

Day / Date: _____

CHAPTER 1 - BOOLEAN RETRIEVAL

Information Retrieval - finding material (usually documents) of an unstructured (usually text) - that satisfies an information need from within large collections.

→ **Unstructured Data** - Data which does not have clear, semantically over, easy-for-comp structure.

* IR is also used to facilitate 'semi-structured data search' (to search a doc where title was java & body has threading)

IR vs DB → Differences

① **Format of data** - DB: structured & clear semantic

IR: semi structure

② **Queries** - DB: SQL (formal)

IR: expressed in natural lang.

③ **Result** - DB: exact result

IR: sometimes relevant

Text vs. Data

Data Mining

Data Obj → Numerical struc. + categorical data

Data represent → straightforward

Goal → to extract info from a dataset & transform it into understandable struc. for further use.

Goal - extract info; process of exploration
data analysis that leads to new info.
Text Mining

Textual data (unstruc + semi-struc)
(but have a hidden structure)

Complex

- ↓ - Item set mining
 - Classification
 - Clustering
 - Recommendation
- 3 types of text mining / method.

Clustering → coming up with good grouping based on content of documents.

Classification → which classes each of the set of doc. belongs to. Dazzle

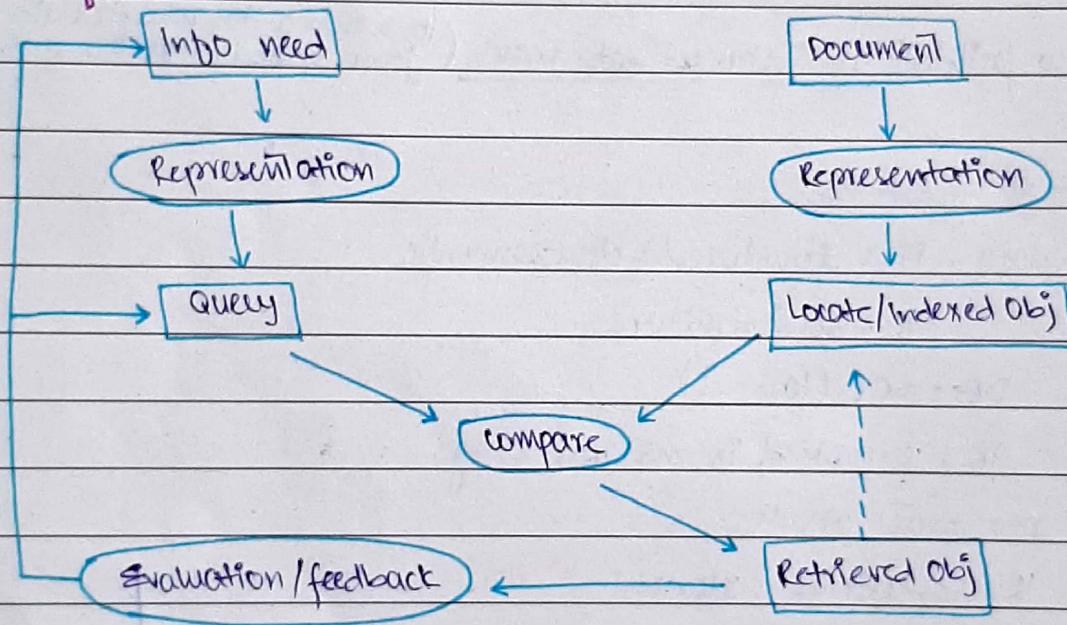
Day / Date: _____

IR system can also be distinguished by the scale at which they operate.

- ① web search
- ② space of enterprise, institutional & domain specific search
- ③ personal info retrieval

Information retrieval vs. Text Mining

Basic IR process



Representation of both can be diff. so we bring them on same level by extracting something common & then compare them

- Indexing is done before the query is applied so that all (eg. frames) can be indexed & its easy to compare.
- work on shallow features of documents — IR .
- Retrieves / works on relevant text / document — text mining .

Content wise IR System

- size (desktop, web scale)
- static (offline) or dynamic
- Type (Text, multimedia)
- Exact match vs. Best match.

Info need → topic about which user desires to know more
Query → what the user conveys to the comp in an attempt to communicate the info need.

Day / Date: _____

Adhoc Retrieval Sys - Sys that aims to provide doc. from within the collection of documents that are relevant to an arbitrary user info need, communicated to sys by means of user-initiated query. e.g. Google, Bing.

Dimensions of IR

- vertical search
- forum search
- P2P search
- literature search

Applications for

- Images
- Audio
- Music
- Bibliographic search

Task

- Filtering
- Clustering
- Content share / collaborative authoring

→ Images algo are diff. from usual search

(1) → QBIC - Query by image content

↳ uses colour, texture, shape features for search & text based search.

(2) → Mirage → supports queries based on color, layout & texture

↳ gives user more control over search process.

(3) → Visual SEEK - research prototype.

↳ considers spatial relationship b/w obj

↳ matching spatial relationship b/w obj & visual features together gives powerful result.

→ Music Search Engine

↳ search by metadata e.g. artists biography, music reviews,

↳ can create playlist, recommend, search for lyrics.

In Boolean model, we can simple scan-through doc by grep (linear scan). But we need more for other purposes

- ① To process large data
- ② To allow more flexible matching operations
- ③ To allow ranked retrieval.

Day / Date:

→ Mathematical Search Engine — to search mathematical content

→ Qs/Ans system — sys that allows extraction of ans to a request.

→ General Search Engine – Eg. google, yahoo, bing.

WHAT IS A MODEL?

• Abstract representation of a process or obj

→ used to study properties, draw conclusion & make prediction

→ Quality of conclusion depends upon how closely the model represents reality.

- A retrieval model describes human & computational processes involved in ad-hoc retrieval.
 - ↓
defines relevance, explicit
 - ↓
resource
(info needed)

- * BOOLEAN MODEL \rightarrow based on set algebra.
 - \hookrightarrow mathematical
 - \hookrightarrow features (dependency check & how we are modelling)
 - \downarrow pose any query which is in the form of a Boolean expression of terms which are combined with AND, OR, NOT.
 - \hookrightarrow the user bring some word query $\&$ we have a resource D.
 - ($w_i \rightarrow$ string)
 - $\{S\} \rightarrow$ collection of well defined & distinct objects
 - $\{D\} \rightarrow$ 2 dimensions

۹۰۸۷۴

3

w₃ mein jitney |
aik jagah aagaye

d_i has (w_1, w_2, \dots, w_n)

`ej` is another document. -

SEARCH RESULT B (q₁, D) = S
result boolean document
we are
retiring

term | w | o | d | document

$w_2 \quad 0 \quad 1 \quad 0$

→ pick each row's air gap

doc doal lein in sets mein

term der matrix

(team * doc)

> dimensions

Dazzle

Day / Date: _____

subsets

$$SR_B = V_i$$

op should be well defined

→ Algebra ($\{\cdot\}$, operation)

↓
(int, decimal)

↓
($x, \div, +$)

set of sel alg.

multiset → (has repetition)

$$\downarrow A = [a, a, b, b]$$

→ if img/sound

Boolean model is based on multiset aka (bag of features/word)

↓ if text

* linear or vector algebra, probabilistic algebra } from table in slides

→ eminent way : important

↓ there is a hidden relationship that we are not expliciting represent but we have captured it.

→ It represented explicitly → so we use models.

2 assumptions for Boolean model

50% documents present → sparsity.

① We can view doc & see features.

80% not present → not a good doc.

② We electronically view doc & see features.

↓ (from dictionary)

* 10^6 → max vector we can take for processing

→ Term Doc matrix mein rows ke and karey ki ya koi aur operation.

Problems of Term Doc Matrix

→ Sparse matrix (most of the arrays are zero)

→ Exact matching

→ Complex Query Formulation.

Solution

→ Inverted Index

Day / Date: _____

Adhoc Retrieval System

Adhoc → user can bring any type of search.

→ Simple interface with a text box to place query & search button.

→ sys that aims to provide doc from within the collection that are relevant to arbitrary user info need

Effectiveness of IR System

↓ quality of the search result

② factors of IR sys.

1. Precision : what frac of returned results are relevant to info need?

$$\frac{\text{relevant}}{\text{total}} = \frac{20}{28}$$

2. Recall : what frac of the relevant doc in the collection were returned by system?

$$\frac{\text{relevant}}{\text{irrelevant}} = \frac{20}{40} = \frac{20}{40}$$

Precision = relevant retrieved / total retrieved

Recall = relevant retrieved / total ~~retrieved~~ collection

(learn)

100% recall → If we send all doc back $\frac{20}{20} \rightarrow \text{retrieved} = 100\%$

100% precision → If we sent 1 relevant result $Y_1 = 100\%$.

8 → relevant

$$\rightarrow \text{precision} = \frac{8}{18} = 0.44$$

10 → irrelevant

$$\rightarrow \text{recall} = \frac{8}{20} = 0.4$$

20 → total relevant

Day / Date: _____

list of doc present

INVERTED MATRIX INDEX

w₁ → Doc 1, Doc 3 → Doc 6

⋮ doc lds

w_n



(Doc read kya uski terms nikalay rahoing ge)

linked list of

terms

dictionary

posting list

Intersect (Posting List Algo) → (works like merge func)

0 - 0 0 0

↓ ↓

0 0 - 0 0

(will compare

jab bhi aik list khota hai na ho jaye)

Query Optimization

process of selecting how to organize work of answering a query so that the least total amount of work needs to be done by the system.

→ major element for this Boolean query is order in which posting list are accessed.

Intersect (Optimized)

(Sari bhi kr raha hain. Jo chhoti list hai wo pehlay aati hei)

re arrange terms so that we perform least amount of work. The no. of work depends on the shortest freq. of posting list.

t₁ t₂ t₃ t₄

1200 115 37 400

(t₁ ∩ t₂) ∩ t₃) ∩ t₄)

(37)

(37)

(37) → shortest

freq among posting list.

Intersect (<t₁, ..., t_n>)

terms ← SortByIncreasingFreq (<t₁, ..., t_n>)

result ← postings (first (terms))

terms ← rest (terms)

while terms ≠ NIL and result ≠ NIL

do result ← Intersect (result, posting (first t))

terms ← rest (terms)

return result.

Dazzle

Day / Date: _____

BOOLEAN MODEL

BM is most simple & most used. Considers that doc contains features (words, phrases) & the user's query is about these features.

Advantages

- ① Clean formalism, easy to implement & intuitive concept.
- ② Results are predictable, relatively easy to explain
- ③ Many diff. features can be incorporated.

Disadvantages

- Exact matching • Query formation is hard • All terms are equally weighted
 - flat result → all docs are at same level;
all are equally important
- ↓
(has feature equal weight ka hai)

→ Drawbacks of Boolean Model

- ① → no criteria for partial matching ② → No ranking of doc (grading scale missing)
 - ③ → info need has to be translated into Boolean expression.
- In conclusion → BM frequently returns either too few or too many doc in response to a query.
- ④ Imposes binary criterion.

→ Extensions of Boolean Model

→ Weighted Boolean Model / Fuzzy Set Model / Probabilistic Model / Vector Space Model / SIMILARITY.

① Extended Boolean Model → to overcome the drawbacks of BM by making use of

② Fuzzy Set Model partial matching & term weights as in the vector space model.
combines characteristics of vector space model with Boolean
Algebra & ranks similarity b/w query & doc & select doc
maybe somewhat relevant to the queried term.

Techniques are based on EBM; to define

flexible sys i.e. sys that can represent &
manage vagueness & subjectivity which

characterizes the process of IR & info representation.

Dazzle

Day / Date: _____

THE TERM VOCABULARY & POSTING LISTS

(Chp 2)

* Word - a delimited string of characters as it appears in the text.

Term - a 'normalized' word (case, morphology, spelling etc); an equivalence class of words.

Error → in IR, collection of emails (dataset) → each file has mails of a single user (example)

* we need to find index of what we want to search.

Challenges of Document Processing

→ what is the format? (pdf, word) → what are hidden fields? ↗ Classification problems

→ what is the language? (can be multilingual).

* we need to convert the byte seq (file) into a linear sequence of characters.

→ what character set is in use? (single or multibyte?).

→ size of document? (A file, several, group of files?)

Tokenization → Given a character sequence & a defined doc unit, tokenization is the task of chopping it up into pieces called tokens, perhaps at same time throwing away certain characters such as punctuations.

Token → an instance of a seq. of characters in doc that are grouped together as a useful semantic for processing.

Issues in Tokenization

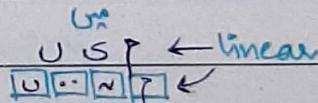
(9/9/2019)

→ Encoding: (o-ed → 1 token; San Francisco → 2 tokens) or (team name of numbers)

③ languages → (urdu, french) ↗ we need to determine the correct encoding..

(we need to decide if a character makes a word so that we must throw a token)

The documents are usually linearized →



Stop words → (to, am, the) → in isolation they don't give any semantic & are very frequent.

• with list, one excludes commonest word from the dictionary.

↓ but we need them for phrase queries (King of Denmark)

* relational queries (flight to London)

→ Indexing a large doc causes issue of index granularity.

↳ (can be handled by implicit/explicit proximity search)

Dazzle

Type → class of all tokens containing the same character sequence.

Term → type that is included in the IR system's dictionary.

Day / Date: → issues of tokenization are lang-specific — lang identification based on classifiers that use short character subsequences as features is highly effective; most langs have distinctive signature patterns.

Normalization (Normalization / Equivalence classing of terms)

→ each word should be changed to lowercase (agreed practice)

→ An abbreviation is mapped to a single word (hyphens are removed)

→ The common or same words are tokenized to same token.

* Process of canonicalizing tokens so that matches occur despite superficial differences in char. seq. to the tokens.

MOST IMPORTANT CRITERIA

→ See how does your user write queries for words, words

→ codilla / diawtia (search)

Morphological Analysis

(who, whom)

↳ Inflection → word changes not to its original / but tokens change to its original

↳ adding a suffix to a word, that doesn't change its grammatical category.

→ each affixes in verbs (-ing, -ed, -s); plural nouns). category

→ don't change type of word (-s), plural noun(s).

↳ (doesn't change type of word)

→ British & American English (colour, color) || use of thesaurus

→ Multiple languages

→ Dealing with complex words; foreign names → have unclear spelling or transliteration standard giving diff. spelling.

→ an alternative is

↳ we use heuristic to equivalence class or phonetic equivalents → best example is soundex algo.

→ Thesaurus & Soundex

STEMMING & LEMMATIZATION

Reduce inflectional forms & derivationally related forms of a word to a common base form.

Stemming → chops off the end of words or removal of derivational affixes to reach goal.

Lemmatization → remove inflectional endings only to return the base / dictionary form of word → lemma.

e.g. SAW → stem : s

lem: saw or see

Dazzle

(*tokens need to be indexed)

Day / Date: _____

walking, walked, walked → tokens is 'walk' (the base one)

→ Each algorithm works on different basis.

→ Stemmer → general + fast → but the output doesn't make sense to user [drawback] use

Poetic's Algo is most common & effective algo for stemming.

Search Engine → Stemmer, summarization of doc → Lemmatization } depends on application.

* Stemmer → using specific rules, but requires less knowledge. on app.

* Lemmatizer → complete vocab & morphology analysis to correctly lemmatize a word.

→ Stemming increases recall while harming precision done by dictionary differences

POSITIONAL & POSTING & PHRASE QUERIES

→ Phrase Queries

Want to be able to answer such queries as 'Stanford University' → as a phrase.

$Q = (t_1, t_2)$ → both should be present & adjacent as well. t_1 should come first.

↓

Post Query Process (will check that t_1 & t_2 are adjacent in every doc)

After result is generated, we need to validate our result (This increases work)

Amount of work \propto no. of docs in result.

Solutions

① Bi-Word Index - consider every pair of consecutive terms in a doc as a phrase.

Eg. friend roman countrymen } 2 bi-words

①

②

Each bi-word is a dictionary term

↳ results are immediate.

→ (long query requires a lot of work)

↳ so we'll divide it & then check from all documents (in dividing & then checking from all documents)

↳ (as a feature index)

having grown index keeping for easy & then search it in posting list & make a posting list.
then scan in doc to make a posting list.

Dazzle

- ③ Make any string of words in form $N \times N$ to be an E. Biword.
- ④ Each such extended biword is made a term in vocab.

Day / Date:

- ② Extended Biword Index -
- ① Tokenize text & perform part-of-speech tagging.
 - ② Group terms in nouns & other classes.

• used for searching some concepts. $T \ X^n \ T \rightarrow T$ search krtay hain par beech mein se kafi words X^n skip krdein ga. skip krdein ga.

e.g. coins are in pockets

$\frac{N}{\text{coins}} \ \frac{X}{are} \ \frac{X}{in} \ \frac{N}{pockets}$

→ take out the nouns (N) & ignore preposition etc (X) etc (a).

→ lookup in index for coin pocket!

Issues in Bi-word

- too many false positives → for which we'll want to do post indexing to validate that result
- no. of indexes increases → memory size increase (carefully place biwords in index)

↓
not a standard solution, but a compound strategy.

② Positional Indexes

→ most commonly used

* need to access the inverted index entries for each distinct term; has more info.

→ docID, freq, & position of the term position of the term.

* For each term in vocab, we store postings of the form $\text{docID} : \langle \text{position}_1, \text{pos}_2, \dots, \text{pos}_{\text{last token appears}} \rangle$

$\langle \text{term}, \text{no. of docs containing term} \rangle$ positions in which token appear.

↓ able to answer length of phrase query

→ check that both terms are present in doc, as well as the positions of appearance in doc

all compatible with phrase query being evaluated; use the following:

$\{w_1, w_2, \dots, w_k\} \rightarrow \{w_1, \{w_2, \dots, w_k\} \text{ tokens apart}\}$

→ Proximity Query - words queried noun & gives the distance it needs to be apart.

(say, x_1, x_2, \dots, x_k) → distance = k gives the distance it needs apart.

→ positional that has a distance of K ; all those docs will be in result.

e.g. employment/H3n:1K means within k words of (on either side) doc. result mein range.

Positional Index Size → It expands required postings storage significantly.

(Binary search (kuch common cheez ko) bohot
sabdy result ayega because unkormege karna
is not preferred (red))

Day / Date:

phrase → phrase → to be or not to be

to []
be []
or []
not []

3 posting list of each for positions
each position

1: []
2: []
3: []
4: []

Searching for valid results in each doc.

positions at
which each word is at
word is at.

Rules of Thumb

→ positional index is good for eng-like queries like queries

→ Positional index is 2-4x as large as a non-positional index.

→ Positional index size 35-50% of vol. of original text.

