

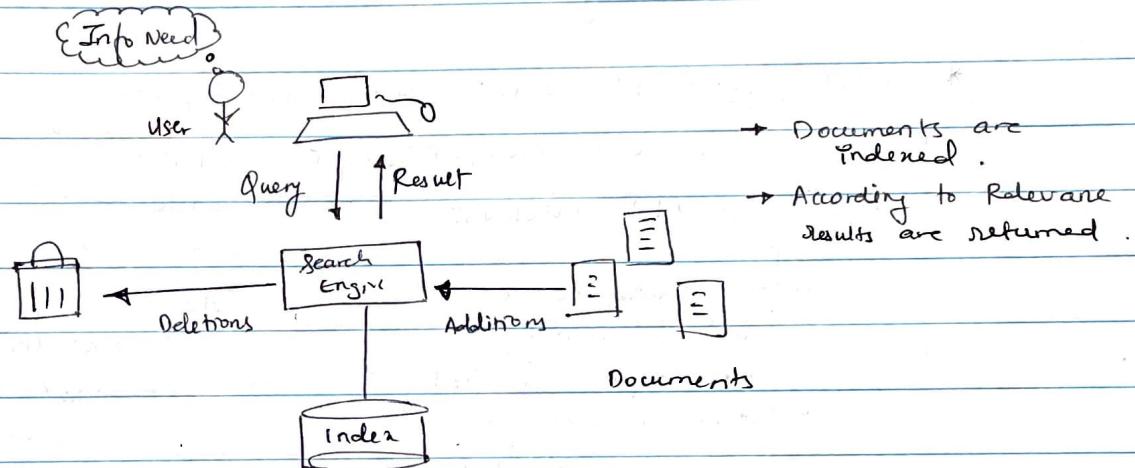
## INFORMATION RETRIEVAL

\* Represent, search, manipulate large collection of data.

\* Finding data from unstructured material.

\* Corpus: large repo. of docs. (majority E documents)  
Info Need: topic about which info reqd

} Need expressed in term of query



\* Types of Data

→ Structured : Relational Database -

- Tabular format - Highly organized
- Rows & cols related - Easy to store, process, access

→ Unstructured Data : textual, number, dates, BLOB (Binary Large Objects)

- not organized in predefined manner.
- 80% worldwide data
- Require more storage.
- Rich media (video, image, audio), analyzable -  
eg: Phone calls, log files, Social media comments, images

Semi Structured Data

- not completely organized
- structure for relational database
- self describing structure
- tags, attributes to separate entities within data
- eg: XML

\* Structured: Technology, Transaction management, version management, ~~scalable~~ scalability,  $\uparrow$  Flexibility, Query Performance  $\uparrow$  (Complex) -

\* Unstructured :  $\uparrow$  scalable, Not complex -

\* to avoid linear scanning of large documents, we need to index the documents for relevance.

TDI:  
Term Document  
Inverted Index  
Hash map  
type

## INVERTED INDEX

TDI Matrix Infeasible when record too large. (usually Sparse Matrix)

Index data structure, stores mapping from content (words or docs) to its location in document or set of documents.

### \* Steps for Building Inverted Index:

→ Token sequencer: sequence of each term in each doc, given doc id.

→ Sorting: alphabetically.

→ Group: by words → by doc id

→ Splitting into dict & Posting: store terms and pointer to posting list for each term.

→ Calculate Frequency:

### \* BOOLEAN QUERY RETRIEVAL: Query in Boolean Expression of terms.

Terms combined with AND, NOR, NOT

e.g.: BRUTUS AND CALPURNIA

Steps:

1. Locate Brutus in Dict
2. Retrieve its posting
3. Locate Calpurnia in Dict
4. Retrieve its posting
5. Intersect both the lists.

ALGO :-

answer ← ↔

while  $p_1 \neq \text{NIL}$  &  $p_2 \neq \text{NIL}$

do

if  $\text{docID}(p_1) = \text{docID}(p_2)$

Add(answer, docID( $p_1$ ))

$p_1 \leftarrow \text{next}(p_1)$

$p_2 \leftarrow \text{next}(p_2)$

else if  $\text{docID}(p_1) < \text{docID}(p_2)$

$p_1 \leftarrow \text{next}(p_1)$

else  $p_2 \leftarrow \text{next}(p_2)$

return answer

### \* RETRIEVAL MODEL:

Process of matching a query and document.

