

* Indicated degree to which words are associated via any type:-

- ① synonym
- ② antonym
- ③ hyponym

e.g.



- ④ hypernym

e.g. colour is hypernym of red.

- ⑤ meronym.

CSE is meronym of egg. inst:-
part of

- 4) F-Score

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

P = Precision

R = Recall.

- 5) Pearson Correlation Coeff.

-1 to 1. ← correlation
Measure of linear dependence b/w 2 vars.

- 6) Con't

Assignment - 2

query 1 Anna Hazare Anti land Acquisition Bill

2 Stock market mutual fund

3 Britney Spears music mp3

4 Khap panchayat honour killing

5 Sql server dbms database

Q1 For the above queries fetch the top 100 doc. retrieved from the Google.

Q2 Calc. the various similarity measures and analyze the relevancy of a particular measure.

(all)

freq mat. for all except Jaccard \rightarrow find if there/ not.

Masquerade User Detection

Masquerade - inside/ outside the system

- disguise.
- attacker pretends to be auth. user of sys.
in order to gain access to it / to gain
greater privileges than auth. for.

13/8/18 10 C, 835 → Commands in Linux: 635

d, [d, f, d, f, d, f, 835]

du

: find freq. of each command for each user

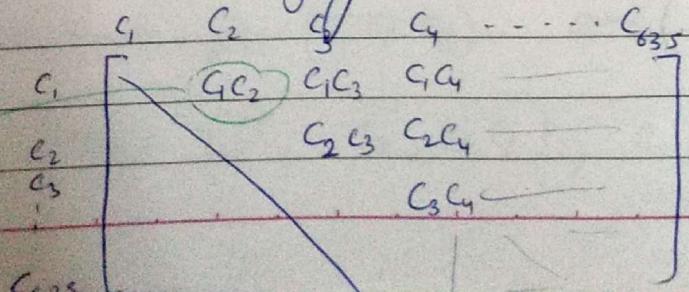
→ Sparse matrix

20 d₅₀ d_{50f} d_{50fss}

user - identification signature is the commands used everyday based
on job.

Hypothesis: deviating from signature commands then that user
is masquerading.

Objective: find similarity b/w commands using some similarity
matrix. → Sparse + each user uses set of related and
everyday = related (order not imp.) (based on user)



Now apply any similarity matrix.

20

 d_{so} d_{soft}

user - identification signature is the commands used everyday based on job.

25

Hypothesis: deviating from signature commands then that user is miscreating.

Objective: find similarity b/w commands using some similarity matrix - sparse + each user uses set of ~~related~~ and everyday = related (order not imp.) (based on user)

30

similarity
between c_1 & c_2
based on
Pearson metric

	c_1	c_2	c_3	c_4	c_{635}
c_1		c_1c_2	c_1c_3	c_1c_4		
c_2			c_2c_3	c_2c_4		
c_3				c_3c_4		
c_{635}						

Now apply any similarity matrix.

Now, we want to gen. cond pattern which is signature of a user.

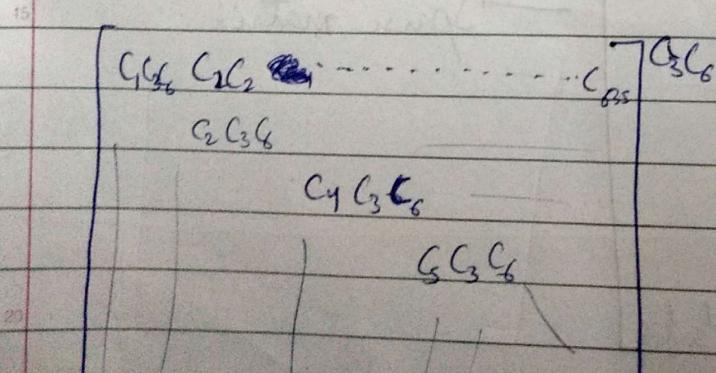
- d- How many cond are to be added ~~to~~ considered as signature)
- A- Experiments

Max. simi. → Related pair of cond.

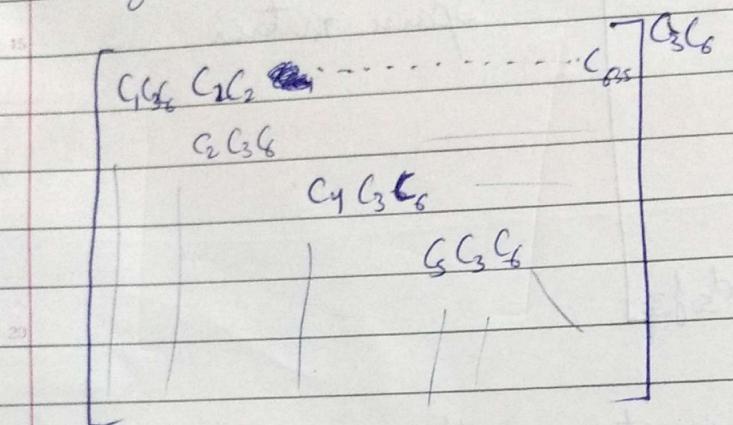
e.g. $C_3 C_6$. — max. simi.

Now treat $C_3 C_6$ as one new cond in place of C_3, C_6 indi.

- 1st pair that is most simi is removed from mat. & new col. is added as a pair.
- Next find ^{NEW} simi. tetra matrix. with a threshold (say 500).
- Again a max. value will come.



- Next find ^{NEW} simi. before maxm. (say 500).
- Again a max. value will come.