Combination Scheme

A combination of binary + positional phrase would do effective search.

→ Need to decide which method should be used for which query

* Where individual words are common but desired phrase is comparatively rare → expansive query

> If posting lists are really long → we'll use skipping pointers at indexing time.

We walk through two posting list simultaneously ; if long → time taken would be $O(m+n)$

Use skip list by augmenting posting list with skip pointers (at indexing time). lots of

↓ Effective shortcuts that allows us to avoid processing part of the postings list that will

not figure in the search result. Few pointer comparison w/ longer span of skip.

① → where to place skip pointers ② how to do efficient merging using skip pointers.

* Presence of skip pointers only helps for AND queries & not for OR queries.

Where do we place skip pointers?

Tradeoff : • More skips means shorter skip spans & that we are more likely to skip. But also

mean lots of comparisons to skip pointers & lots of space storing skip pointers.

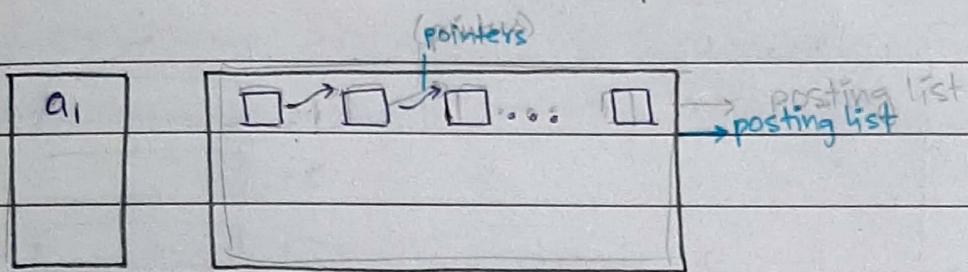
• Fewer skips means fewer pointer comparisons, but then long skip spans which means that there will be fewer opportunities to skip.

Heuristic → P_2 length of posting list ; use SP evenly spaced skip pointers.

Dazzle

DICTIONARIES & TOLERANT RETRIEALS

(Chp 3)



* dictionary is kept in memory
for faster access

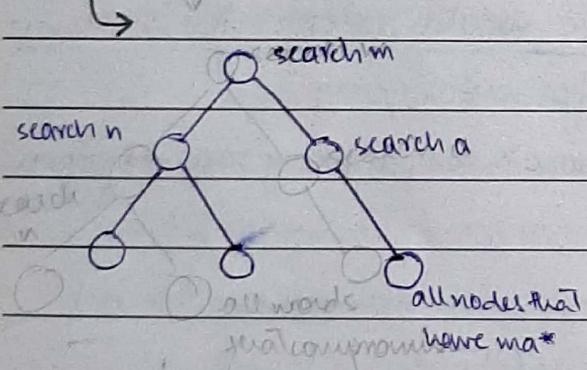
Dictionary

↳ determine whether each query term exists in the vocab, if so identify the pointer to the corresponding posting.

The vocab lookup uses the data struc. "dictionary" → two sol: ① Hashing ② Search Trees.

↓
The entries in the vocab (terms) are referred to as keys. efficiently perform these queries

Data Structure for Dictionary

① Trees② HashingTrees overcome issues by:

→ allowing us to deal prefix

* Best tree is Binary Tree ~~or BTree~~

↓ each internal node has 2 children

↓ searches with $O(\log m)$

drawback → rebalancing

↓ as terms are inserted, or deleted from

binary search tree, it needs to be rebalanced.

Hashtable

Each vocab term is hashed in an integer

Pros → lookup is faster than for trees ($= O(1)$)

Cons → ① No easy way to find minor variant

② No prefix search (words starting with 'pre-')

To reduce rebalancing, use B-Trees.

↳ a search tree in which every internal node has a no. of children in interval $[a, b]$

where $a \leq b$ are two integers; B-Tree can be viewed as 'collapsing' multiple levels of binary tree

In one → advantageous when dictionary is disk-resident || serves as prefetching Dazzle

→ Search tree demands that char. used in doc have a prescribed ordering. binary test.

Day / Date:

Hello \$ \rightarrow ending boundary of word
Hello \$ H \rightarrow start from H, \$ Hello is prefix so Hello is a word
Hello is prefix

→ WILDCARD QUERIES

① Trailing wildcard query \rightarrow query (mon*) mon*

suffix \rightarrow lo \$ Hel
prefix \rightarrow starting from this

② Leading wildcard \rightarrow query (*on) * on

suffix \rightarrow ending at this so lo is prefix for Hel. \rightarrow This

③ General \rightarrow query (m*n*tion) n*ition

so lo is prefix for Hel.

→ efficiently perform these queries ; Why DO WE NEED WILD QUERIES?

→ ① User is uncertain of the spelling of a query term (unintentional)

Sydney vs Sidney \rightarrow S*dney

→ ② User is aware of multiple variants of spelling term & (intentionally) seek doc containing any of the variants (e.g. colour vs. color)

→ ③ User seek doc containing variants of a term that would be caught by stemming, but unsure whether the search engine performs stemming (e.g. judicial vs judiciary \rightarrow judicia*)

→ ④ User is uncertain of the correct rendition of a foreign word or phrase (e.g. Universit* stuttgart)

mon* \rightarrow use a simple tree following symbol m, o & n ; lookup in inverted index to retrieve all doc starting with mon.

mon \rightarrow use terms to reverse tree create a reverse B-Tree & lookup then.

**se*mon \rightarrow se* \rightarrow do by normal tree ; *mon \rightarrow do by reversing B-Tree.

mon* \rightarrow simply do it \rightarrow we'll get desired result

General Wildcard Queries to handle like (dictionary reverse) & create a

Two techniques to handle them : tree.

① Permuterm Index \rightarrow aNb \rightarrow we'll get desired set of doc. * Dictionary \rightarrow dictionary size becomes large (as it has all rotations)

② \$: to mark the end of a term ; hello is written as hello \$

create a permuterm index having various rotations of each term with \$ that links to original vocab

set of rotated terms \rightarrow permuterm vocab. \rightarrow will keep permuterm index in the memory along with dictionary so we don't have to access it from disk

* flag to rotate one letter at a time

rotations need to map on one word

* use B-tree as data structure. \rightarrow cat\$ \rightarrow cat* \rightarrow cat\$ c

create an index above our dictionary

one item at a time that allows us to select rotation with a given prefix

cat \rightarrow cat\$ \rightarrow cat* \rightarrow cat\$ c \rightarrow t\$ ca \rightarrow t\$ ca

create an index above our dictionary

map rotations on a word (prefix/suffix) on our dictionary

Dazzle

Day / Date: _____

② Bigram: k-gram indexes for Wildcard Queries

- chop entire term in k-characters at a time ($k=2$ if bi) {sequence of k-characters}
- In a k-gram index, the dict. contains all k-grams that occur in any term in vocab. Each posting list points from a k-gram to all vocab terms containing that k-gram.

Disadvantage: Needs post filtering (post validation of result doc) b/c gives false positives
↳ results are individually checked against query. ↳ low precision

- Processing is expensive b/c of added lookup in special index, filtering & then standard inverted index.

SPELLING CORRECTIONS

- Forms of Spelling Corrections

① Isolated Term → Query = (w_1, w_2, w_3) → correct single query term at a time.

→ check each word on its own for mis-spelling

→ will not catch typos resulting in correctly spelled word

e.g. from → form

② Context sensitive

(alternate is to use biword statistics in collection, use log of queries issued by user)

correct the typographical errors such as flew from or flew form

→ enumerate corrections of each query term, substitute the correction word

& run query for matching results → (Expensive b/c too many combinations)

→ fundamental premise → There is a lexicon ^{max} _(dictionary) from which the correcting spellings come from

↓ if not a lexicon based → find closest → If only 5 or few doc → we correct query & find closest.

(dictionary ke har word se query ke word ki relation banayi & then give the closest)

Techniques for addressing isolated term correction:

① Edit distance (Levenshtein distance)

② Weighted Edit Distance

③ k-gram overlap

Two basic principles:

① for mis-spelled query, choose the nearest one.

② when two correctly spelled queries are tied (or nearly tied), select one with more common.

grnit & grnt → for grnt

choose the one that occurred more times in collection.

Day / Date:

Edit Distance:

e.g. change cat \rightarrow dog \rightarrow dog; weight will be 3 b/c 3 words changed
fast \rightarrow past past; 1 weight.

Weighted Edit Distance:

B/w 2 strings in time $O(|s_1| \times |s_2|)$ $\rightarrow |s_1| \rightarrow$ length of string s_i . Use dynamic Programming.

K-gram Overlap:

fast fa as st
past pa as st

$$\frac{|s_1 \cap s_2|}{|s_1 \cup s_2|} = \frac{2}{4} = 0.5$$

jaccard distance
jaccard coefficient
jaccard similarity

added distance \rightarrow subsequence matching prevalence max (people have used it in query at most)

added distance max = no. of letters in word. (query at most)

1 3 check most frequent from this

2 3 from this

3 3 In k gram, k is usually 2 or 3.

4 3 jaccard value \uparrow so words will appear at top

jaccard value \downarrow so words will appear at bottom

3 ranking of documents

is done on the basis

of jaccard coefficient

PHONETIC CORRECTIONS

Mis-spellings that arise b/c the user types a query that sounds like the target term.

main idea \rightarrow to generate, for each term, a phonetic hash so that similar sounding terms hash to the same value.

\rightarrow Algos for such phonetic hashing are commonly collectively known as 'soundex Algos'

far defined

Soundex Algo is built on following scheme:

- ① Turn every term to be indexed into a 4-char. reduced form. Build an inverted index from these reduced forms to the original terms; call this the soundex index.
- ② Do the same with query terms.
- ③ When query calls for soundex match, search this soundex index.

Dazzle

Day / Date: _____

→ Issues in Spelling Corrections

→ Spelling correction is computationally expensive

↓ only runned on docs which have result less than (close) to threshold.

→ Soundex

Uses of heuristics to expand a query into phonetic equivalent.

↓ soundex algo used (study it). → classic algo → not very useful for info retrieval.

↓ eg. Hermann maps to → H655

* Given a query we compute its Soundex code & then retrieve all vocab term matching this soundex code from soundex index, before running the resulting query on the standard inverted index.

context sensitive Analysis Spelling Corrections

↓ Hit Based Spelling Correction:

Suggest the alternatives that has lots of hits.

$$M = K^{\frac{1}{B}}$$

cte $\propto \frac{1}{g}$

INDEX COMPRESSION

(Chp 5)

→ WHY DO WE NEED COMPRESSION?

- ① use less disk space (saves money)
- ② increased used for caching; search sys uses some pairs more frequently. Fit more info in main mem. (increased speed)
- ③ faster transfer of data from disk to memory (disk → memory)

→ sufficient decompression algs run fast; Premise: Decompression algs are fast (true).

↳ dictionary compression is effective for effective IR sys esp. large scale.

(simple relation b/w collection size & vocab size is linear in lg n)

→ **Hemp's Law:** Estimating the no. of terms

Typical val for parameters

$$30 \leq K \leq 100$$

$$b \approx 0.5$$

→ Helps in estimate vocabulary size as a function of collection size.

$$M = K T^b$$
 (vocab size M as a func collection size T)

→ parameter K is quite variable b/c vocab growth depends on collection (nature + how compressed)

→ case folding & stemming reduce growth rate of vocab, including num & spelling errors increases it.

→ Regardless of the value of parameters, Hemp's law suggest ① dictionary size continues to increase with more docs in collection, rather than max vocab size being reached. ② size of dictionary is quite large for large collections.

→ **Zippf's Law:** Modeling the distribution of terms

To understand how terms are distributed across doc; characterize the properties of algo for compressing posting lists.

Zippf's law: If t_1 is the most common term in collection & t_2 is the next most common & so on, then collection freq c_i of i^{th} most common term is proportional to $\frac{1}{i}$: $c_i \propto \frac{1}{i}$ ($c_i = a/i$)

• RCV1 Collection → used as benchmark data set; intuition is that freq decreases very rapidly with rank.

1 document has approx 200 token, $n \rightarrow 400$, avg byte per token 6 bits. (space)

→ DICTIONARY COMPRESSION

Main goal → To fit it in the main memory or atleast large portion of it to support high query search throughput.

- ① To be able to design limited sys for limited hardware eg. phonetic → one term of dicti.
- ② Fast startup time by having to share resources with other apps. → entries

$$400 \times 28 = 11.2 \text{ MB}$$

① DICTIONARY AS A STRING

Sort vocab lexicographically & store it in an array of fixed width entries.

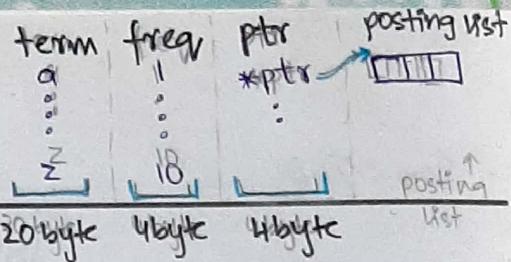
- 4 bytes → doc freq, 4 bytes → pointer to its posting list, 20 bytes → term

→ $M \times (20 + 4 + 4) = 400,000 \times 28 = 11.2 \text{ MB}$ for storing dictionary in this scheme. Dazzle for Reuters-RCV1

- using fixed-width entries is wasteful
- * avg term length in Eng = 8 char (12 wasted)
- no way of storing terms of more than 20 char.

↓ Day / Date: _____

overcome these by



- Storing dictionary terms as one long string of characters

- * The pointer to next term is also used to mark the end of current term.
- * saves us 60% space compared to fixed width

- now also need to store term pointers.

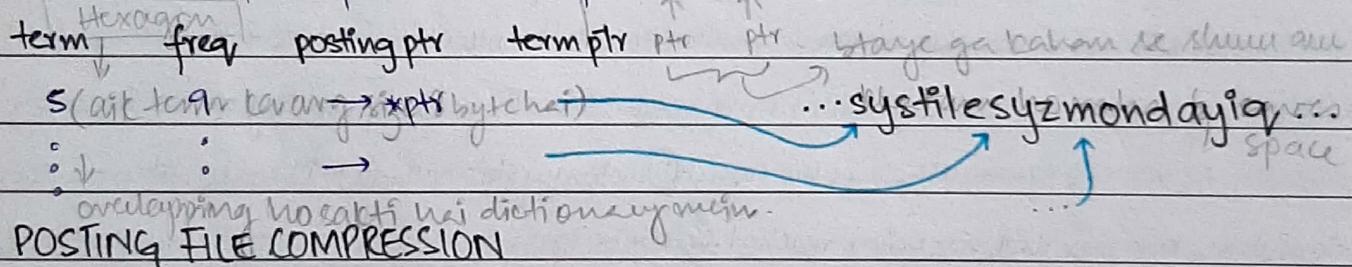
→ term pointer requires $400,000 \times 8 = 3.2 \times 10^6$ positions $\rightarrow \lceil \log_2 3.2 \times 10^6 \rceil = 22$ bytes

↳ kahan aik word khatam aur doosra shuru

so for Reuters-RCV1 which 4 bytes for freq and posting pointer, 3 bytes for term pointer, 8 byte for term

$$\rightarrow 400,000 \times (4 + 4 + 3 + 8) = 7.6 \text{ MB.}$$

→ deduced space by one-third ie from 11.2 MB to 7.6 MB.



POSTING FILE COMPRESSION

→ 4 byte docID

→ offline indexing criteria

→ Postings list > dictionary

→ offline indexing criteria

→ since we'll pick one docID at a time, so we'll do increment of docID after scanning

* IT for the query dictionary

• Alternate: Variable Length Encoding we'll do increment of docID after

scanning. * uses fewer bits for encoding short gaps. ↗ ; short codes for small no.

→ for an economical representation of gaps

→ By using maps source symbol to a variable no. of bits.

→ difference b/w docID will be placed so that the gap to reach next docID will be

be $(\text{docID} + \text{difference b/w 2 docIDs}) - \text{hoga}$

* Avg gap for a term is 9, we want to use $\sim \log_2 9$ bits/gap entry

INDEX COMPRESSION → improves time & space efficiency of indexes by amount of disk needed

* amount of info that can be kept in cache, speeding up data transfer from disk to memory

(only 4% size of total size of collection & 10-15% of total size of text in collection) Dazzle

Day / Date:

INDEX CONSTRUCTION

(Chp 4)

→ How to efficiently make index

→ HARDWARE BASICS

While building an IR sys, many decisions are based on characteristics of comp hardware on which sys runs.

① Access to data in memory > Access to data on disk.

② When R/W on disk, it takes a while for disk head to move to the part of the disk where data is located.

This is called seek time; avg is ms. No data being transferred during seek so to maximize data transfer rate, chunks of data will be read together should be stored contiguously on disk. (that is fixed).

③ Transferring blocks of large data is easier than transferring small chunks; buffer → whole block is being R/W

④ Fault tolerance is expensive. RAID is used to avoid fault tolerance (backups)

⑤ Processor is available to process data during disk I/O so we use it to speed up data transfer by storing compressed data on disk.

→ Basic Inverted Index (already studied)

• Block Sorted Based Index (BSBI)

→ If disk read some keyword & do ID buffer me data hain & when buffer full we put it on the disk after sorting. So we have blocks of sorted words sorting. So we have blocks. Then we merge it with help of pointers in such a way that repetition does not take place.

↑ better than previous process be through merge kringge in such a way that

→ use when dataset cannot be kept in memory. Block by block we can easily fetch query result from the disk. This process.

→ use when data set to memory mein mat rakh sakega.

• Single-Pass-In-Memory indexing → fetch the query result from disk.

Key Idea ~ Generate separate dictionaries

Single word already dictionary mein hai to bas usko posting mein + kringge.

BSBI se fast and efficient (hashing ke nahe)

Agar ki word already present hai to bas posting list mein + kringge bas

faster than BSBI b/c soi mai kar nahe (hashing ke nahe)

Day / Date: _____

• Distributed Index (Search Engine do this)

→ Maintain a master machine

① Parallel Task ② Parser ③ Index

Blocks are taken from master machine by various indexing machine
machine takes blocks from indexing machine to select way in a sorted way = sorted way
of taking chunks back from indexing machine

→ all decisions are taken by master machine.

Day / Date: _____

Chp 6

Scoring, TERM WEIGHTING & VECTOR SPACE MODEL (VSM)

Limitations of Boolean Model

- takes a lot of skills to come up with a query that produces a manageable no. of hits.
- and gives too few or too many documents.