Basis of Ranking in search engines.

Techniques: query suggestion, query expansion, relevance feedback, user interaction and context to refine

key aspects of relevance

\* Important for retrieval models and evaluation measures.

- topical and user relevance.

- whether it is binary or multivalued.

binary: relevant / non-relevant.

multivalued: relevant, non-relevant, unsure.

### \* BOOLEAN RETRIEVAL: Exact match retrieval

Assumption: All docs equivalent in terms of relevance.

Merit :- Results predictable, easy to explain

- Any feature of documents can be used as operand.

- Efficient in fast elimination

Demerit :- Effectiveness on user.

∴ Equivalent relevance results in sometimes low or no retrieved  
of highly relevant doc.

IR image 1 word doc.

may return  
less user relevant

## VECTOR SPACE MODEL

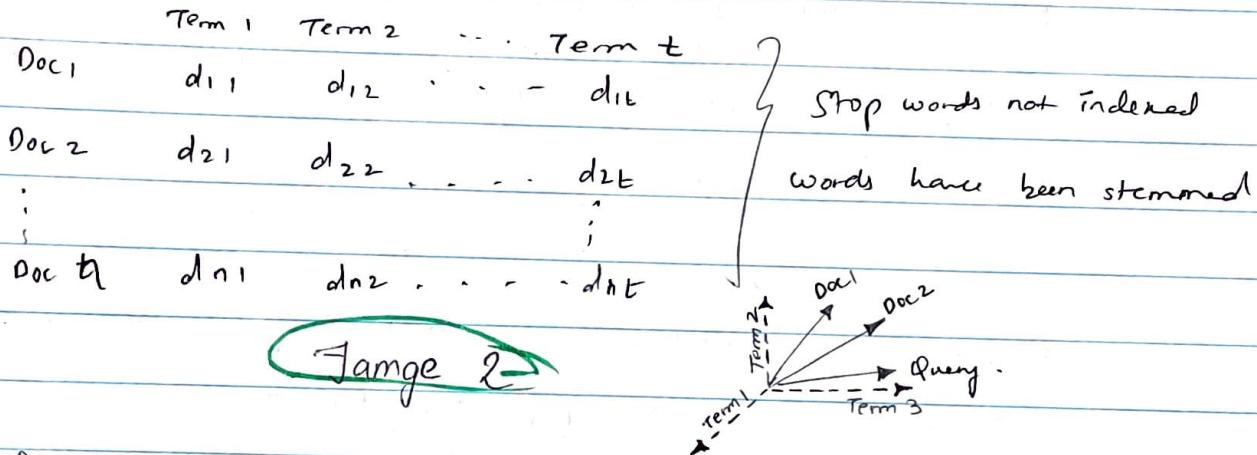
query assumed as part of  $\mathbb{f}$  dimensional vector space.

Document  $D_i$  represented by vector of index terms.  
 $D_i = (d_{i1}, d_{i2}, d_{i3}, \dots, d_{it})$   
 $d_{ij} = \text{weight of } j^{\text{th}} \text{ term}$

Query = vector of weights.

$$Q = (q_1, q_2, \dots, q_t)$$

$q_j$  = weight of  $j^{\text{th}}$  term in query.



Documents can be ranked by computing distance b/w points representing the documents and query.

Similarity measure used to calculate distance.

Cosine: angle b/w query and doc vectors.

vectors normalized, represented by vector of equal length  
 $\cos \theta$  / angle b/w two identical  $\theta = 0^\circ$       angle  $= 90^\circ$

Cosine of angle b/w two dissimilar vectors (non zero term) = 0, angle  $90^\circ$ .

$$\cosine(D_i, Q) = \frac{\sum_{j=1}^t d_{ij} \cdot q_j}{\sqrt{\sum_{j=1}^t d_{ij}^2 \cdot \sum_{j=1}^t q_j^2}}$$

### \* Tokenization:

Chopping text into tokens (pieces), throwing away certain characters.

Token: useful semantic unit

Type: class of all tokens containing same char sequence.