\uparrow Hit rate: Ability to declare magnified user as magnified user

\downarrow Miss rate: " " " " " anti-user

\downarrow False alarm rate: Declares ^{anti.} magnified user as ^{magnified} user.

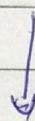
Goal

(find for all tuning para.)

Tuning para. — ① threshold of no. of simil. and,
 ② frequency of generated and,
 ③ threshold of simil.

5 Use pattern over test data.

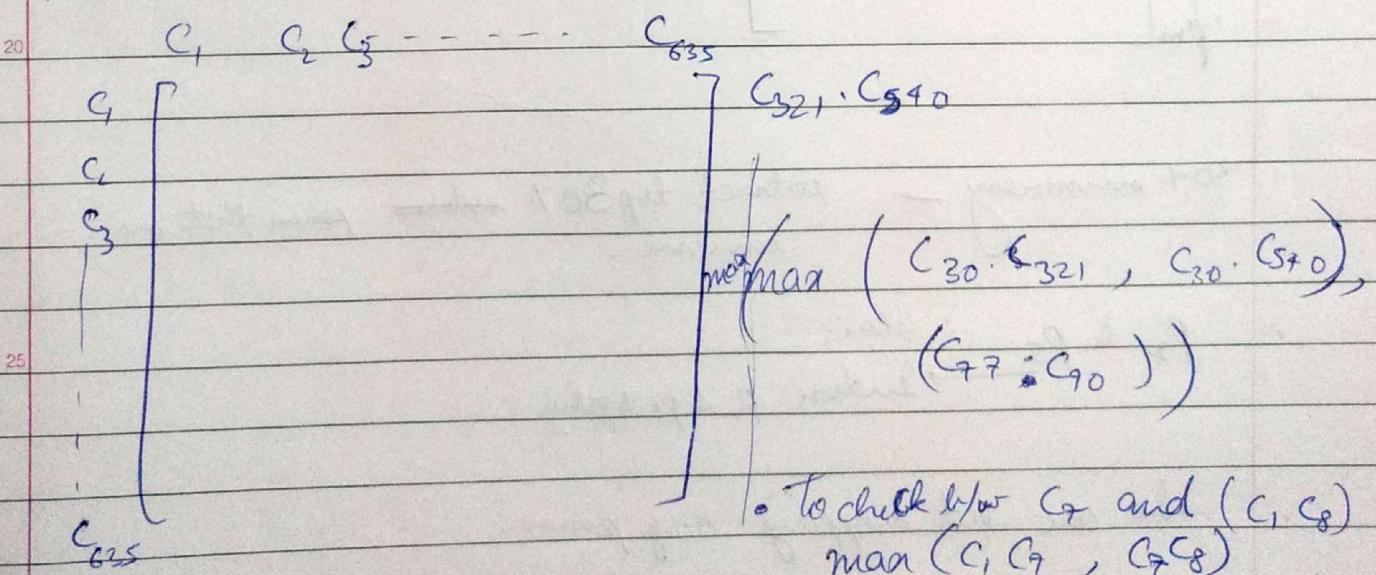
Day user_i — out of 100 — 80 in pattern
 — 20 out of n.



Set threshold (say 70 in pattern — OK)

(Based on assumption that users get good and can try their and not to harm sys./access maliciously) (Human Nature has to be taken into acc)

~~17/08/18~~



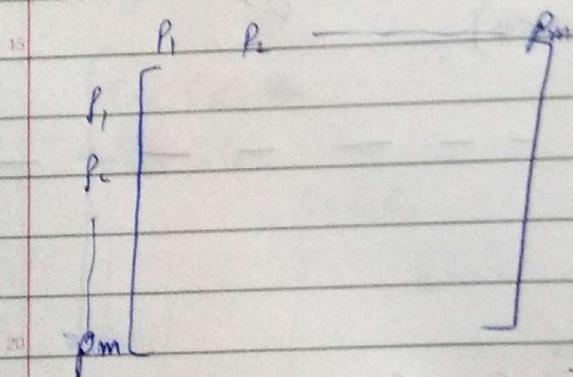
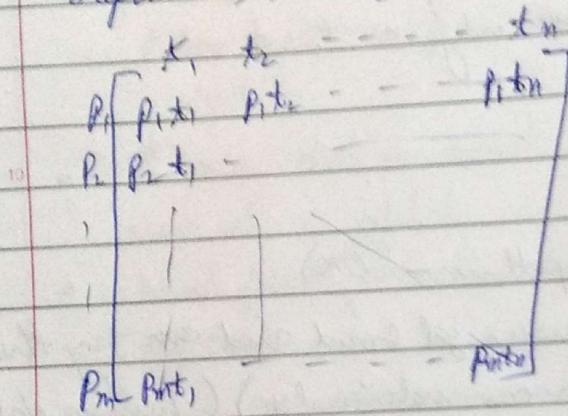
30 • To check for $(C, G_7 G_8)$ and $(C_2 C_3)$

$$\max ((C_1, C_7, G_8), C_2 C_3)$$

Text Summarization

- comprehensive view of doc.
- finding key-words

stopwords, lemmatization, stemming etc. (8/1)



30% summary - extract top 30% ~~similar~~ paras. that are unique.

- $P_2 \& P_5$: simi
→ Index P_2 & P_5 together.
- We are not dropping any paras.

- If 2 paras are more simi - consider them as 1.
If dissimilar - maybe not v related to para j so

semantic simi. on basis of para - not good idea

para can be converted to sentences.

- Better to find simi b/w sentences

Arrange sentences in descending order of simi.
→ displays more relevant