Ranked Retrieval

- Result is shown in decreasing order of ranking
- If score is equal, the doc are of both equal importance
- $f = q, d_1) = \{d_1, \dots, d_n\}^{\uparrow}$ → in precomputed way
every collection of doc order → permutation
and we get max no. of result terms main
- Free text queries:
- function of score-ranking
 - max → if identical (not same)
 - min → if totally different
- Score is based on [0, 1] → totally same
totally different

Query-Matching Document

(→ freq. matters)

Assign a score to a query/doc pair.

$$\text{Jaccard Query} \rightarrow \frac{|D_q \cap q|}{|D_q \cup q|}$$

$w_p^q w_q^p$ multiset

$w_p^q w_q^p$ every doc is considered as a bag of words

The more freq. the query term, the higher the score

$$\text{Sim}(d_i, q) \quad \boxed{q = w_q; w_p}$$
$$R = \frac{5}{8}$$

w_i	w_p	w_p
w_p	w_z	w_p
w_q	w_q	

Dazzle

Day / Date: _____

→ limitation

→ favoured for smaller doc ; biased for shorter doc, short doc at top.

Parametric Search / Meta Search

→ common in digital libraries

↓ provide search on basis of parameters.

→ weighted Zone Scoring

$$\sum_{i=1}^n g_i = 1 \quad \text{sum of all scores of zones.} \rightarrow \text{zone is abstract, can be on any basis}$$

$$S(Q, D_i) = \sum_i g_i \times \text{score}(q_i, D)$$

↓
(weight of a zone)

↑ we try to learn weight by supervised learning mechanism

→ freq. plays a role but not this algo

(partial ordering) (Search) B
(+ ranking + arrangement of result)

$$\text{Train} \rightarrow f(x) = y$$

$$\text{Test} \rightarrow f(x_t) = y \rightarrow y$$

Learning weights for zone Scoring

use
query doc rank → explicit relevance feedback

q₁ 0 → use to query all doc which are relevant to
ranking for query q₁

multiple user to overcome biasness.

(supervised learning)

x → y
humanly annotated
work → logistic judge

Training data
Testing data

labeled supervised data → Dazzle

$$x_i = \{f_1, \dots, f_n\} = R/R$$

Day / Date: _____

Learning weights for Scoring zone

$$Q = t$$

$$D = d_i = [\begin{array}{c} \text{title} \\ + \\ \text{body} \end{array}]$$

$$S(q_j, d_i) = g S_T(t_i, d_i) + (1-g) S_B(t_i, d_i) \leftarrow \text{how to determine } g \text{ (constant)}$$



$$\varepsilon(g, \phi_j) = (r(d_i, q_j) - \text{score}(d_i, q_j))^2$$



$$\sum \varepsilon(g, \phi_j) \leftarrow \text{need to sum all}$$

↓ jithna kam error bane wala optimized sol.

$n_{ij}/N_{01} \rightarrow$ any term

$n_{10}/N_{01} \rightarrow$ title main present body main nei

$r \rightarrow$ relevant

$n \rightarrow$ non-relevant
(in subscript)

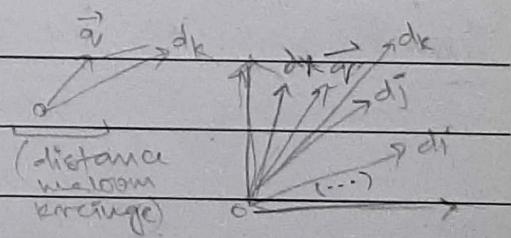
$$\frac{n_{10r} + n_{01r}}{n_{10r} + n_{10n} + n_{01r} + n_{01n}}$$

Bag of words

Score = Common terms in doc / query over total terms (use jaccard query)

Issues -

① Doc. size affects answer



$$\vec{a} = c_1 \hat{n} + c_2 \hat{y} + c_3 \hat{z}$$

Avg.

Model for repr. tex doc as vectors of identifiers
unique.

Every term considered as a dimension;

terms bi binary

dictionary is in order → dictionary is of d dimension. $\frac{1}{2}$ d dimension hyper space.

Used in info filtering, info retrieval, index

Day / Date: _____

Vector representation of Query

- Query (\vec{q}) is also represented similar to doc.
- Use distance b/w vector as a measure of closeness
- N-dimensional vector is also used for query only the term present in query are non-zero vector for query

$$\cos(d_i, q) = \frac{\vec{d}_i \cdot \vec{q}}{\|\vec{d}_i\| \cdot \|\vec{q}\|}$$

Cosine similarity

- ① If angle close, doc is closer (similar) [angle b/w q & doc]
- vector distance reject b/c sensitive

- ② hal doc to unit vector se represent

$$\begin{aligned} \vec{v} &= \vec{v}/\|\vec{v}\| \\ \vec{u} &= \vec{u}/\|\vec{u}\| \end{aligned}$$

$$\cos(\theta) \approx 0$$

→ jo term query mein aa rahi hai, sirf uske liye calculate kringa.

Term Frequency

→ How many times a term has appeared in a doc.

$$t_{(d,t)} \Rightarrow \text{love} (1,3)$$

doc ↓ ↑
no. freq

Sublinear 'tf' scaling

→ Term freq. can be scaled in different ways (tf)

* log linear $\rightarrow 1 + \log(tf_{t,d})$

* fragmented \rightarrow weighting scheme $\rightarrow 0.5 + \frac{0.5 \times tf_{t,d}}{\max_t (tf_{t,d})}$

Dazzle

Day / Date: _____

* Maximum Freq. normalisation (ha k term aur uska max)

Doc freq.

→ collection ke kitney doc mein aik term aya hai

$df = 3 \rightarrow 3$ doc mein aya

↓ (aka. inverse doc. freq.)

① Scaling the df : $idf_t = \log \frac{N}{df_t}$ → wo term jo kam doc mein aya hai wll have higher idf — discriminating power.
• agar koi term kab doc mein hai to no advantage in retrieval.

② Augmented — $\max\{0, \log \frac{N - df_t}{df_t}\}$ { → sawal doc ko weight de takay they can partially affect score }

Best weighting Scheme

$$= (\text{term freq}) \times (\text{inverse doc freq})$$

$$= \frac{tf}{t,d} \times \frac{idf}{t,d}$$

have ↓
 have dimension doc $\frac{1}{d}$
 mein aik term aya hai,

assigns a weight to term

① highest when t occurs many times within a small no. of doc.

② lower when term occurs few times in a doc or occurs in many doc

③ lowest when term occurs in almost all doc.

3 imp properties

If we are keeping doc freq. = 1 so means that we have not considered any doc's freq.

SMART SYSTEM — 60 possible config.

(\downarrow add. query)

- Advantages of VSM
- Disadvantages of VSM

Day / Date: _____

singular similarity

Ex: 1 → 6010

→ BDCxidfe

{ information retrieval }

simple

stable + efficient stable sort

n log n

a	doc1	after	a	doc1
b	doc1	⇒	b	doc2
c	doc1	sort	c	doc3
:			:	
z	doc2		z	doc2
:			;	doc4
:				

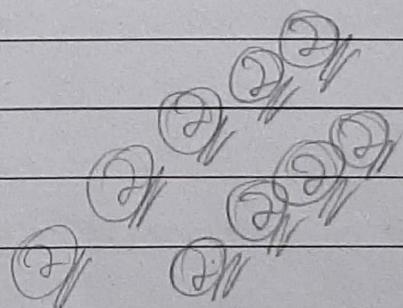


efficient
way is

to pack

a doc
BET & after
buffer is
filled shift

to disk
in blocks



Day / Date: _____

- BOOK NOTES -

(Chpt 1)

→ Building an Inverted Index

→ (to gain speed benefits of indexing at retrieval-time, build index in advance:)

① Collect docs to be indexed.

Friends, Countrymen, Roman So let it be Caesar...

② Tokenize text by turning each doc in tokens.

list of Friends Countrymen Roman So let it ...

③ Do linguistic preprocessing by producing a list of normalized tokens.

Friend Roman Countryman ..

④ Index the doc that each terms occurs in by creating an inverted indexing consisting of dictionary & posting list.

Brutus	1	2	4	11	...
Caesar	1	2	4	5	...
:					

* Each doc has a unique docID.

* dictionary also stores doc freq = length of posting list.

↑
increases efficiency of search engine
at query-time + used for ranked retrieval model.

dictionary

posting list

→ This is then later sorted & grouped to make an index.

→ The dictionary stores the terms & has a pointer to the posting list for each term.

→ Inverted Index is the most efficient structure for supporting adhoc text search.

* We pay for storage of both; dictionary kept in memory while posting list in disk (since large)
what data structure should be used for posting list?

- ① singly linked list ② variable length array

→ Processing Boolean Queries

* Merge Algo

→ Brutus AND calpurnia

Intersect (p1, p2) {

① locate Brutus from dictionary

answer ← ()

② Retireve its posting list

while p1 ≠ NIL and p2 ≠ NIL

③ locate calpurnia in dictionary

do if docID(p1) = docID(p2)

④ Retireve its posting list

then ADD (answer, docID(p1))

⑤ Intersect the two posting list

p1 ← next(p1)

Brutus → 1 → 2 → 4 → 11 → 31 → 45 → 173 → 174

else if docID(p1) < docID(p2)

Calpurnia → 2 → 31 → 54 → 101

then p1 ← next(p1)

Intersection → 2 → 31

else p2 ← next(p2)

Dazzle
return answer ?

Day / Date: _____

(Chp 3)

→ SPELLING CORRECTIONS

① EDIT DISTANCE

2 character strings s_1 & s_2 ; edit distance b/w them is minimum no. of edit operations required to transform s_1 into s_2 .

common edit operations allowed:

- ① insert a character in string
- ② delete a character
- ③ replace a character with another.

edit distance aka Levenshtein distance

e.g. cat & dog ; distance is 3.

→ Weighted Edit Distance

→ B/w two strings in time $O(|s_1| \times |s_2|)$; use dynamic programming.

→ char. in s_1 & s_2 are in array form. Algo fill the (integer) entries in a matrix m whose dimensions equal the lengths of two strings whose edit distance is being computed.

→ (i, j) entry of matrix will hold the strings consisting of first i char of s_1 & first j char of s_2 .

EDIT DISTANCE(s_1, s_2)

```
int m[i][j] = 0  
for i ← 1 to |s1|  
do m[i, 0] = i  
for j ← 1 to |s2|  
do m[0, j] = j  
for i ← 1 to |s1|  
do for j ← 1 to |s2|
```

do m[i, j] = min {min[i-1, j-1] + if ($s_1[i] = s_2[j]$) then 0 else 1 if
m[i-1, j] + 1,
m[i, j-1] + 1}

→ The spelling correction demands more than

computing edit distance.

↓
This exhaustive search is expensive.

return m[|s₁|, |s₂|]

Dazzle

Day / Date: _____

② K-gram Indexes

→ to further limit the set of vocab terms for which we compute edit distances to query term.

↓ K-gram: to assist with retrieving vocab terms with low edit distance to query q.

x — x

Chp 54

→ BLOCKED SORT-BASED INDEXING (Excellent scaling properties, but needs a DS for mapping terms to termIDs. For large collection, DS doesn't fit in memory)

→ to make index construction more efficient, represent terms as term IDs (where each is unique sno.)

→ build mapping from terms to termIDs on the fly while we are processing the collection.

→ Ruter-RCV1 collection as our model collection — 1GB text → 800,000 documents

BSBI:

BSBIIndex()

(i) segments the collection into parts of equal size.

$n \leftarrow 0$

(ii) sorts the term ID - doc ID pairs of each part in memory.

while (all doc have not been processed)

(iii) stores intermediate sorted results on disk

$do \# \leftarrow n + 1$

(iv) merges all intermediate results into final index

$block \leftarrow \text{ParseNextBlock}()$

BSBI-Invert(block)

① parse doc in termID-docID pairs & accumulates the pair in memory until a block of fixed size is full.

WriteBlockToDisk(block, fn)

MergeBlocks(f₁, ..., f_n; fmerged)

② we choose block size to fit comfortably into memory to permit a fast in-memory sort.

③ Block is then inverted & written to disk; inversion has 2 steps

(i) we sort termID-docID pairs (ii) we collect all termID-docID pairs with same termID into a posting list where posting is simply a docID.

④ The result is then written to the disk

→ with RCV1 we end up with 10 blocks, each an inverted index of one part of the collection.

⑤ Finally, we merge 10 blocks into 1 large merged index.

Time Complexity → O(T log T)

Dazzle

Day / Date: _____

→ SINGLE-PASS IN-MEMORY INDEXING

More scalable terms ^{alternative} is SPIMI. → uses terms instead of termIDs, write each block's dictionary to disk & then starts a new dictionary for the next block.

→ SPIMI can index collections of any size as long as there is enough disk space available.

→ SPIMI - Invert is called repeatedly on token stream until the entire collection has been processed.

→ Tokens are processed one by one during each successive call of SPIMI - invert.

→ When a term occurs for first time → added to dictionary (best as hash), and a new posting list ^(line 6) created. → the call in line 7 returns this posting list for subsequent occurrences of the term.

→ Difference b/w BSBI & SPIMI → SPIMI adds a posting directly to its posting list (line 10).

Each posting list is dynamic & is immediately available to collect postings.

Advantages:

① faster b/c there is no sorting required

② saves memory b/c we keep track of ~~to~~ term a postings list belongs to so termID doesn't need to be stored.

→ whole lot process is more efficient

Three components of SPIMI:

① Constructing new dictionaries for each block

② Eliminates expensive sorting

③ Compression; posting & dictionary can be

stored compactly on disk if we

employ compression.

(increase efficiency)

TIME COMPLEXITY : $O(T)$ b/c no sorting &

almost all operations are linear.

Day / Date: _____

→ DISTRIBUTED INDEXING

- Web scale search engine cannot build index on a single server — Large data.
- The result of construction process is a distributed index partitioned ^{across} on the several machines — either acc. to term or acc. to doc → we describe for acc. to term

- Application MapReduce → designed for large comp clusters.

- A master machine (node) directs the process of assigning task to individual nodes
 - map & reduce phases of MapReduce split up computing job into chunks that standard machine can process in short time.

- (1) Input data (collection of web pgs) split into n splits where the size of the split is chosen to ensure that the work can be distributed evenly & efficiently (16 or 64 MB size)
 - Splits are not assigned, they are assigned by master node on an ongoing basis.
 - As a machine finishes processing one split, it is assigned another.
 - If a machine dies or become lagged, the split it is working on is assigned to another machine
 - 'key-value pair' is → term → termID mapping.

Map phase:

- Mapping splits of input data to key-val pair. Same parsing task as BSB1 & SPIMI, so called parser
 - machines are
 - each parser writes its output to local intermediate files, the segments files

a-f	g-p	q-z
-----	-----	-----

Reduce Phase:

- we want all values for a given key to be stored close together → read & processed quickly.
 - Achieved by partitioning key into j term partitions & having parser write key-val for each term partition into a separate segment file
 - j=3 ←
- Collecting all val(docID) for a given key (term ID) into one list is task of inveter.
 - Master assigns each term partition to a different inveter.
- same machine can be parser in map phase & inveter in reduce phase depends on if idle.
- Each parser writes its segment file to its local disk; minimize amount of traffic in indexing.

Conclusion

- simple & robust framework

Dazzle

Day / Date: _____

→ DYNAMIC INDEXING

Daz

Day / Date: Someone!

q

dot

d

$$d_1 q_1 + d_2 q_2 + \dots$$

$$\sqrt{d_1^2 + d_2^2 + \dots}$$

vocab = $\{w_1, w_2, w_3, w_4, w_5\}$

$$d_1 = \{w_2, w_3, w_1, w_3, w_2, w_1, w_5\} = 5$$

$$d_2 = \{w_3, w_2, w_1, w_4\}$$

$$d_3 = \{w_1, w_3, w_2\}$$

$$q = \{w_2, w_1, w_3\} = 2$$

Similarity \Rightarrow

Similar df \rightarrow Similar

fauto, company, moto?

$$q = tf * idf$$

$d_1 = \{a\}$

$d_2 = \{c, m\}$

$$d_1 \langle 1, 0, 0 \rangle$$

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

$$d_1 \langle 1, 2, 2, 0, 1, 1 \rangle$$

$$d_1 \langle 0, -0.24, -0.24, 0.18, 0.48 \rangle$$

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

$$\text{sim}(d_1, q) = \vec{v}(d) \cdot \vec{v}(q)$$

$$d_3 \langle 1, 1, 1, 0, 0, 0 \rangle$$

$$|\vec{v}(d)| |\vec{v}(q)|$$

$$q \langle 2, 1, 0, 0, 0 \rangle$$

$$N = 3$$

tf

$$idf = \log \frac{N}{df}$$

$$\log \frac{N}{df}$$

ayesha \Rightarrow

seriously

~~d1 d2~~ N df

$$\log \frac{3}{3} = \log 1 = 0$$

w1

31

[0]

tf*idf

w2

4

-0.12

w3 Monday

3

-0.12

summer

w4

2

0.18

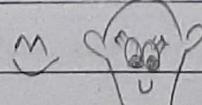
w5

1

0.47

sunday! \approx

ayesha



J(C)

ayesha.

ayesha.

$$JC(d_1, q) = \frac{d_1 \cdot q}{d_1 \cdot q} = \frac{2}{5}$$

Day / Date: _____

$$d_1 = \{ w_1, w_2, \underline{w_3}, w_4, \underline{w_5}, w_6 \}$$

$N=3$

$$d_2 = \{ w_1, w_2, w_6, w_4 \}$$

$$TF = 1 + \log(t_f)$$

$$d_3 = \{ w_2, w_5, w_6, w_7 \}$$

$$IDF = \log\left(\frac{N}{df_t}\right)$$

$$q = \{ w_5, w_4, w_7 \}$$

$$\boxed{tf * IDF}$$

$$IDF \quad tf - D1 \quad tf - D2 \quad tf - D3 \quad tf - q \quad EF * IDF - D1$$

$$w_1 \quad 0.176 \quad 1 \quad 1 \quad 0 \quad 0 \quad 0.176$$

$$w_2 \quad 0 \quad 1 \quad 1 \quad 0 \quad 0$$

$$w_3 \quad 0.77 \quad 1.301 \quad 0 \quad 1 \quad 0$$

$$\vec{v}(D) \cdot \vec{v}(q)$$

$$w_4 \quad 0.176 \quad 1 \quad 1 \quad 0 \quad 1$$

$$|\vec{v}(D_2)| \cdot |\vec{v}(q)|$$

$$w_5 \quad 0.176 \quad 1 \quad 0 \quad 1 \quad 1$$

$$\frac{\vec{v}(D_2) \cdot \vec{v}(q)}{|\vec{v}(D_2)| \cdot |\vec{v}(q)|}$$

$$w_6 \quad 0.176 \quad 0 \quad 1 \quad 1 \quad 0$$

$$w_7 \quad 0.477 \quad 0 \quad 0 \quad 1 \quad 1$$

$$\text{sim}(D_1, q) \quad \text{sim}(D_2, q) \quad \text{sim}(D_3, q)$$

$$w_1 \quad D_2 q_V +$$

$$w_2$$

$$w_3$$

$$w_4 \quad \text{sim}(D_1, q) = \frac{(0.176)^2 + (0.176)^2}{(0.176)^2 + (0.62)^2 + (0.176)^2 + (0.176)^2}$$

$$w_5 \quad \left(\sqrt{(0.176)^2 + (0.62)^2 + (0.176)^2 + (0.176)^2} \right)$$

$$w_6$$

$$w_7$$

Day / Date:

EVALUATION IN INFORMATION RETRIEVAL

(Chp 8)

* Intro:

8.1 — Information Retrieval System Evaluation

To measure adhoc info retrieval effectiveness, we need test collection consisting of 3 things:

① A doc collection. ② A test suite of info needs, expressible as queries.