Term: Type that is included in IR system dictionary

Stemming

## \* Issues with Tokenization:

- characters like , ' : are thrown away, this may result in Data loss ex: apostrophe O'Neill.
- Language issues: eg: French, German writes compound nouns without spaces -  
Compound  
splitter module  
should be used  
(subdivide into multiple words, for words present in vocabulary)
- ↳ (East Asian languages): no space between words.  
word Segmentation can be used as prior linguistic processing

## \* Stemming: Reducing inflected words to word stem, or root form.

Applications: IR, determine domain vocabulary in domain analysis -

Algo: Porter Stemmer (5 rule algo)

Dawson "

Loving "

Xerox "

Knowetz "

## \* Stop words: common words, ignored when searching and retrieving result.

using stop list reduces number of postings. (positive)

## \* Term Frequency - Inverse Document Frequency (statistical measure):

evaluates how relevant a word is to a document in collection of docs.

metrics:-

tf  $\Rightarrow$  term frequency of word in doc.

idf  $\Rightarrow$  word across a set of documents. Closer it is to 0, more common a word is.

If word not common in docs, tf value comes near to 1

$tf * idf = score$  (weightage of particular doc)

↑ score ↑ relevant word in particular document  
particular doc set of docs

$$tf \cdot idf(t, d, D) = tf(t, d) \cdot idf(t, D)$$

↑  
term

$$tf(t, d) = \log(1 + freq(t, d))$$

t : term

d : particular doc

D : set/collection of docs

$$idf(t, D) = \log\left(\frac{N}{\text{Count}(d \in D : t \in d)}\right)$$

N : total documents

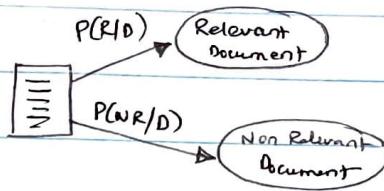
## \* Language Model : Probabilistic Model :

when info is returned to user, highly ranked docs will be displayed first, i.e. they have highest probability. This probability goes on decreasing as we move down the response.

How to calculate probability of relevance?

Treats info as classification problem.

Model classifies, then returns if relevant.



$P(R/D)$  = conditional prob (Prob of relevance, given representation of document). Similarly  $P(NR/D)$ .

$P(R/D) \rightarrow P(D/R)$  } Prob of doc being in relevant set.  
Bayes Rule.

$$P(R/D) = \frac{P(D|R) \cdot P(R)}{P(D)} \leftarrow \text{Normalizing constant}$$

Relevant if  $P(D|R) P(R) > P(D|NR) P(NR)$ .

$$\frac{P(D|R)}{P(NR)} > \frac{P(D|NR)}{P(NR)}$$

Score  $\uparrow$   
likelihood ratio  $\rightarrow$   $\frac{(P(D|NR))}{(P(D|R))}$

Relevant set

## \* Language Model :

represent text in various language technologies; speech recognition, handwriting recognition...

Simplest language model unigram language model.

If we have prob of each term in a document, we can predict what the next word in sequence would be.

n gram predict words on basis of previous  $(n-1)$  words.

words imp for topic will have unusually high probability.

\* Query Likelihood Ranking :- probability of query generation is the measure of how likely is that document is about same topic as the query.

Issue if word missing & it will give possibility

$$P(D/Q) \stackrel{\text{rank}}{=} P(Q/D) \cdot P(D)$$

$\uparrow$   
uniform

$$P(Q/D) = \prod_{i=1}^n P(q_i/D)$$

Ignore  $P(Q)$  as denominator normally const

$$P(q_i/D) = \frac{f_{q_i,D}}{|D|} \leftarrow \begin{array}{l} \text{frequency of query words in D} \\ |D| \leftarrow \text{number of words in D} \end{array}$$

Refer Language modeling ppt

## \* Okapi: The BM25 Ranking Algorithm

BM: Best Match

Ranking Algo / scoring value: (how query related to document) -  
Relevance of doc to search query.

BM25 used in probabilistic framework as core part

Simplest score for document d = idf weighting of query terms present in d.

Refer Image under 4

$$K = k_1((1-b) + b(\text{ld}/\text{Lave}))$$

↑ score ↑ relevance

$$\text{RSVd} = \frac{\log(N)}{df} \cdot \frac{(k_1+1)tf_{id}}{K((1-b) + b(\text{ld}/\text{Lave})) + tf_{id}}$$

$$= \frac{\log(r_{i+0.5})/(R-r_{i+0.5})}{(n_i-r_{i+0.5})/(N-n_i-R+r_{i+0.5})} \cdot \frac{(k_1+1)f_i}{K+f_i} \cdot \frac{(k_2+1)qf_i}{k_2+qf_i}$$

## \* RELEVANCE FEEDBACK

Polysemy, synonyms

↑ same word diff meaning

Helps overcoming the problem

Eg: plane fuel, fuel for aircraft, jet, ?