③ A set of relevance judgements, standardly a binary assessment of either relevant/non-rel. for each query-doc pair.

> Wrt user info need, a doc in collection is given a binary collection as either rel/nonrel. This decision is referred as gold standard or ground-truth judgement of relevance.

> Relevance is assessed relative to an info need, not a query.

→ Many sys contain weights (known as parameters) which can be adjusted to tune sys performance

↓
It is wrong to ~~test~~ ~~report~~ results on a test collection which were obtained by tuning these parameters to maximize performance on collection; we should have development test collection to tune the parameters dev test collection. The tests run those weights on test collection gives a result ~~which~~ as unbiased estimate of performance.

8.2 — Standard Test Collections

→ Intro: Different IR Models:

Many retrieval models are present, the best component for (Ranking func → dot-mod, cosine), Term selection → stopword removal, stemming) & (term weighting → TF, TF-IDF).

→ Difficulty in IR models:

Effectiveness is related to relevancy of retrieved items; typically binary.

8.2 — Standard Test Collections

→ Cranfield → TREC (Text Retrieval Conference) → GOV 2 → NTCIR
→ CLEF → Reuters → 20 News group

(Read details from book 154)

Day / Date: _____

8.3 - Evaluation of unranked retrieval sets

> Most basic & frequent IR effectiveness measures are: ① precision ② recall

Precision - fraction of retrieved doc that are relevant.

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

Recall - fraction of relevant doc that are retrieved

$$R = \frac{\#(\text{relevant items retrieved})}{\#(\text{relevant items})} = P(\text{retrieved} | \text{relevant})$$

∴ $P = tp / (tp + fp)$ where, $tp = \text{rel. + retrieved}$, $fp = \text{non-rel + retrieved}$
∴ $R = tp / (tp + fn)$ where, $fn = \text{not retrieved + rel.}$

* since result of a query set is unranked (flat), we use $P \& R$ to evaluate.

		Rel	Not Rel		
		Ret	Not Ret	$P(0,1)$	$P = \frac{\text{Ret}_{\text{rel}}}{\text{Ret}_{\text{rel}} + \text{Ret}_{\text{not rel}}}$
Retrieved	Ret	Ret _{rel}	Ret _{not rel}	$R(0,1)$	$R = \frac{\text{Ret}_{\text{rel}}}{\text{Ret}_{\text{rel}} + \text{Ret}_{\text{not rel}}}$
	Not Ret	Not Ret _{rel}	Not Ret _{not rel}		

- an alternative to judge an IR sys by its accuracy is fraction of its collection that are correct.
> Accuracy is not a good measure for IR prob b/c data is extremely skewed (many times)
> Precision & recall concentrate the evaluation on the return of true +ve.

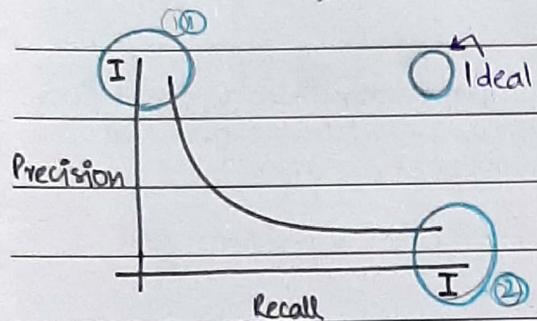
Precision vs. Recall

sys makes 2 types of errors:

- false +ve → retrieves a doc that is non-relevant (should not have been retrieved)
false -ve → fails to retrieve a doc that is relevant (" have been retrieved)

Day / Date: _____

> false +ve affects \rightarrow precision ; false -ve affects \rightarrow recall



① Retrieves relevant doc but misses many useful too

② Retrieves most relevant doc but includes lots of junk

> Precision & Recall are inversely proportional.

• In some cases, precision matters more than recall & vice versa.

> **F-measure** — A single measure which trades off precision vs. recall.

F-measure is the weighted harmonic mean of precision & recall.

$$F = \frac{1}{\alpha \frac{1}{P} + (1-\alpha) \frac{1}{R}} = \frac{(B^2+1)PR}{B^2P+R} \quad \text{where } B^2 = \frac{1-\alpha}{\alpha} \quad \text{where } \alpha \in [0,1] \text{ & } B \in [0,\infty]$$

→ Default F measure equal weights P & Recall; $\alpha=1/2$ or $B=1$ if $F_{B=1} = \frac{2PR}{P+R}$.

→ However if $B < 1$ → emphasize P & if $B > 1$ → emphasizes R.

→ recall, precision & f-measure are b/w 0 & 1.

Precision Critical Tasks

→ Time matters a lot, Tolerance to missed doc, Redundant: many equal info resources

Eg: Web search & demand is very high. General Optimization

Recall Critical Tasks

→ Time matters less, Non-tolerance to missed doc, Less Redundant info: only one or few resources

Eg: Legal/patent search & moderate demand. Specific optimization.

> $F_{B=1}$ is called harmonic mean of P & R.

↓ Harmonic mean is advantageous to use b/c it is more closer to $\min(P, R)$, whereas arithmetic mean is closer to the highest of two values of P & R.

Day / Date: _____

Fallout Rate

Problems with Precision & Recall:

- ① No. of DNR not taken into account. ② Recall undefined → no relevant doc in collection
 - ③ Precision undefined → no doc is retrieved.

Fallout = $\frac{\text{No. of DNR retrieved}}{\text{total no. of DNR in collection}}$

8.4 Evaluation of ranked retrieval results

- P, R & F-measure are set based measures & computed using unordered set of doc.
 - In a ranked retrieval context, appropriate sets of retrieved doc are naturally given by top k retrieved doc.
 - For each ^{such} set, precision & recall values can be plotted to give a precision-recall curve.
 - (P/R graph)
 - (used computer ranked result)

→ Saw tooth shape.

→ If $(k+1)^{th}$ doc retrieved as non-rel. then recall is same as for the top k doc, but precision dropped.

→ If relevant, then both $P \& R$ increase & curve jags up and to the right.

Interpolated precision P_{interp} is used to remove giggles; P_{interp} at a certain R level r is defined as the highest precision found for any recall level $r' \geq r$:

$$P_{\text{mleap}}(r) = \max_{r' \geq r} p(r')$$

n	doc #	relevant	#D _i =6	R = 1/6 = 0.176	P = 1/1 = 1
1	588	x			
2	589	x		R = 2/6 = 0.333	P = 2/2 = 1
3	576				
4	590	x		R = 3/6 = 0.5	P = 3/4 = 0.75
5	620				
6	990	x		R = 4/6 = 0.667	P = 4/6 = 0.667
⋮	⋮	⋮	⋮	⋮	⋮
13	772	x		R = 5/6 = 0.833	P = 5/13 = 0.385

Day / Date: _____

Average Precision

█ → relevant docs → total = 6.



R	0.17	0.17	0.33	0.5	0.67	0.83	0.83	0.83	0.1	
P	1.0	0.5	0.67	0.75	0.8	0.83	0.71	0.63	0.56	0.6

$$\text{Avg ranking} = (1.0 + 0.67 + 0.75 + 0.8 + 0.83 + 0.6) / 6 = 0.78$$

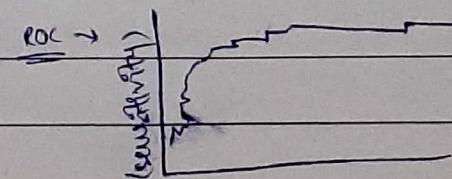
Mean Average Precision (MAP)

- > provides single-figure measure of quality across recall levels.
- > MAP shows good discrimination & stability; good for web search
- for single avg precision is the avg of precision val. obtained for the set of top k docs existing after each rel. doc is retrieved, and this val. is then averaged over info needs.
- > If rel. doc never gets retrieved, we assume precision corresponding to that rel. doc = 0.
- > MAP is macro-averaging - each query counts equally.
- > MAP assumes user is interested in finding rel. doc for each query.
- > MAP requires many relevance judgements in test collection.
- has adv. of not requiring any size of doc set; drawback: less stable, not good avg.
- Precision at k : measuring P at low levels of retrieved result.
- R-Precision → requires set of known relevant doc returned from which we can calculate the precision of top D_r.

Break-even point → when R-Precision & the recall are same $\frac{1}{D_r + N}$; r → relevant.

ROC - Curve - Receiver Operating Characteristics

- > Plots true positive rate or sensitivity against false pos rate or (1 - specificity); sensitivity = recall
- > false pos rate is given by $fp / (fp + tn)$



(1-specificity)

Dazzle

Day / Date: _____

Kappa Statistics (from slide 14 W7)

$$K = (P(A) - P(E)) / (1 - P(E))$$

; where $P(A)$ = Observed proportion of the times judges agreed
 $P(A) = \frac{\text{YES} + \text{NO}}{\text{YES} + \text{NO} + \text{NO}} / \text{total}$

(YES / YES)

(NO / NO)

0.8 \Rightarrow good agreement

- $K \approx 0.0 - 1$; $0.6 - 0.8 \rightarrow$ fair agreement
(consistently agree with main judges)

$$P(E) = P(\text{non-relevant})^2 + P(\text{relevant})^2$$

\hookrightarrow proportion of times the judges would be expected to agree by chance

\rightarrow A common measure for agreement b/w judges is the kappa statistics.

\rightarrow Cumulative Gain (CG)

Sum of graded relevance val of all results in a search result list; judges score the doc.

Ex: for query $\rightarrow q$, $D_1 D_2 D_3 D_4 D_5 D_6$ are relevant docs.

$$3 \quad 2 \quad 3 \quad 0 \quad 1 \quad 2 \quad \Rightarrow [CG = 11]$$

\rightarrow Only change in doc brings a change in CG ;

Drawback - changing of order doesn't affect CG measure.

\rightarrow Discounted CG (DCG)

DCG is total gain accumulated at a particular rank p: $DCG_p = \text{rel}_1 + \sum_{i=2}^p \frac{\text{rel}_i}{\log_2 i} \text{ or } \sum_{i=1}^p \frac{2^{\text{rel}_i} - 1}{\log(1+i)}$

\rightarrow emphasizes on retrieving highly relevant docs.

\rightarrow If $\text{rel} = 1$; order = $\log(1+i) \rightarrow DCG = 3$ (slide 16 of Week 7 - solved example)

\rightarrow Normalized DCG at rank n

NDCG at rank n by DCG value at rank n of the ideal ranking

Ideal result \rightarrow sorted in decreasing order of doc relevancy.

Useful \rightarrow for contrasting queries with varying no. of relevant results; Used in web search.

$$nDCG = \frac{DCG}{\text{Ideal DCG}}$$

(see slide 18 for more solved sol - Week 7)

Day / Date: _____

8.5 - Assessing Relevance

- > To properly evaluate sys, given info needs & doc needs, you need to collect relevance assessments
- > Standard approach is **pooling** → relevance is assessed over a subset of collection that is formed from the top k doc returns by a no. of different IR sys

8.6 - A broader perspective: System quality & user utility

Systems that allow quantitative evaluation & the issue of user utility.

- System Issues

retrieval

Several benchmarks to rate an info ~~inform~~ system beyond its retrieval quality, they are:

- ① How fast does it index, how many doc/hr it index over doc length? [measurable & can quantify size or speed]
- ② How fast does it search ie what is latency as a func. of index size?
- ③ How expressive is it query lang. & how fast is it on complex queries?
- ④ How large is doc collection?

- User Utility

- Way of quantifying aggregate user happiness based on relevance, speed & user interf. of sys.
- Take survey from user (feedback); if user → happy, it will use sys again. (return to same engine)

- Refining a deployed system

• A/B testing

- > Precisely one thing is changed b/w the current sys & proposed sys, and a small proportion of traffic (1-10% users) ^{is} randomly directed to the variant sys, while most users use current sys.
- > click through log analysis or click stream mining. To see whether user like A or not.
- > The basis of A/B testing is running a bunch of single variable tests (either in seq or in parallel): for each test only one parameter is varied from the control (i.e. current sys)
- > A/B finds it challenging to which part it should evaluate.

Day / Date: _____

8.7 — Result snippets

> Search snippet is useful for reviewing the search results.

> Two kind of summaries:

Static — which are always the same regardless of query → drawback
— fast computing → (advantage)

Dynamic — (aka Query dependant) which are customized acc. to the user's info need as deduced from a query. ~~Dynamic~~ Dynamic summaries attempt to explain why a particular doc was retrieved for the query at hand.

* where the word appears in doc, it takes 100 words before & after the doc and processes a summary of it.

- Q dependant → user explains query very well.
- Advantage → give user full idea of doc.
- Drawback → time consuming eg. google.

→ Keyword -in- context (KwIC) snippets

Dynamic summaries display one or more "windows" on doc, aiming to present the pieces that have the most utility to user in evaluating doc w.r.t user's need of info

→ usually these windows contains one or several of q terms, and so are referred to as keyword -in- context (KWIC) snippets.

Note: find everything relevant — high recall

find only retrieve what is relevant — high precision.

Day / Date: _____

RELEVANCE FEEDBACK & QUERY EXPANSION

(Chp 9)

- PROBLEM with IR System

- Also known as Synonymy → has an impact on recall of most IR systems.
- * users address this prob by refining query themselves ; Query refinement / Expansion.
- Method for tackling is → ① Global Methods ② Local Methods
- Global Method: Reformulating query terms independent of the query & result returned from the query. Changes will cause it to match other semantically similar terms.
 - Using Thesaurus → via automatic Thesaurus generation, spelling corrections
- Local Method: Adjust query relative to a doc that initially appear to match the query.
 - Relevance Feedback → involve user feedback from results return for a query.
 - Pseudo Relevance Feedback → blind relevance feedback ; • Global (Indirect) R feedback.

9.1 - Relevance Feedback & Pseudo Relevance Feedback.

- Idea of relevance feedback → involve user in the retrieval process as to improve final result set
- In particular, user gives rf on initial set of docs ; Rf can go through a no. of iterations.
- Basic procedure:
 - ① user issues a simple/short query.
 - ② sys returns initial set of retrieval result
 - ③ user marks doc as R or NR.
 - ④ sys computes a better representation of info need based on user feedback.
 - ⑤ sys displays a revised set of retrieval result.
- * we assume user is not able to formulate query properly so we involve him to get perfect result.

9.1.1 - Rocchio Algorithm

- classic algo for implementing relevance feedback ; model of a way to incorporate rf into in vector space model (VSM).
- compute new query in VS.
- need to find query vector, $\text{dist}(q'_v)$ that maximizes similarity with R doc & min. similarity with NR.

Dazzle

centroid \rightarrow vs ka common point.

$$\vec{m}(C) = \frac{1}{|C|} \sum_{d \in C} \vec{d}$$
 where, $C = \text{set of docs}$.

Day / Date:

(sum of all vectors) / no. of doc marked R by user)

\vec{q}_v = query vector

$\alpha > 0 \rightarrow$ user have good judgement of initial query

$C_r = \text{set of relevant docs}$

$D_r \geq S$ - marked by user then

$C_{nr} = \text{set of NR docs}$

$B > 0 \& B > Y$

$\alpha \rightarrow$ based on initial q

$B \rightarrow w \text{ of } D_r, Y \rightarrow w \text{ of } D_{nr}$

-ve weight \rightarrow if $B < Y$

\hookrightarrow term ignored

find, $\vec{q}_{opt} = \arg \max_{\vec{q}} [\text{sim}(\vec{q}_v, C_r) - \text{sim}(\vec{q}_v, C_{nr})]$

under cosine similarity, $\vec{q}_{opt} = \frac{1}{|C_r|} \sum_{d_j \in C_r} \vec{d}_j - \frac{1}{|C_{nr}|} \sum_{d_j \in C_{nr}} \vec{d}_j$

$$\vec{m}(C) = \frac{1}{|C|} \sum_{d \in C} \vec{d}$$

\leftarrow optimal query is vector difference b/w The centroids of the R & NR doc.

(\hookrightarrow set of docs)

(not precise cause we do the full set of R doc)

push towards +ve

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

\downarrow pulls away from -ve

where, \vec{q}_m = modified query, \vec{q}_o = original query; $\alpha, \beta, \gamma \rightarrow$ weights of terms (cause user gave ambiguous)

\rightarrow query more toward R doc & away from NR.

\rightarrow tradeoff α vs. β/γ = if we have a lot of judged docs, we'd want a higher $\beta \& \gamma$.

\rightarrow Negative term weights are ignored i.e. term weight is set to 0.

\rightarrow Rf can improve both precision & recall. is most useful for increasing recall in certain situations

9.1.2 Probabilistic Relevance Feedback

Rather than reweighting query is VSM, we can use Naive Bayes on the basis of R & NR doc.

$$P(x_t=1|R=1) = |V_{Rt}|/|VR|$$

R \rightarrow relevant or not.

$$P(x_t=1|R=0) = (df_t - |V_{Rt}|)/N - |VR|$$

N \rightarrow total no. of docs

df_t \rightarrow no. that contains t.

VR \rightarrow set of known relevant docs

x — x

Explicit fd \rightarrow user ~~doesn't~~ know that he is given feedback

Implicit / Indirect fd \rightarrow user doesn't know; less reliable but more useful than Dazzle pseudo fd.

Day / Date: _____

9.1.3 - When does Relevance Feedback Works?

Success of Rf depends on few assumptions :

- A1 : User has sufficient knowledge for initial query (user formulated a good q i.e. $\alpha > 0$)
- A2 : Distribution of term in Dr is similar ie. relevance prototypes are well behaved'
 - Term distri. in NR doc will be different from those in Dr.
 - either all doc R are tightly clustered around a single prototype or there is a vocab overlap.

Violations

A(1) : User does not have sufficient initial knowledge.