- Global methods for expanding or reformulating query.

Thesaurus generation for query expansion  
& spelling correction.

Query expansion with Thesauruses.

- Local method: adjust query relative to the doc that initially appeared

- Relevance feedback

- Blind/Pseudo Rel feedback

- (Global) Indirect Relevance feedback

## Relevance Feedback

user query

Initial set returned

user selects relevant

System recomputes better results based on user's feedback

System displays revised results

Assume:  
Ordered content unknown to user.  
→ Sys need changes after user's query or result.

## Relevance Feedback : Rocchio Alg

Find query vector  $\vec{q}$ , maximise similarity with relevant docs.

Rocchio:

$$\vec{q}_{opt} = \arg \max [ \underbrace{\text{sim}(\vec{q}, c_r)}_{\substack{\text{Similarity bet' } \\ \vec{q} \text{ and set of relevant docs.}}} - \underbrace{\text{sim}(\vec{q}, c_{nr})}_{\substack{\text{Similarity bet' } \\ \vec{q} \text{ and non-relevant docs.}}} ]$$

↑ recall

(+ve) & -ve feedback

valence

$\gamma < \beta$  using cosine similarity

$$\vec{q}_{opt} = \frac{1}{|c_r|} \sum_{d_j \in c_r} \vec{d}_j - \frac{1}{|c_{nr}|} \sum_{d_j \in c_{nr}} \vec{d}_j$$

diff bet' centroid of relevant and non relevant doc vectors

In practice

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

$\alpha, \beta, \gamma$  weights attached

RF is expensive,

long queries expensive to process.

often hard to understand why some particular doc retrieved after RF.

### \* Evaluation:

Required before deploying

we need test collection

1. Document collection

2. Info need, expressible as query

3. A set of relevance judgements, traditionally a binary assessment

Gold Std / Ground Truth Judgement of Relevance

No parameter Tuning.

Relevance is assessed relative to information need and not query.

Final decision by user

### \* Why evaluate?

- effectiveness & efficiency : to find right information in the time quickest.

- cost (processor, memory, disk space, networking).

- Comparison bet' techniques.

## Traditional Effectiveness Measures :-

### Assumption

- A doc in test collection would be relevant or irrelevant
- Relevance depends on info need and d, and independent of other docs.

$$\text{Precision} = \frac{\#(\text{relevant items retrieved})}{\#(\text{retrieved items})} = P(\text{relevant/retrieved}).$$

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

$\uparrow$ Search Engine	$\uparrow$ Analysis document	Relev	Non Relev
$P = \frac{tp}{(tp+fp)}$	$R = \frac{tp}{(tp+fn)}$	Retr tp Not Retr fn	fp tn

$$\text{Accuracy} = \frac{tp+tn}{tp+fp+tn+fn}$$

In almost all circumstances, data is extremely skewed., 99-9% non relevant

\* F-Measure : weight harmonic mean of P & R .

$$F = \frac{1}{\alpha \frac{1}{P} + (1-\alpha) \frac{1}{R}} = \frac{(B^2+1) PR}{B^2 P + R} \quad \text{where } B^2 = \frac{1-\alpha}{\alpha}$$

$$\alpha \in [0,1] \quad B \in [0, \infty]$$

default balanced F measure equally

$$\alpha = \frac{1}{2} \quad B = 1$$

$$F_{B=1} = \frac{2PR}{P+R}$$

## \* Latent Semantic Indexing

due to inherent  
vagueness  
with retrieval  
via keywords.

- \* Summarizing via index terms can lead to poor retrieval due to two effects:
  - many unrelated docs might be included in result
  - relevant documents not indexed by any query or keyword are not retrieved.

- \* LSI is an attempt to address the above issue.

It tries to match the documents based on the concept matching instead of index term matching term.

In this retrieval Matrix computation is used as the basis for information retrieval. Matrix computation in particular Single Value Decomposition (SVD) is used.