Eg. (i) Misspellings (Britney Spears) (ii) cross lang. IR (Miguel)

(iii) Mismatch of searcher's vocab vs. collection vocab (cosmonaut / astronaut)

A(2) There are several relevance prototypes. (the doc ki representation bhi hain kisी vocab can move our query there)

Eg. (i) Burma / Myanmar (ii) pop stars that worked at Bregal King.

→ Often: instances of a general concept ; good editorial content can address prob.

9.1.4 - Relevance Feedback On Web

Some web search engines offer a similar/related pgs. feature - the user indicates a doc in the results set as exemplary from standpoint of meeting his info need & requests more doc like it.

Eg. Sprinkle et. al (2000) presents result using Excite Search Engine.

9.1.5 - Evaluation of relevance feedback strategies

① Interactive Rf can give very substantial gains in retrieval performance.

→ One round of rf is useful, second round (2) are more useful. At least five judged doc is recommended

② Use doc in residual collection (set of doc minus those assessed relevant) for 2nd round of evaluation

↑ neither of the two is fully satisfactory

∴ best is do a time based comparison i.e. how fast does a user find Rel doc with rf vs another strategy (eg. query reformation)

Dazzle

Day / Date: _____

9.1.6 — Pseudo Relevance Feedback

- Aka blind rf provides a method for ~~to~~ automatic local analysis.
- Automates the manual part of rf so the user gets improved result without an extended interaction.
- Method → do a normal retrieval to find an initial set of top K ranked doc are relevant
do relevance fd as before on this assumption.

9.1.7 — Indirect Relevance feedback (aka implicit rf)

Relevance feedback Problems

- ① long queries are inefficient for typical IR engine
- ② long time taken for response by users
- ③ higher cost for retrieval sys ⇒ Partial sets
- ④ users are often reluctant to provide explicit feedback.
- ⑤ Often hard to understand why a particular doc was retrieved after applying rf

9.2 — Global Methods for Query Reformulation

Expanding query by simply adding aiding the user in doing so, by using Thesaurus.

9.2.1 — Vocabulary Tools for Query Reformulation

IR sys ~~would~~ might suggest search terms by means of a thesaurus or a controlled vocab.

9.2.2 — Query Expansion

- User gives additional input on query words or phrases, possibly suggesting additional query terms (good/bad search term)
- Query expansion help in extending user queries with related terms in order to solve the lexical gap prob in IR
- Most common form of query expansion is global analysis using some thesaurus.
- * Use of Thesaurus can be combined with ideas of term weighting.

Day / Date: _____

- How to augment user query?

→ Methods for building a thesaurus for query expansion includes:

- ① Use of controlled vocab is maintained by human editors (can be query rather than just synonyms)
- ② An automatic derived -thesaurus
- ③ Query Reformulations based on query log mining.

→ 'Thesaurus-based query expansion' has the advantage of not requiring any input.

→ Query exp. increases recall → effective ; may significantly ↓ precision.

→ Global analysis (static of all docs in C)

• Automatically derived thesaurus → generate a thesaurus automatically by analyzing a collection of docs. There are two ways:

- ① * Exploit word cooccurrence (word text statistics) to find most similar words
- ② Use shallow grammatical analysis of text and to exploit grammatical relations

In a nutshell,

Query Expansion is less successful than relevance feedback; however,
it is more understandable to sys user.

Day / Date: _____

Note: This chp is till 11.03

(Chp 11)

PROBABILISTIC INFORMATION RETRIEVAL

11.1 - REVIEW OF BASIC PROBABILITY THEORY

→ Joint probability = $P(A, B)$

→ Relation b/w both & Chain Rule:

→ Conditional probability = $P(A|B)$

$$\left. \begin{array}{l} P(A, B) = P(A \cap B) = P(A|B)P(B) = P(B|A)P(A) \end{array} \right\}$$

→ Partition Rule $\Rightarrow P(B) = P(A, B) + P(\bar{A}, B)$

$$\left. \begin{array}{l} * P(A_1) = P(A_1, B_1) + P(A_1, B_2) \\ * P(A_2) = P(A_2, B_1) + P(A_2, B_2) \end{array} \right\}$$

→ Bayes's Rule $\Rightarrow P(A|B) = \frac{P(B|A)P(A)}{P(B)} = \left[\frac{P(B|A)}{\sum_{x \in \{A, \bar{A}\}} P(B|x)P(x)} \right] \cdot P(A)$

↳ this eq. can be thought of as a way of updating probabilities.

(i) Start off with prior prob. $P(A)$ {initial estimate of how likely event A is in absence of any other info}

(ii) Derive a posterior prob $P(A|B)$ after having seen the evidence B, based on the likelihood of B occurring in 2 cases that A does / doesn't hold.

→ Odds of an event, a kind of multiplier for how prob changes:

$$\text{Odds} - O(A) = \frac{P(A)}{P(\bar{A})} = \frac{P(A)}{1-P(A)} \quad \leftarrow \begin{cases} P(A) \approx \text{odd}(A) \\ \text{success : failure ratio} \end{cases}$$

Why probability in IR?

① → IR deals with uncertain info ; Prob → measure of uncertainty.

② → Probabilistic Ranking Principle → • Provable • Minimization of risk.

③ → Probabilistic Influence → To justify your decision

Basic probability

• Probability is study of randomness & uncertainty.

• favourable events / total event space in consideration.

• no. assigned to each event (subset) of a sample space - n.

- ① * If $A \in \Sigma \rightarrow P(A) \geq 0$
 ② * $A, B \in \Sigma A \subseteq B \rightarrow P(A) \leq P(B) \geq 0$

Note: $E \nmid E'$ are mutually exclusive
 $\therefore P(E) = 1 - P(E')$

Day / Date: _____



Axioms of Probability

- for any event A : $0 \leq P(A) \leq 1$
- $P(\Omega) = 1$ (entire sample space)
- If $A = A_1 + \dots + A_n$, then $P(A) = P(A_1) + \dots + P(A_n)$
- $P(\emptyset) = 0$.

Properties of Probability

- $P(A^c) = 1 - P(A)$
- If $A \subseteq B$, then $P(A) \leq P(B)$
- $P(A \cup B) = P(A) + P(B) - P(A \cap B)$
- $P(A \cup B \cup C) = P(A) + P(B) + P(C) - P(A \cap B) - P(A \cap C) - P(B \cap C) + P(A \cap B \cap C)$

→ Independence $\Rightarrow P(A, B, C) = P(A) + P(B) + P(C)$; all are independent events.

11.2 - THE PROBABILITY RANKING PRINCIPLE

11.2.1 - 1/0 loss case

We assume that for a query q and doc. d in collection, value = 1 iff relevant else val=0.

In prob. model, docs are ranked in decreasing order of their probability, with prob. estimated accurately wrt info need: $P(R=1 | d, q)$. The overall effectiveness of the sys to its user will be the best that is obtainable on the basis of those data.

In simplest case of PRP, there is no retrieval costs or other utility that would differentially weight actions or errors. You lose a point for either returning a rel doc or failing to return a doc. The goal is to return best possible top results as top k doc, for any val. of k (user's choice)

Such binary situation where you are evaluated on your accuracy is called 1/0 loss.

If set of retrieved doc is to be returned rather than ordering, then Bayes Optimal Decision Rule is to be applied. The decision which minimizes the risk of loss \rightarrow return doc than relevant (more likely)

d is relevant iff $P(R=1 | d, q) > P(R=0 | d, q)$

$$P(R|n) = P(x|R) \cdot P(R) / P(n); P(NR|n) = p(x|NR) \cdot P(NR) / P(n) \rightarrow \text{if } P(R|n) > P(NR|n) \rightarrow \text{then } n \text{ is relevant}$$

11.2.2 - The PRP with retrieval costs

Assume model of retrieval cost $\rightarrow c_0 = \text{cost of not retrieving a rel doc} \rightarrow c_1 = \text{cost of retrieving a non-rel doc.}$

for doc 'd' & for all docs $d' \rightarrow$ not yet retrieval

$$c_0 \cdot P(R=0 | d) - c_1 \cdot P(R=1 | d) \leq c_0 \cdot P(R=0 | d') - c_1 \cdot P(R=1 | d')$$

Applications of PRP:

Day / Date:

then it is the next one to be addressed. Such model gives a formal framework where we can model differential costs of false +ve & false -ve & even sys performance at modelling stage.

11.3 – BINARY INDEPENDENCE MODEL

→ BMI has been traditionally used in PRP.

→ Assume: Binary = boolean ; doc & queries both represented as binary ~~term~~^{term} incidence vectors.

* doc d is represented by vector $\vec{x} = (x_1, \dots, x_m)$, where $x_t=1$ if term t present in doc else $x_t=0$

* η is represented by \overrightarrow{q} .

* Independence means that terms are modelled as occurring in doc independently - No association

* The representation is an ordered set of Boolean variable.

* Many does can have same vector representation with their simplification.

~~→ we assume that relevance of each doc is independent of relevance of other doc~~

→ Under BM1 we model the prob $P(R|d, q)$ that a doc is relevant via prob in terms of term incidence vectors $P(R | \vec{n}, \vec{q})$. Then using Baye's Rule:

$$P(R=1 \mid \vec{x}, \vec{q}_v) = \frac{P(\vec{x} \mid R=1, \vec{q}_v) \cdot P(R=1 \mid \vec{q}_v)}{P(\vec{x} \mid \vec{q}_v)} \quad \begin{matrix} \text{if rel or nrel is} \\ \text{retrieved} \\ \text{then representation} \\ \text{of doc is } \vec{x} \end{matrix}$$

$$\text{also, } P(R=1 | \vec{q}_h, \vec{q}_v) + P(R=0 | \vec{q}_h, \vec{q}_v) = 1$$

$\times \text{---} \times$, continued after deriving P.R.P.

11.3.1 — Deriving a ranking function for query terms

We wish to order returned doc by descending \rightarrow Under BMF, it is modeled as $P(R=1 | \vec{x}, \vec{y})$.

$$O(R|\vec{u}, \vec{q}) = \frac{P(R=1|\vec{u}, \vec{q})}{P(R=0|\vec{u}, \vec{q})} = \frac{P(R=1|\vec{q}) \cdot P(\vec{u}|R=1, \vec{q})}{P(R=0|\vec{q}) \cdot P(\vec{u}|R=0, \vec{q})} \quad [\text{relevance odd}]$$

↑
rank difference

(odd main transform
kiya hai)

non rel. odd (expansion
of bayes rule)

Day / Date: _____

→ Deriving PRP ranking principle.

11.3.1 → Deriving a ranking function for query terms

- Assumptions that we'll follow:

A1 - one random variable for each term (word)

; A2 - $d(w)$ are mutually independent given R .

A3 - $P(O|R=1) = P(O|R=0) = 0$

A4 - If word is not in query, it is equally likely to occur in $R \notin NR$ doc; only need to calculate prob of common words in query & doc.

A5 - On avg, a query word will occur in half of the rel doc; $P(w) \approx 1 - P(w)$ will cancel out.

A6 - NR doc approximated by collection as a whole (most docs are NR).

- Derivation of PRP

→ Assume, $D = \{d_1, d_2, \dots, d_n\}$ & q is a fixed query. ; $|D|=n$

→ As per PRP, for a given q , we can rank doc in decreasing order of probability of relevance.

rank discrepancy $P[R=1/d_i, q]$ — formula of ranking (compare query doc with others to see if rel. or not)

representing in vector \vec{x} $\rightarrow P[R=1/\vec{x}, \vec{q}] = P[R=1/\langle t_1, \dots, t_n \rangle, \vec{q}]$

→ This ranking of doc can be represented in prob. of odd:

$O[R/\vec{x}, \vec{q}] = P[R=1/\vec{x}, \vec{q}] \rightarrow$ true event

\downarrow odd prob. $\overline{P[R=0/\vec{x}, \vec{q}]} \rightarrow$ reciprocal event

Now expanding using baye's rule:

relevance odd $(P[R=1/\vec{q}]) \cdot P[\vec{x}/R=1, q] / P[\vec{x}/\vec{q}]$

→ Here, $P(\vec{x}/\vec{q})$ will cancel out whereas,

NR odd $(P[R=0/\vec{q}]) \cdot P[\vec{x}/R=0, q] / P[\vec{x}/\vec{q}]$

$P[R=1/\vec{q}] \approx P[R=0/\vec{q}]$ are constant as per assumption (A4)

$$\therefore O[R/\vec{x}, \vec{q}] = \frac{P[\vec{x}/R=1, \vec{q}]}{P[\vec{x}/R=0, \vec{q}]}$$

writing in vector form:

total no. of terms $\sum_{t=1}^n$

$$P[\vec{x}_t/R=1, \vec{q}] \cdot O(R/\vec{q}) \xrightarrow{\text{fixed (can be removed)}}$$

$\xrightarrow{\text{for each term that occurs in rel doc for query}}$
 $\xrightarrow{\text{in NR doc}}$

Day / Date: _____

Let, $P_t = P(x_t=1 / R=1, \vec{q})$; $P_t \rightarrow$ (count of) kitney rel doc mein term has appeared / total relevant

$u_t = P(x_t=1 / R=0, \vec{q})$; $u_t \rightarrow$ kitney non-rel doc mein term has appeared.

We'll get the ~~confusion~~ matrix as : $R=1$ $R=0$

term appeared in doc $\leftarrow x_t=1$

(distrn't) $\leftarrow x_t=0$ P_t u_t

1 - P_t 1 - u_t

Now, we have :

$$O[R | \vec{\pi}, \vec{q}] = \pi_{t=x_t=1} P_t + \pi_{t=x_t=0} u_t$$

$\pi_{t=x_t=1}$ term appeared in doc $\pi_{t=x_t=0}$ term not in doc but in query
↳ query both

Now multiply above eq. by :

$$\pi_{t=x_t=0} u_t \quad 1 - P_t \text{ is a constant}$$

$$O[R | \vec{\pi}, \vec{q}] = \pi_{t=x_t=1} \frac{P_t(1-u_t)}{u_t(1-P_t)} \cdot \pi_{t=x_t=0} u_t \quad \frac{1-P_t}{1-u_t} \text{ cancels out b/c of AS}$$

So, This odd $[R=1 / \vec{\pi}, \vec{q}]$ can be used as a Retrieval status value.

$$RSV = \pi_{t=x_t=1} P_t(1-u_t)$$

$\pi_{t=x_t=1}$ is ratio ko (* multiply kardlein
↳ terms common main unke)

∴ finally,

$\pi_{t=x_t=1}$	$P_t / (1-P_t)$
$\pi_{t=x_t=0}$	$u_t / (1-u_t)$

this formula will be used for calculation & implementation

→ we can use logarithm to transfer this into summation :

(see video of April '02 → 30 mins onward for application (numerical qn))

Day / Date: _____

→ Estimating $P_t \hat{=} U_t$

→ let N be no. of docs in corpus ; → R be # rel doc in sample.

→ n_t → no. of docs that contain term t ; → x_t → # rel doc that contains t .

$$P_t = \frac{x_t}{\text{total no. of rel doc}} \quad U_t = \frac{n_t - x_t}{N - R} \rightarrow \text{NR doc but have } t$$

for NR doc \downarrow $N - R \rightarrow$ total no. of NR doc

Lindstone Smoothing → $(\lambda = 0.5)$ → to avoid val = 0.

$$P_t = \frac{x_t + 0.5}{R + 1} \quad U_t = \frac{n_t - x_t + 0.5}{N - R + 1} \rightarrow \text{opar} + 0.5 \rightarrow \text{smooth} + 1$$

→ (introduced in general term)

BINARY INDEPENDENCE MODEL

• we get the below eqn after doing basic steps of PRP:

$$O(R | \vec{q}, \vec{u}) = O(R | \vec{q}) \cdot \prod_{t=1}^m \frac{P(x_t=1 | R=1, \vec{q})}{P(n_t | R=0, \vec{q})}$$

divide the terms to give:

$$O(R | \vec{q}, \vec{u}) = O(R | \vec{q}) \cdot \overbrace{\prod_{t: x_t=1}^m \frac{P(x_t=1 | R=1, \vec{q})}{P(n_t=1 | R=0, \vec{q})}}^{(U_t) \text{ term in rel doc}} \cdot \overbrace{\prod_{t: x_t=0}^m \frac{P(x_t=0 | R=1, \vec{q})}{P(n_t=0 | R=0, \vec{q})}}^{(\text{term appearing in doc})} \cdot \overbrace{\prod_{t: x_t=0}^m \frac{P(x_t=0 | R=0, \vec{q})}{P(n_t=0 | R=0, \vec{q})}}^{(\text{term not appearing in doc})}$$

We get contingency table:

	doc	relevant ($R=1$)	non-rel ($R=0$)
Term present	$x_t=1$	P_t	U_t
Term absent	$x_t=0$	$1 - P_t$	$1 - U_t$

constant for the particular query

$$\therefore O(R | \vec{q}, \vec{u}) = \underbrace{O(R | \vec{q})}_{\text{query term found in doc}} \cdot \prod_{t: x_t=1, q_t=1} \frac{P_t / U_t}{\prod_{t: x_t=0, q_t=1} \frac{1 - P_t}{1 - U_t}}$$

Dazzle

Day / Date:

$$h(i,j) = \text{term appearing in doc } i \rightarrow \begin{cases} 1 & \text{if } j \in \text{rel.} \\ 0 & \text{if } j \in \text{non-rel.} \end{cases}$$

$$\text{RSV}_d = \log \left[\frac{\pi}{t: x_t = q} \frac{P_t(1-U_t)}{U_t(1-P_t)} \right] \Rightarrow c_t = \log \frac{P_t}{1-P_t} + \log \frac{1-U_t}{U_t}$$

where, $c_t \rightarrow \log \text{odd ratios for term in } q$.

$\rightarrow c_t = 0$ if t has equal odds of appearing in rel & non-rel docs

$\rightarrow c_t = +\infty$ if t has more odds of appearing in a rel doc.

Probability estimates in theory

	doc	rel	non-rel	total
Term present	$x_t=1$	s	$df_t - s$	df_t
Term absent	$x_t=0$	$s-s$	$(N-df_t) - (s-s)$	$N-df_t$
Total	S	$N-S$		N

$$f/(N) = g(u_1) + R/n_1$$

$df_t \rightarrow$ no. of doc that contains t

$s \rightarrow$ relevant

$df_t - s \rightarrow$ non-relevant

Using this,

$$P_t = s/S \quad \text{if } U_t = (df_t - s)/(N-S)$$

$N \rightarrow$ total no. of doc

$$\text{RSV}_d \rightarrow c_t = K(N, df_t, S, s) = \log \left| \frac{s/S - s}{(df_t - s)/((N-df_t) - (s-s))} \right|$$

$N - df_t \rightarrow$ doc in which t does not appear

$s-s \rightarrow$ rel. doc but t does not appear

→ RELEVANCE FEEDBACK IN PROB IR

① Guess initial estimates of P_t & U_t . e.g. assume P_t is constant across x_t in q so $P_t = 1/2$

② Use current estimate to guess best set of rel doc $\rightarrow R = \{d : R_{d,q} = 1\}$; present this rel doc to user.

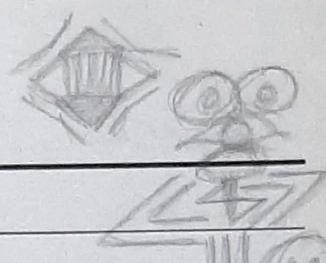
③ Interact with user to refine model of R \rightarrow we do this by learning from user's rel. judgement for subset of doc.

V is partitioned in 2 subsets $\rightarrow VR = \{d \in V, R_{d,q} = 1\} \subset R \neq VNR = \{d \in V, R_{d,q} = 0\} \rightarrow$ disjoint from R .

④ Re-estimate P_t & U_t on basis of known rel & non-rel doc. If sets VR & VNR are large enough, we may be able to estimate quantity directly from those docs $\rightarrow P_t = |VR_t|/|VR|$; where $VR_t \rightarrow$ set of docs in VR containing t .

$$P_t = \frac{|VR_t| + 1/2}{|VR| + 1} \rightarrow$$
 general practice of doing smoothing.

Day / Date: _____



→ However, set of docs judged by user (V) is very small so resulting stats is unreliable. ∴ we combine new

info with original guess in a process of Bayesian Update: $P_t^{(k+1)} = \frac{1|VR_t| + kP_t^{(k)}}{|VR_t| + k}$

Hence $P_t^{(k)}$ → software interface

$P_t^{(k)}$ is k^{th} estimate for P_t in an iterative updating process & is used as Bayesian prior in next iteration with the weighting of k .

→ Prior is strongly weighted so that the estimate doesn't change too much from the evidence provided by a very small no. of docs.

⑤ Repeat above process from step 2, generating a succession of approximations to R & hence P_t , until user is satisfied.

RELEVANCE

→ PSEUDO FEEDBACK IN PROB IR

① Assume initial estimate for P_t & $|U_t|$ as above.

② Determine a guess for size of rel doc set. If unsure, a conservative (small size) guess is likely to be best. This motivates use of a fixed size set of N of highest ranked doc.

③ Improve guess for P_t & $|U_t|$. Re-estimate P_t based on set V instead of VR . Let V_t be subset of doc in V .

$$P_t = |V_t| + \frac{1}{2} \quad \text{and if we assume that docs that are Non-relevant are non-rel.}$$

then we can update our $|U_t|$ estimate as: $|U_t| = \frac{df_t - |V_t| + \frac{1}{2}}{N - |V_t| + 1}$

④ Go to step 2 until ranking of returned results converges.

→ Once we have real estimated val. for P_t , then c_t used in RSV val. looks almost like tf-idf value.

ISSUES WITH PROBABILISTIC IR

→ Term Independence (assuming all terms are independent which is not true. Some terms might appear (joint prob not calculated) in doc but not others)

→ Terms not in query don't affect the outcome (not true; synonyms or some concept may affect)

→ Boolean representation of docs / queries / relevance (Binary representation only; no partial)

→ Doc relevance values are independent (not true, because docs might be inter-related & relevancy might change in relation to other)

Day / Date:

(c) crossover - main point
more diversity in one iteration

④ (AI)

GENETIC ALGO EXAMPLE

Example of PP

sum
(1) of w < 30

TEXT CLASSIFICATION & NAIVE BAYES

(Chp 13)

② max survival chance

D = {d₁, d₂, ..., d_n} m

C = {C₁, C₂} → no. of classes are fixed ; given

+ $\rightarrow c_j \left\{ \begin{array}{l} \text{for all doc classify to which class do they belong} \\ \text{(calculation survival pt)} \end{array} \right. \quad \begin{array}{l} \text{(d)} \\ \text{choose 3 max} \end{array}$

e.g. classify mails if spam/harm

① 0 1 1 1 0 1 0 + 7 + 10 + 5 + 0 + 17

Ad Hoc Query - Free text query - search engine supports it (any query need of user)

Standing Query - Executed on a collection of changing doc e.g. new cases of corona virus.
(can multiple times)

Text / Doc Classification

(general topics)

Text → can be anything ; news, email, book, researches (no boundary).

Doc → has a particular agenda.

(unit text ka aik doc, doc → lone tweet or a research paper)

Text/doc classification is process of assigning over predefined set of classes.

* we refer to this classification using standing query is also called routing / filtering.

50 items in 100 generation

TYPES OF CLASSIFICATION

→ ① Manual

Manually mark (now you block a specific no. or email id)

It is challenging b/c we manually cater.

Manual Hand Code Rules → Rules based on combination of words or other features.

→ ② Supervised

A very large dataset available. We know every doc class & all is easily available.

$$d_i = \vec{x}_i \quad ; \quad \vec{x}_i = \langle t_1, \dots, t_m \rangle \quad ; \quad \vec{x} \rightarrow \hat{y}$$

↓ doc ↓ class

Dazzle

Day / Date:

SE (Clip) (not in book)

① USER STORY

① Training req → ② Testing → test history train is per how efficient our result is

$f(d) \rightarrow c$: who + what + why → imp

* labelled data is available (rich) & we build classifier by making our program learn through this data. mode owner or customer start writing cards which is a high

→ A training set of m hand labelled dataset
($\in \mathbb{N}$ gives ans)

→ ③ Semi-supervised

→ Start as supervised. us analyst & developer together ask the customer

→ Tune in such a way that it will continue to learn from functionality of its working
< classification > on

* if labelled dataset has us see new doc to classify bring & check that data & train further
→ Format of func. in card

→ ④ Unsupervised (type I want to) I want to perform an action so that

Training dataset has labels but we give only doc to classifier & check if it has accurately classified.

→ Learn on its own how to classify text/doc by implicitly learning the features

13.2 - NAIVE BAYES TEXT CLASSIFICATION

(N.B.)
→ first & basic supervised learning method - Multinomial Naive Bayes - Probabilistic learning method

→ The probability of a doc d being in class c is:

$$P(c|d) \propto P(c) \prod_{k=1}^{n_d} P(t_k|c)$$

where,

$$P(t_k|c) \rightarrow \text{prob of term } t_k \text{ occurring in doc of class } c. \rightarrow \text{measure}$$

$P(c)$ → prior probability of doc occurring in class c

n_d → no. of tokens in d

→ simple classifier based / works on prior prob. of classes' relationship.

Day / Date: _____

Similar

In text classification, our goal is to find best class for doc. Best class in NB is mostly likely or maximum posterior (MAP) class c_{map} :

$$c_{\text{map}} = \arg \max_{c \in C} \hat{P}(c|d) = \arg \max_{c \in C} \hat{P}(c) \propto \prod_{k=1}^n \hat{P}(t_k|c) \quad \text{feb '20}$$

for the predictions.

\hat{P} → represent $P(b/c)$ we don't know true values of parameters $P(c) \& P(t_k|c)$, but estimate them from training set.

for argmax, we'll calculate for everyone, all c belongs to class (c) . So max \log that will be the class most likely!

MAP = $\arg \max_{c \in C} P(c|d)$ a chance to everyone to avoid this opportunity

to showcase their talents in the auditions. The auditions will be held in

audition → $= \arg \max_{c \in C} P(d|c) / P(c) \rightarrow$ applying Baye's Rule on denominator.

$P(d)$ → constant → doesn't affect ranking

Audition = $\arg \max_{c \in C} P(d|c) / P(c) \leftarrow$ to well drop the denominators

→ How To CALCULATE THESE PROBABILITIES?

To first attempt → max likelihood estimates → simply use freq in data (counting)

P for a given class $j \leftarrow \hat{P}(c_j) = \text{doc_count}(c=c_j) \rightarrow$ for given class in training dataset

(so we'll calculate prob. of every class) $N_{\text{doc}} \rightarrow$ total no. of doc in training dataset.

prob. of a word w_i belonging to class $j \leftarrow \hat{P}(w_i|c_j) = \frac{\text{count}(w_i, c_j)}{\sum_{w \in V} \text{count}(w, c_j)} \rightarrow$ count of word i in class j

$\sum_{w \in V} \text{count}(w, c_j) \rightarrow$ total count of word i in all classes

set of words = set of vocab

① → formula for classes

after smoothing → $\hat{P}(c) = \frac{\text{count}(c) + 1}{\sum_{c \in C} (\text{count}(c)) + 1}$

② → formula for words

$= \frac{\text{count}(w, c) + 1}{\sum_{w \in V} [\text{count}(w, c)] + |V|} \rightarrow$ count of vocab. Dazzle

min circuit download time
 \downarrow
 time = F/d_{min}

Time to distribute F to N clients using c-serv || D approach $\rightarrow \max(F/U_s, F/d_{min})$
 file dist. time = P2P everyone is downloading & uploading their own upload link or relay
 server trans only 1 copy $\rightarrow F/U_s$; each client must download $\rightarrow F/d_{min}$
 client as Aggregate max download F/N bit = max upload rate $\rightarrow U_s + \sum U_s$ (assuming max download rate
 $(\text{total download data})$) \rightarrow Time to dis. F to N clients using P2P
 $\max(F/U_s + F/d_{min}, FN/U_s + \sum U_s)$

→ DNS (Domain Name Sys)

Host name \rightarrow converter or ke IP address ke data hai

host time \uparrow

- ① distributed database \rightarrow implemented in hierarchy of many name servers.
- ② App-layer protocol \rightarrow host name service communicate to resolve name (address & name data)

• DNS services + structure

① \rightarrow hostname to IP address translation karta hai.

② host aliasing : canonical name

\downarrow host with a complicated host name can have one or more alias name.

③ Mail server aliasing

(many IP add. corresponds to one name)

canonical hostname.

④ Load distribution - for replicated web servers, a set of IP addresses is thus associated with one \rightarrow DNS rotation distributes the traffic among the replicated servers.
 \rightarrow also used for email servers.

Why not centralize DNS?

\rightarrow single point of failure \rightarrow traffic volume \rightarrow maintenance

DNS : root name service

\rightarrow contacted by local name servers that cannot resolve name.

\rightarrow root name service \rightarrow contact authoritative servers. \rightarrow top level domain (TLD) servers (forward TLD to query)

authoritative DNS server.

\downarrow auth. to IP address

will search authoritative
K to represent kr
hierarchy.

*local DNS service

\rightarrow has local cache of recent name to add. translation; acts as a proxy \rightarrow forward query to

hierarchy.

*DNS name resol.

b) host at loc. pol. \rightarrow DNS caching, updating record

c) iterative query \rightarrow once (any) name server learns mapping, it caches mapping.

once (any) name server learns mapping, it caches mapping.

\rightarrow cache entries timeout (disappear) after some time (TTL)

c) recursive query \rightarrow cache updating

• DNS records \rightarrow distributed db storing resource records (RR) \rightarrow RR format: (name, val, type, ttl)

• DNS protocol \rightarrow query & reply. both have same msg format. \rightarrow 16-bit ID & flags

• Inserting record in DNS \rightarrow reg name at DNS registrar; provide name, IP add. of author. name server
 \rightarrow create authoritative server type A for records for networkutopia.com.

pure

P2P operates there & no authority on server; authority and sys directly connect. peers are immediately.

File distribution : client server vs. P2P \rightarrow ask sys pc dependency ne hogi — connected peer 2 peer n server hain to n copies.

Time to send 1 file \rightarrow $\frac{F}{U_s} \rightarrow$ file size \rightarrow speed.

$\left\{ \begin{array}{l} F/U_s \rightarrow \text{for n files} \\ \dots \end{array} \right.$

Dazzle

11/11 Ayesha Singular => simple => similar.
 Day / Date: Ayesha => similar.

" (similar to - simple) website use cookie 4 components.

User - Server state : Cookie (saves user state)

~ Server app header file mein client ko id dega

→ server se sari conn. us id ke against hogi - history maintain.

(database browser aur server dono sides pc ho raha hai)

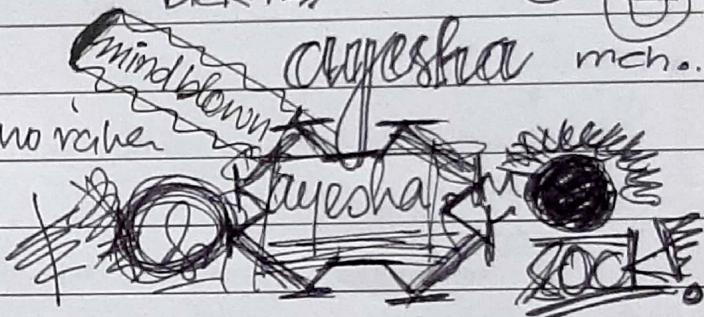
for each web a new id is given to user.)



http transaction ki state maintain ho rahi

HTTP → stateless

↓ http msg carry state.



Web Cache (proxy server)

get method se web pg search kring
 → if first time, no copy on proxy → it will send pg
 → if copy available for cache, cache send obj
 get → if normal web page on browser
 post → if send a form to server.
 * as a client when copy not available
 * cache act both as server & client

(web cache reduces response time)

→ cache example read from book

→ RTT → honot zyada hoga hai

↓
web cache → speed high & cheaper.

(2 types of delay)

sign

signature

size

signature

size

signature

size

signature

size

signature

size

signature

size

CONDITIONAL GET → Proxy server ko itna rakhay agr page ki copy nai
to user req. age send na kren server ko page.

Electronic Mail

→ user agent (email generate)

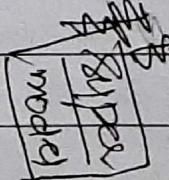
→ mail server (saving email)

→ simple mail (will communicate with main server)

transfer protocol:

SMTP

→ (uses TCP to reliably send)



3 main component

mail access protocol

POP3 + IMAP

Day / Date:

$$P(c_j | d) = P(d | c_j) \cdot P(c_j)$$

MULTINOMIAL NB (Algorithm)

→ Training → $D = \{d_1, d_2, \dots, d_n\}$; $C = \{c_1, c_2, \dots, c_n\}$ → given

① Extract vocab (clean data) from D as V .

② Count no. of docs as N . $\rightarrow n = 101$

③ for each class, calculate posterior prob from prior prob.

④ Take out ~~the~~ count of doc belonging to that class as N_c .

⑤ prior prob $\rightarrow N_c / N$. $\hookrightarrow P(c_j)$

⑥ text_c mean combine all the text of those docs that belong to c .

⑦ then for each term that belongs to Vocab.

⑧ calculate count of term that belongs to text_c as T_{ct}

⑨ then for each term that belongs to Vocab.

⑩ calculate conditional prob $[t][c]$ i.e. $\frac{T_{ct} + 1}{\sum_t (T_{ct} + 1)}$. $\hookrightarrow P(d | c_j)$

⑪ Returns V (vocab), prior prob & conditional prob for all classes - return V , prior, condpr

→ Testing → $C, V, \text{prior}, \text{condprob}$, & new doc d given

① Extract words from the doc using vocab as V .

② for each class i that belongs to C .

③ calculate prior prob of every class as score

④ for each term t that belongs to word (V)

⑤ calculate condprob in score & save it in score

⑥ return the max score for showing to which class d belongs to.

$\text{Train}(C, D)$

$V \leftarrow \text{ExtractVocab}(D)$

$N \leftarrow \text{CountDocs}(D)$

for each $c \in C$

do $N_c \leftarrow \text{CountDocInClass}(D, c)$
 $\{\text{prior}[c] \leftarrow N_c / N\}$

$\text{text}_c \leftarrow \text{Concatenate}(D, c)$

for each $t \in V$

do $T_{ct} \leftarrow \text{CountTerm}(\text{text}_c, t)$

for each $t \in V$

do $\text{condprob}[t][c] \leftarrow$

$\leftarrow \frac{T_{ct} + 1}{\sum_t (T_{ct} + 1)}$

$\text{ApplyNB}(C, V, \text{prior}, \text{condprob}, d)$

$W \leftarrow \text{ExtractTokensFromDoc}(d)$

for each $c \in C$

do $\text{score}[c] \leftarrow \log \text{prior}[c]$

for each $t \in W$

do $\text{score}[c] + \log \text{condprob}[t][c]$

return $\underset{c \in C}{\operatorname{argmax}} \text{score}[c]$

for example ques (numerical → see April 9 video .35 min onwards)

→ agar testing mein koi aisi term a jaye which doesn't belong to training data, we calculate its prob. by $P\left(\frac{1}{V+1}\right)$ → add 1 in vocab. e.g. unha. $|V| \text{vocab} = 6$ so $P(\text{unha}) = \frac{1}{6+1} = \frac{1}{7}$

Day / Date: _____

VARIATIONS OF NAIVE BAYESIAN

(computationally lengthy)

doc ke token se class ka andaza

→ Two ways in which we can setup NB classifier : ① Multinomial NB ② Multivariate Bernoulli model

→ MULTIVARIATE BERNOULLI MODEL

→ Works on doc count ghol on features.

Training : TrainBernoulli(C, D)

1. $\{ V \leftarrow \text{ExtractVocab}(D)$
2. $N \leftarrow \text{count_docs}(D)$
3. for each $c \in C$
4. do $N_c \leftarrow \text{count_docs_in_class}(D, c)$
5. $\text{prior}[c] \leftarrow N_c/N$
6. for each $t \in V$
7. do $N_{ct} \leftarrow \text{count_docs_in_class_containing_term}(D, c, t)$
8. $\text{condprob}[t][c] \leftarrow \frac{N_{ct} + 1}{N_c + 2}$ *(smoothing)*
9. return $V, \text{prior}, \text{condprob}$

Testing : ApplyBernoulli($C, V, \text{prior}, \text{condprob}, d$) {

1. $V_d \leftarrow \text{ExtractTermfromDoc}(V, d)$
2. for each $c \in C$
3. do $\text{score}[c] \leftarrow \log \text{prior}[c]$
4. for each $t \in V_d$
5. do if $t \in V_d$
then $\text{score}[c] += \log \text{condprob}[t][c]$
6. else $\text{score}[c] += \log(1 - \text{condprob}[t][c])$
7. return $\arg\max_{c \in C} \text{score}[c]$

→ For example numerical see WII slides no. 11, and for differences b/w both types too.