This filters out noise found in documents such that two documents have the same semantics. In a multidimensional space these docs will be located close to each other despite of not having common terms.

Term doc matrix  $A$ ,  $a_{ij}$  represent, term  $i$  appears in doc  $j$ .

SVD of  $A$  results in matrices  $U\Sigma V^T$ .  $\Sigma$  is diagonal matrix.

values in  $\Sigma$  = singular values. } sorted by magnitude, top  $k$  values taken

as means of developing latent semantic representation of matrix  $A$ .

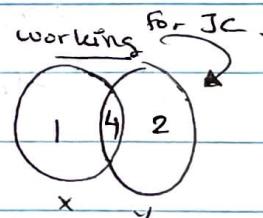
Remaining singular values set to 0.

first  $k$  cols are kept in  $U_k$

first  $k$  rows kept in  $V_k^T$

After setting results to 0.

$A'$  generated to approximate  $A = U\Sigma V^T$



## \* Ranked Retrieval :-

- system returns an ordering over the (top) documents in the collection for a query.
- Free text queries: just one or more words in human lang. query.
- No operators or expressions.

Assign a score say in  $[0, 1]$  to each document. } shows how well document & query match.

Assign score to query document pair -

If query term not in document, score should be 0.

More freq the term higher the score -

## \* JACCARD COEFFICIENT:

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

$$\text{jaccard}(A, A) = 1 \quad \text{jaccard}(A, B) = 0, \text{ if } A \cap B = 0.$$

### \* Issues with Jaccard

- privileges shorter docs.

More sophisticated normalizing of length reqd  $|A \cap B| / \sqrt{|A \cup B|}$

- does not consider ~~tf~~ term frequency.

- does not account for term informativeness.

↑  
Imp of term in doc.

### \* Log frequency weighting:-

- ★ weight of term  $t$  in  $d$  is

$$w_{t,d} = \begin{cases} 1 + \log_{10} \frac{tf_{t,d}}{N} & \text{if } tf_{t,d} > 0 \\ 0 & \text{otherwise.} \end{cases}$$

Score for doc query pair ; sum over terms  $t$  in both  $q$  and  $d$ :

$$\text{score} = \sum_{t \in q \cap d} (1 + \log_{10} tf_{t,d})$$

- ★ idf weight

$\underbrace{df_t}_{\text{no. of docs containing } t}$

↳ inverse measure of informativeness of  $t$ .

$$df_t \leq N$$

$$idf_t = \log_{10} (N / df_t) \quad \left. \right\} \text{log dampens the effect of idf}$$

tf - idf weight is the product of its tf weight and its idf weight

$$w_{t,d} = \frac{\log(1 + tf_{t,d})}{1 + \log tf_{t,d}} \times \log_{10} (N / df_t)$$

↳ it affects queries containing  $\geq 2$  terms.

Score for doc given a query -

$$\text{Score}(q, d) = \sum_{t \in q \cap d} tf \cdot idf_{t,d}$$

## \* Language Modelling (Cont)

Query likelihood: major problem, if any words missing from doc, score given will be zero.

Smoothing used to avoid this:

lower probabilities of word seen in doc, and assign the rest of words not seen

estimate for unseen words based on freq occurrence of words in doc collection.

$P(q_i/c) = P$  for query word  $i$  in collection of docs  $c$ .

estimate =  $\alpha D P(q_i/c)$

$\alpha D$  = coeffs of controlling prob assigned to unseen words

Prob estimate for word seen is  $(1 - \alpha D) P(q_i/D)$

$$\text{Dirichlet smoothing} \quad \alpha D = \frac{|D| + \mu}{|c|} \quad \begin{cases} \text{depends on doc length} \\ \text{more effective} \end{cases}$$

it set empirically } best results when  $1000 < \mu < 2000$

$$P(Q/D) = \sum_{i=1}^n \log \frac{f_{q_i/D} + \mu \frac{c_{q_i}}{|c|}}{|D| + \mu}$$

$c_{q_i}$  = no. of times query word occurs in collection of docs.

$|c|$  total no. of word occurrences in collect

$f_{q_i/D}$ , freq of  $q_i$  in  $D$ .

$|D|$  no. of words in  $D$ .

### \* Errors in stemming

Over stemming: university & universe

Under stemming: data & datum

↑ could be reduced to data & wrong

### \* Access time less for memory

\* Storage high for Disk

Image 5