Day / Date: _____

(* accuracy increases w/c of side usage allowed)

②

EVALUATION OF CLASSIFICATION TASK

	spam	Ham	Accuracy = $\frac{35+54}{120}$
Predicted spam	35	16	Error = $\frac{15+16}{120}$
Predicted Ham	15	54	

Recall — fraction of docs in class i classified correctly (kita accurately class to identify tr railway)
 $\rightarrow \text{rel retrieved} / \text{total relevant} ; \frac{c_{ij}}{\sum c_{ij}}$

Precision — fraction of docs assigned class i that are actually about class i. $\rightarrow \frac{\text{rel}}{\text{total retrieved}}$.
 $= \frac{c_{ii}}{\sum c_{ij}} \rightarrow \frac{0+0+0+0}{\cancel{1}+\cancel{1}+\cancel{1}+\cancel{1}} \rightarrow \frac{\text{sum}}{\text{sum} \cdot (\text{diagonals ka sum wd. } k_i)}$

Accuracy — Fraction of docs classified correctly ie ($1 - \text{error rate}$)

\rightarrow For calculating combined multiple performance measure, we use micro / macro avg

- Macroaveraging — compute performance for each class, then avg.

- Microaveraging — collect decisions for all classes, compute contingency table, evaluate.

e.g.

c_1		c_2		microavgtable		macroavg
yes	no	yes	no	yes	No	
yes	10	10	yes	90	10	Microavg $P = \frac{(0.5 + 0.9)}{2} = 0.7$
no	10	910	no	10	890	\downarrow (calculate separately)

$$\rightarrow \text{Microavg } P = \frac{100}{20} = 0.83 \quad ; \text{ microavg score is dominated by score on common classes}$$

SUPERVISED LEARNING — DATASETS

e.g. $X \rightarrow Y$

$e_i = f_1, \dots, f_m \rightarrow Y$
 \downarrow
 features

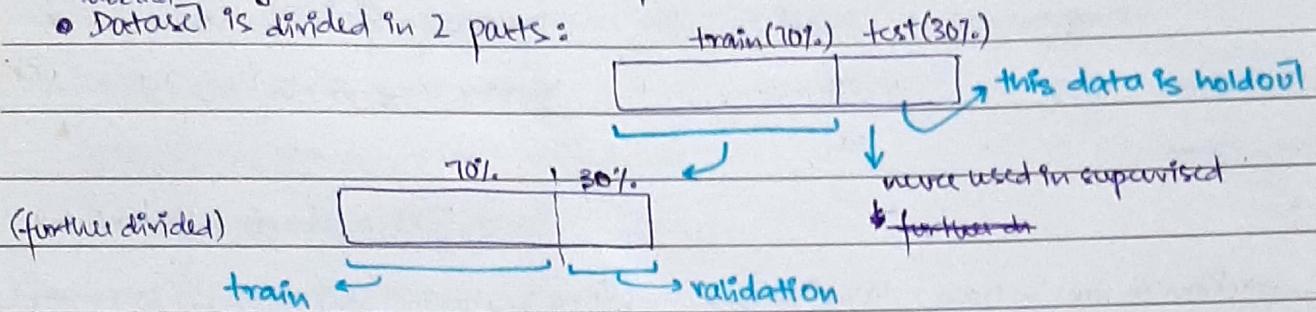
; for binary classification e.g. Spam / Ham

Dazzle

Day / Date: _____

→ Model fitting → feature selection / parameter estimates.

- Dataset is divided in 2 parts:



→ this approach is called Holdout aka split approach.

→ 30% validation is used for testing whereas rest is used for training.

→ for cross validation & better model performance, we use k-fold technique.

$k \rightarrow$ no. of parts; use $k-1$ parts to train model on k different parts of data & use 1 part for testing & avg out the performance of all models.

Variance vs. Bias → two types of errors

Error due to:

Biasness → Difference b/w expected & actual

Variance → variability of model.

e.g. bulls eye. → If low var. & low biasness so near target performance.

CLASSIFICATION MODEL & ACCURACY

→ Classification model learning is a challenging task.

• Textual docs are rich in feature → take word to its relevant space. cluster near words. (make a group)

• Feature Selection is helpful in simplifying the model & improving accuracy of model.

→ Textual features: Words/Tokens, Phrases/Bigram, Graph, NLP based features.

FEATURES

An individual measurable property of a phenomenon being observed.

→ Choosing informative, discriminating & independent features is a crucial step for effective algo in pattern recognition, classification & regression.

Dazzle

Day / Date: _____

→ Features are usually numeric, but structural features such as strings & graphs are used in syntactic pattern recognition.

FEATURE SELECTION

→ Purpose of selection of subset of features for training set & using only this subset as features in text classification.

→ Feature selection serves two main purposes:

- ① It makes training & applying a classifier more efficient by decreasing the size of effective vocab.
- ② Feature selection often increases classification accuracy as noise features are eliminated through it.

→ Types of Selection Feature Method:

① WRAPPERS

→ Jo features nikalay hain us se dekhin ge ke wo target ko kitna bcharta predict kar rheay.
We make a classifier using these features & then measure the accuracy of the feature, & select the features that are highly accurate.

→ Wrappers use classification accuracy of some learning algo as their evaluation function.

→ When no. of features is high, they are time consuming b/c wrappers train a classifier for evaluating each subset feature.

→ Wrappers not suitable for text classification.

② FILTERS

→ Perform feature selection independently of learning algo that will use selected features.
(first randomly filter features & then check accuracy)

→ For evaluation of a feature, filter use an evaluation metric that measures the ability of feature to differentiate each class.

→ Much less time consuming than wrappers & widely used in text classifi. Dazzle

Day / Date: _____

(3) EMBEDDING (new technique.)

- Use a group of features. We take projection of a particular word from relevant space & place the group of features at that place.
(Aik word ko uski relevant space se utha ke use aapne words ka group hona ke us word ki jagah place krdetay hain) • Different target space mein different features.

• Forward Feature Selection (FFS) — Aik feature select kerte uski accuracy, phr 2, then 3 and so on. Make numerous combinations.

• Backward Feature Selection — n no. of features ki accuracy, then n-1, then n-2 & so on.

HOW FEATURE SELECTION WORKS?

- FS replaces a complex classifier (using all features) with a simple one (using subset of features).
- Purpose of feature selection algo is to select only those features that to compute utility measure $A(t, c)$ for each term of vocab & select the k terms that have highest val. of $A(t, c)$.
- * $A(t, c)$ → function A to calculate val. of feature t & class c to see if t should be added in classifier or not.
- Algo of feature selection in (fig-13.6)

Underfitting → classifier has less knowledge due to which it classifies wrong.

Overfitting → Model is trying to adopt all features ; makes function complex (many parameters)
(none of these are desirable) * we require best fitting. Selection of relevant features

→ size, colour, shape & material are imp features.

Feature Extraction ~ Similar to feature selection but we extract relevant feature from an image that is relevant to the domain & train model with that sufficient knowledge.

Day / Date: _____

THREE APPROACHES TO FEATURE SELECTION

- ① Mutual Info (MI), $A(t, c) = I(U_t; C_c)$
- ② Chi square, χ^2 , $A(t, c) = \chi^2(t, c)$; general lib \rightarrow chi2, χ^2 , xsq()
- ③ Freq based features, $A(t, c) = N(t, c)$ ie ($tf * idf$)

Noise removal / Outliers \rightarrow Values that are outside our decision boundary

↳ challenge in Text

• MUTUAL INFORMATION

(feature class se kitha related hai aur kitha help kr raha hai class identify kaise mein)

→ MI measures how much info the presence/absence of a term condition contributes to making the correct classification decision on c .

$I(U_t; C_c)$ — N_{11} : term present & belongs to class
random var. for \leftrightarrow term \leftrightarrow classes
 N_{10} : term " but doesn't belong to class
 N_{01} $\not\in N_{00}$.

→ Example 13.3 for numerical.

• CHI SQUARE METHOD

→ Important for text classification.

Null hypothesis: $H_0 \rightarrow (t, c) \rightarrow$ feature t & class c are ^{not} significantly related.

$H_1 \rightarrow (t, c) \rightarrow$ " " " " " are significantly related.

→ t : occurrence of term ; c → occurrence of class . Rank terms wrt :

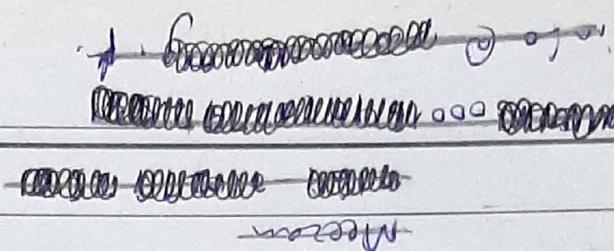
$$\chi^2(D, t, c) = \sum_{t \in \{0, 1\}} \sum_{c \in \{0, 1\}} \frac{(N_{etc} - E_{etc})^2}{E_{etc}}$$

→ χ^2 is a measure of how much expected counts E & observed count N deviate from each other. A high val. of χ^2 indicates that the hypothesis of independence, which implies $E \approx N$ are similar, is incorrect.

→ Example 13.4 for numerical.

Dazzle

Day / Date:



FS increases & at one point it reaches max. Then it decreases

redundant — less info features // irrelevant to our prob

↑
effect of feature set size on accuracy

Feature Generation

The idea is to generate high level features (more abstract) from low level features.

Example: sentence dependency graph from sentence.

↓ i.e. If height & mass given so we can add them for calculating BMI.

Add.
~~ANIMAL~~
~~ANIMAL~~ ...
(airbit) +
(anybody)

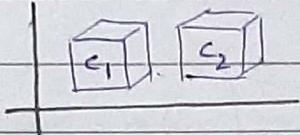
INFORMATION RETRIEVAL (pt. 2)

Day / Date: _____

VECTOR SPACE MODEL

- A doc is represented as a feature vector in feature-dimensional space.
- Distance is generally the notion of similarity for this model.
- General assumptions for classification:

Contiguity Hypothesis — Docs in same class form a contiguous region & regions of different classes don't overlap.



CLASSIFICATION — VSM

$$D = \{D_1, D_2, \dots, D_n\}$$

$$C = \{C_1, C_2\}$$

- Corpus — Classification dataset

→ Preprocessing

- All docs of C_1 : $\square \square \square$

$f_{1c_1}, f_{2c_1} \rightarrow$ features of these docs

$$S_1 = \square \square$$

$$\downarrow f_{1c_2}, f_{2c_2}$$

create a vector space model of combined features

$$[f_1 | f_2 | \dots |]$$

Two algos discussed:

① **ROCCHIO'S ALGORITHM** → can be used for text classification as well.

→ divides the vector space into regions centered on centroids or prototypes, one for each class computed as the centre of mass of all docs in the class.

→ simple & efficient classification, but inaccurate if classes are not approximately spheres with similar radii.

How algo works

→ Has class be doc to separate kernels & merge to make one doc.

→ Make a high dimension space for each doc & calculate unit/centroid vector of all docs.

i.e. $\vec{U}_j = \frac{\vec{d}_1 + \vec{d}_2 + \vec{d}_3}{3}$ → for each class calculate mean vector.

→ training func returns centroid of each class.

Dazzle

Day / Date: _____

Intesting → calculate distance of the test doc with each mean (centroid) of training doc & if it is the least (minimum) distance then, that would be returned & that class would be assigned to the doc.

→ centroid vector (\vec{u}_i) is also called prototype.

② K-NEAREST NEIGHBOUR LEARNING

→ K-NN learning is an example based learning / instance based learning. Memory based method in which learning is just storing the representations of training (e.g. in D). (example to use neighbouring example to compare training in the class to see how similar they are. Jisse zyada similar wohi assign kardein ga).

→ Testing instance x : Compute similarity b/w x & all examples in D. Assign x the category of the most similar example in D.

→ It does not explicitly calculate a class/category prototype descriptor.

* Algo for train & test in slides.

* Example of numerical from book / slide. * k should be odd (preferable)

→ Advantages

- Simple, intuitive, easy to implement

- Only one hyper parameter i.e. k

→ Disadvantages

- No training time, but large computation time for large datasets (more examples & more dimensions)

- Lazy learner → no precompute model

distance noise

- Sensitive to value of k , distance func & noisy data (example is placed in wrong class, labelled wrongly or $\frac{\text{value of feature of}}{\text{example has changed}}$)
attribute noise

Day / Date: _____

(Week 12)

TEXT CLUSTERING

(Chp 16 + 17)

CLUSTERING — Unsupervised machine learning technique. It separates out similar obj in a heterogeneous collection automatically.

→ The features of the obj implicitly identify for clustering.

→ An objective func is used to optimize similarity b/w obj in a cluster & reduce the similarity b/w different clusters.

↳

* aim group to be each mean (apar mein) similarity all doc group to be each major differences have
↓ challenging prob → dual optimization → high intra-cluster similar & low inter-cluster similarity.

→ Doc clustering is a special prob → obj to cluster are docs.

- Natural Lang Text
- Knowing exact no. of distinct groups (clusters)
- Producing an optimal clustering arrangement
- Label the groups
- Evaluation of clustering results.

CLASSIFICATION Vs. CLUSTERING

Classification → supervised learning → classes are human-defined & part of input to the learning algorithm.

Clustering → unsupervised → clusters are inferred from the data without human input

CLUSTER HYPOTHESIS

→ Docs in same cluster behave similarly wrt relevance to info need.

→ All applications of clustering in IR are based (directly/indirectly) on cluster hypothesis.

→ "closely associated docs tend to be relevant to same requests".

Day / Date: _____

CLUSTERING APPLICATIONS (wide application in IR)

- ① Search result clustering
 - ② Scatter-Gather
 - ③ Collection clustering
 - ④ Cluster based retrieval.

Scatter-Gather → used for news refinement.

→ first scatter cluster (gather), then disperse (scatter) → into topics related to it, The cluster & scatter again for more refined subtopics of user's interest.

FLAT VS. HIERARCHICAL CLUSTERING

→ Flat Algorithms

→ clustering space is partitioned, partitions data (no relation) \nmid treat all doc at same level.

→ usually start with a random (partial) partitioning of docs into groups.

↓ refine relatively into distinct / non-common groups (naturally)

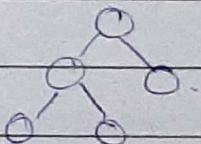
Main algo \rightarrow K-means

→ Hierarchical Algorithms

→ creates a hierarchy of all possible cluster levels.

→ Can be bottom up → agglomerative ; → can be top down → divisive

Main algo : HAC



HARD VS. SOFT CLUSTERING

Hard \rightarrow Each doc belongs to exactly one cluster.

Soft → A doc can belong to more than one cluster.

K-MEAN CLUSTERING

→ Best algo. Simple, works well in many cases.

→ Used as a default / baseline for clustering doc.

→ Performance is determined by initialisation & appropriate distance measure. Dazzle

Day / Date: _____

3 main steps in doc clustering:

(1) → Representation of docs (2) → Similarity measure (3) → Clustering approach

Similarity func or measure

$$\rightarrow \text{cosine}(d_i, d_j) = \frac{\vec{d}_i \cdot \vec{d}_j}{\|d_i\| \cdot \|d_j\|}$$

measure → any calculation that is transformed in number, not necessary to have a unit.

metric → should have 3 properties

① Equal self similarity ② Minimality

③ Symmetry ④ Triangle inequality

Equal Self Similarity → $d(A, A) = d(B, B) \rightarrow s(A, A) = s(B, B)$

Minimality → $d(A, B) > d(A, A) \rightarrow s(A, B) < s(A, A)$

Symmetry → $d(A, B) = d(B, A) \rightarrow s(A, B) = s(B, A)$

Triangle Inequality → $d(A, B) + d(B, C) \geq d(A, C)$

where A, B, C are any doc assumed as obj.

K-mean: (Algo in book)

① Selects K docs

① Selects random seeds (s_k)

② Calculate centroid \bar{m}_k of each cluster

③ And until cluster does not become same or distance is not minimum, we

will compare every doc with cluster's mean (\bar{m}_k) & align min. with it (reassignment of vector)

& recalculate mean (\bar{m}_k) → recomputation of centroids.

→ 3 criterion to stop clustering in k-mean:

① → arrangement is consistent. ② Give specific no. of iteration e.g. i=5.

some specific

(fix the loop)

③ → can stop it on the value of centroid.

→ variants of K-means

• K-medoids : resistance to noise/outliers b/c works using median.

• K-modes : extension to categorical data clustering analysis

Dazzle

• CLARA : extension to deal with large datasets

Day / Date: _____

CONVERGENCE OF K-MEAN

→ first there are at most K^n ways to partition N data points into k

For each iteration of algo, we produce a new clustering based only on old clustering.

→ If old clustering is same as new, then next clustering will again be same.

→ If new clustering is different from old, then new one has a lower cost.

→ Since algo iterates a function whose domain is a finite set $\rightarrow \therefore$ The iteration must eventually enter a cycle. The cycle must have length exactly 1. Hence k-means converges in a finite no. of iterations.

* N is generally a very great no. as compared to k. ($N \ggg k$)

There are K^n clusters. & running time of K-mean $\rightarrow \underbrace{K+n}_\text{time is bounded on no. of doc} \propto INI$

(partitions)

DRAWBACKS OF K-MEAN

① Sensitive to initial seeds (agrs seeds durr hui ya outlier hui)

② Difficult to compare results

③ Fixed no. of clusters, difficult to predict the actual no. from datasets.

EVALUATING CLUSTERING TASK

→ The task is very challenging \therefore evaluation is also challenging.

→ Evaluation of 2 types:

Internal — Only data is used to evaluate the clustering result

(When data is not labelled, and we don't have any way to evaluate clustering alone, so we calculate distance of obj among a cluster & distance b/w clusters.)

If obj among clusters are tightly packed (less distance) & distance b/w clusters is greater. It means 'good clustering'

Dazzle

Day / Date: _____

(ground truth)
External - Evaluation with gold standard. (labelled data given) ~~both~~

→ we hide labels and cluster the docs. Then check if they are clustered acc. to their correct label. If 100% correct, then best clustering.

WHAT IS A GOOD CLUSTERING?

→ Internal Criterion → Intra-class ie. Intra-cluster similarity is high.

- Inter-class similarity is low.

→ The measured quality of a clustering depends on both the doc representation & similarity measure used.

→ External Criterion → clustering should reproduce the classes in gold standard

• First measure for how well we were able to reproduce the classes — Purity.

• Purity → tells how homogeneous clusters are.

$$\text{purity}(\mathcal{R}, \mathcal{C}) = \frac{1}{N} \sum_k \max_j |w_k \cap c_j| ; \begin{array}{l} \text{should be in range } [0, 1] \\ 0 \rightarrow \text{not homogeneous} \\ 1 \rightarrow \text{max. homogeneity} \end{array}$$

$\mathcal{R} = \{w_1, \dots, w_K\} \rightarrow$ set of clusters, $\mathcal{C} = \{c_1, \dots, c_n\} \rightarrow$ set of classes.

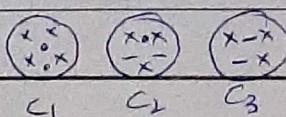
for each cluster $w_k \rightarrow$ find class c_j with most member n_{kj} in w_k .

$N =$ total no. of points.

• Random Index (How many attempts are made to do 100% correct clustering \rightarrow randomly assign 2 docs to same cluster if $\frac{2}{N}$ only if they are similar.)
if obj are picked 2 at a time

→ true + rc (TP) decision assigns 2 similar doc to same cluster. True - rc (TN) \rightarrow dissimilar docs to different cluster.

→ 2 types of error we can commit. (FP) → assigns 2 dissimilar doc to same cluster
(FN) → assigns 2 similar docs to different cluster.



→ 17 obj \rightarrow so we have $\binom{17}{2} = 136$ ways to make clusters
2 obj at a time

Dazzle

→ Numerical example in book.

Day / Date: _____

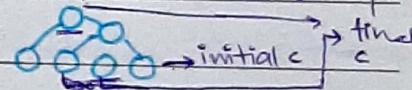
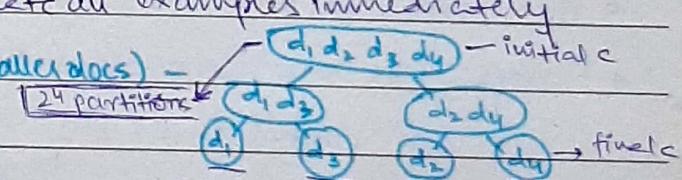
• Normalized Mutual Information (NMI)

- How much info does the clustering contain about the classification.
- Singleton clusters (no. of cluster = no. of docs) have max MI → means totally homogeneous
- one to one correspondence check ^{of} how much similarity with ground truth.
- relation classes with cluster.

• F Measure

Precision & recall are weighted.

AGGLOMERATIVE VS. DIVISIVE CLUSTERING

- Agglomerative (bottom up) method starts with each example in its own cluster & iteratively combine them to form larger & larger clusters. 
- Divisive (partitioning, top-down) separates all examples immediately into clusters (cluster to breaking in smaller docs) - 
- Simple HAC algo from book.
 - ① Take out similarity ~~from~~ ^{of} each doc with one another (similarity of 2 instances)
→ starts with instances in a separate cluster & repeatedly
 - ② Then merge krtai hai maximum similar doc ko. To shrink no. of clusters until there is only one cluster.
→ The history of merging forms a binary tree or hierarchy.

Variations of HAC

Dendrogram

- A binary tree that shows how clusters are merged/split hierarchically.
- Each node on tree is a cluster; each leaf node is a singleton cluster.
- computational complexity
- Running time of HAC if perfectly executed $\rightarrow O(n^2)$. If not perfectly, can be:
 $\rightarrow O(n^3)$
 $\rightarrow O(n^2 \log n)$
∴ datastructure used matters a lot & affect efficiency of algo.
- different variations of HAC depending on the max function we are using.
- In order to maintain an all over $O(n^2)$ performance, computing similarity to each cluster must be done in constant time.

Dazzle

Day / Date:

→ Variations of HAC

① Single Link Clustering (aka connectedness, the minimum or nearest neighbor method)

→ consider distance b/w 2 clusters to be equal to shortest distance from any member of one cluster to member of other cluster.

→ If data consist of similarities, the similarity b/w a pair of cluster is considered to be equal to the greatest similarity from any member of one cluster to member of other cluster.

→ The distance b/w two clusters is represented by - the dist. of 'closest pair of data obj' belonging to different clusters i.e. $d_{\min}(c_i, c_j) = \min_{p \in c_i, q \in c_j} d(p, q)$.

② Complete Link Clustering (aka diameter, the maximum method or farthest neighbour method)

→ consider the distance b/w two clusters to be equal to longest distance from any member of one cluster to member of other cluster.

→ farthest pair of data obj $\rightarrow d_{\max}(c_i, c_j) = \max_{p \in c_i, q \in c_j} d(p, q)$ → max distance & min. similarity b/w those points

③ Average Link Clustering (aka minimum variance method)

→ consider distance b/w two clusters to be equal to avg distance from any member of one cluster to member of other cluster.

→ avg dist. of all pairs of data obj $\rightarrow d_{\text{avg}}(c_i, c_j) = \text{avg}_{p \in c_i, q \in c_j} d(p, q)$

④ Centroid Distance

The distance b/w two clusters is represented by distance b/w the means of clusters.

$d_{\text{mean}}(c_i, c_j) = d(m_i, m_j)$ where (m_i & m_j are means of c_i & c_j)

→ There are many (20-25) variations. All variations will give different clusters.

Advantage of HAC → we don't have to predefine $k \rightarrow$ how many clusters we want.

• once clustering is done, we can then extract k ^{no. of} clusters from it.

Day / Date: _____

WEB SEARCH BASIC

(Chp 19)

Web Search - Client Server

Client — generally a browser — an app with GUI (environment)

Server — the server communication with a client via a protocol HTTP.

→ lightweight, simple & asynchronously carrying a variety of payloads (text, img) encoded in a simple markup lang. called HTML.

→ resources & webpages are available & we browse it by inputting URL.

- HTTP like as a responder → server ; • HTTP like as a requester → client.

HTTP

- An app protocol for distributed, collaborative & hypermedia info.

- HTTP header contains a lot of fields for effective transfer of info.

- Persistence connection — do not let connection die. If was client/server connection & things are transferred b/w it.

→ Major Advancement

- Session state can be maintained — session info (client-server) can be served through cookies.

- Support authentication mechanism.

→ HTTP Status Code

Client sends a GET url command then ~~<* blank lines *>~~ ^{((CRLF))} then Carriage Return (CR), line feed (LF)

↳ (protocol command completed)

In return, server returns a pg & forwards index.html pg to client.

• Informational 1XX

• Redirection 3XX

• Server Error 5XX

• Successful 2XX

• Client Error 4XX

Dazzle

Day / Date: _____

Server Side Scripting

Client pc continuously sends requests to server to fetch how many changes have taken place (A no. of server side scripting available).

Client Side Scripting

Generally UI & interaction with local machine, mostly Java Script.

Cascading Style Sheet (CSS)

A lang. that describes the style of an HTML doc. ; → depending on browser displays output differently on a single web pg.

* Server side pc data fixed hai. can make new pgs from each client's request & render it acc. to user.

WEB AS A GRAPH

• Web can be viewed as a graph of connected web pgs.

• Anchor text & link to a web resource is used to link these pgs.

↓
jo link ke opre text hota hai — resource info.

• Outlink → hyperlink out of a pg.

• Inlink → hyperlink into a pg.

BOW TIE MODEL

→ First pg (home pg. is known usually)

→ Home pg → Inpage → (goes through seq. of pg ssc (round)) & then you reach outpg.

* we can go back from out to ssc.
(Surfing through different web pg)

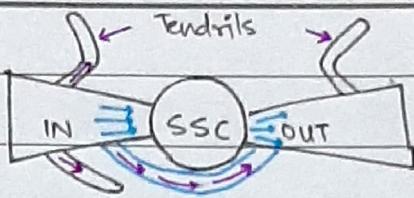
→ there's a tube & → classified channel (dark / classified web pgs) → not all users can access (pornography, data selling)

→ Tendrils → incompatible / not known

→ Has pg. mcn i links hotay hain & they grow exponentially on avg.

↓ Power law states that these inlinks & outlinks are governed with a power functions. Inlink proportional to $\frac{1}{\sqrt{i}}$.

Day / Date: _____



SSC → highly connected pgs.
that surf you among pgs.

Tendrils → dead end web pgs → jahan se return nahi karta hain; can only be accessed by only few people — not everyone allowed.
Web as a graph

• Spam (in context of web search) is the manipulation of web pg content for the purpose of appearing up high in search result for selected keywords.

• Web pgs rank highly on selected keywords.

• Adversarial Info Retrieval is the area of study making a balance b/w SEO & SE
→ competitor ki info build using search engineering.

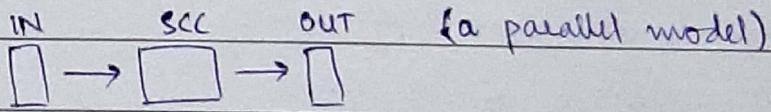
SEO → to provide consultancy services for clients who seek to have their web pgs ranked on selected keywords.

IN → home pg → initial

SCC → tightly coupled pgs → they move us to OUT

→ Stochastic Model

Prob of opening a link



→ Adversarial IR : Area of study making a balance b/w SE & SEO

Organic → free web / content only Paid → Adds advertisement too

Advertising as an Economic Model

→ Ad words are located in query & the ad are pushed to make some business value.

→ Typically, the ads are priced on cost per mile (CPM).

→ cost per click (CPC) → cost per user action (CPA)

→ Understanding how search engines do the ranking & how to allocate marketing campaign budgets to different keywords & to different sponsored search engines for maximum (SEM)

→ Click Spam — a phenomenon for advertising competition (SEM)

Dazzle

Day / Date:

- Pure vs. Sponsored: Pure — only displaying content of web, not paid content (ads)
extra
Mix — pure + sponsored — ads + content
- Combining Pure vs. sponsored
- Blog vs. Product reviews

USER NEEDS

- Informational — want to learn about something (40–65%) e.g. how to
- Navigational — want to go to some page (25–15%) e.g. PIA
- Transactional — want to do something (35–20%) • Access service, download, shop
- Gray areas — find a good hub; exploratory search 'see what's there'. • rental agent

WEB INFORMATION DISCOVERY

→ Directories

- Taxonomies populated with web pages in categories e.g. Yahoo!
- The user to browse through a hierarchical tree of category label.

→ Search Engine

- Full index search such as Infoseek.
- The user with a keyword search interface supported by inverted indexes & ranking mechanism.

Web Characteristics

① Web User Interaction

② Web as a Graph

③ Web Spam

↓
Anchor text (along with link)

Top 10 Search Engines

① Google

② Mahalo

③ Yahoo

④ Bing

Index Size & Estimate

→ Capture / Recapture Method — Pick a random page for index of EI & check/test whether it

Dazzle

Day / Date: _____

is in E2's index & vice versa. These experiments give us fraction $\alpha \leq 1$ such that α frac. of pgs in E1 are in E2 while $1 - \alpha$ frac. of pgs in E2 are in E1.

→ Then $|E_1|$ denotes the size of index of search engine E; $\therefore \alpha |E_1| \leq |E_2|$.
(All search engine ke mugalabhy mein kis ² main se konse engine engines, pc zyada data index kar raha)

$$\frac{|E_1|}{|E_2|} \approx \frac{1}{\alpha} \rightarrow \text{konse engine zyada pgs index kr raha}$$

Sampling Methods (How to choose pgs for testing)

- ① Random Searches (Randomly selected; have biasness)
- ② Random IP addresses
- ③ Random walks
- ④ Random queries

Duplicate / Near Duplicate Detection

→ While indexing, we may come up for duplicate. Checksum is common method to detect a duplicate. (Full duplicate have 90-100% similar data)

Near duplicate — not identical, but a portion is common, based on preset threshold we can filter out the near duplicates (50% + content same)

Shingling — Given a trc integer K & a seq. of terms in doc d , K -shingling is a set of all consecutive seq. of K terms in d .
used to find near duplicate.

e.g. A = a rose is a rose is a rose

shingling set S = { a rose is a ; rose is a rose ; is a rose is .. } ; $|S| = 5$.

Jaccard = $2/5$ b/w $|A| = 4$ (took 4 terms at a time, then next 4).

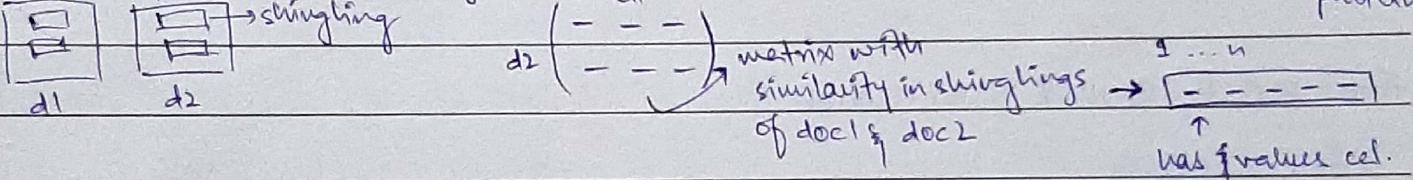
→ Duplicate can comparatively be easily found by checksum, whereas it is challenging to find near duplicate.

Day / Date: _____

Method → calculate shingles of each pg/doc. Then put them in a hash-table & calculate similarity b/w them to see if the pgs/doc are near duplicates or not.

These pgs are used for
Doorway Pages → Search engine use these pgs for optimization.

→ Near Duplicate Scaled Approach

- A pair wise approach — using shingling overlap to detect near duplicates.
- We can perform better by using a large integer hashing func & doing hashing for shingling pattern.


The diagram illustrates the process of generating shingles from two documents, d1 and d2. Each document is represented by a box containing a set of symbols. These symbols are then converted into shingles, which are represented as sets of bits (0s and 1s). Below this, a matrix is shown with rows labeled d1 and d2. The columns represent the shingles generated from d1 and d2 respectively. The matrix entries are binary values (0 or 1), indicating whether a specific shingle from d1 appears in a specific position of d2's shingles. This matrix is described as having a "similarity in shinglings" between doc1 & doc2. An arrow points from the matrix to a row labeled "hash value calc. using matrix".

shingling

d1 d2

shingling

d1 d2

matrix with similarity in shinglings of doc1 & doc2

hash value calc. using matrix.
- next time jab koi similar pair aye ga, we can easily search in hash table.
- we need hashing of minimum 2^{64} (High computational cost)

Important Topic

→ User Needs (Query types)

→ Index size & estimate

→ Duplicate & near duplicate

Day / Date: _____

WEB CRAWLER

(Chp 20)

- Web crawling is the process by which we gather pgs from the web to index them & support a search engine.
 - Objective of crawling is to quickly & efficiently gather as many useful pgs as possible, together with the link structure that interconnects them.
 - web crawler is sometimes referred to as Spider
- Steps* process starts by URL queue, then url generator generates a DNS & send to exact IP addresses.
this is a two way communication.

FEATURES A CRAWLER MUST PROVIDE

① ROBUSTNESS

The crawler must be robust to deal with a large no. of linked pgs from a website.
Sometimes a server traps a crawler, and the crawler is stuck b/w giving two pgs.
∴ The crawler must identify these traps.

② POLITENESS

Web servers have both implicit & explicit policies regulating the rate at which a crawler can visit them. These politeness policies must be respected.

③ DISTRIBUTED

The crawler should have the ability to execute in a distributed fashion across multiple machines.

④ SCALABLE

The crawler architecture should permit scaling up crawl rate by adding extra machines & bandwidth.

Day / Date: _____

⑤ PERFORMANCE & EFFICIENCY

The crawl sys should make efficient use of various sys resources including processor, storage & network bandwidth.

⑥ QUALITY

Given that a significant fraction of all web pgs are of poor quality/utility for serving user query needs → the crawler should be biased toward fetching 'useful' pgs first.

⑦ FRESHNESS

In many apps, crawler should operate in continuous mode. It should obtain fresh copies of previously fetched pgs.

⑧ EXTENSIBLE

Crawler should be designed to be extensible in many ways — to cope with new data formats, new fetch protocols & so on. Crawler archit. should be modular.

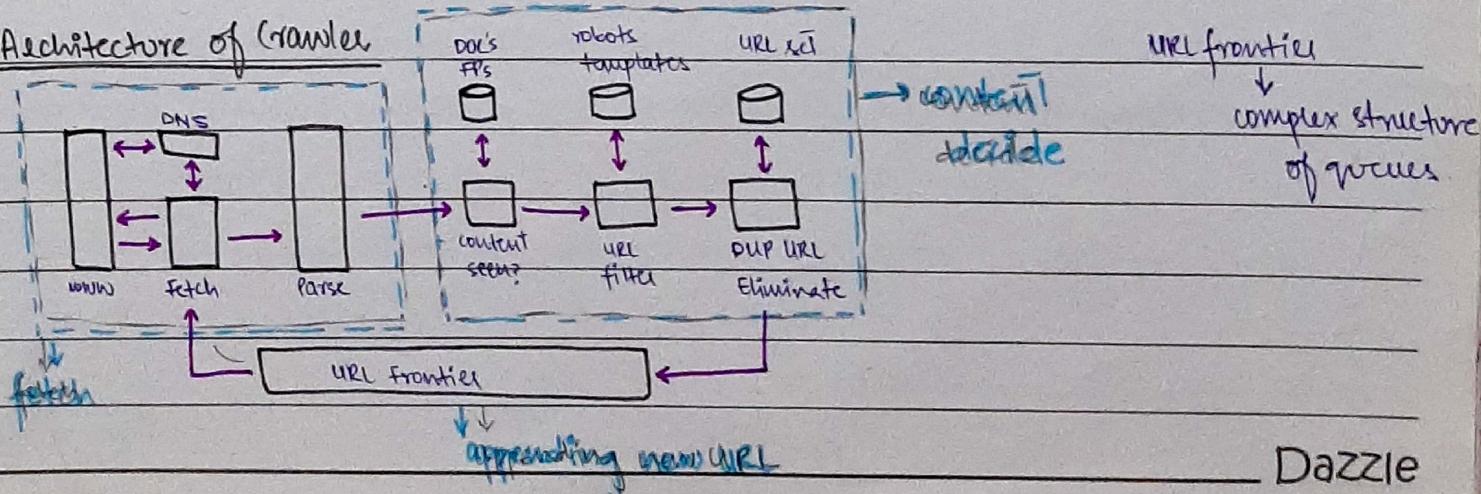
TYPES OF CRAWLER

- Path ascending crawler
- Topic focused crawler
- Academic focused crawler
- Semantic focused crawler → has aspect like search like dega (360°)

only focused on certain content (e.g. google scholar)

Academic focused crawler

Architecture of Crawler



Dazzle

Day / Date: _____

URL frontier → containing URLs yet to be fetched in the current crawl.

At first, a seed set is stored in URL frontier, & a crawler begins by taking a URL from the seed set.

DNS → domain name service resolution. Lookup IP address for domain names.

Fetch → generally use the HTTP protocol to fetch the URL.

Parse → The pg is parsed. Texts (imgs, videos) and links are extracted.

Note: Read the slip from book for more details.

Day / Date: _____

LINK ANALYSIS

(Chp 21)

WEB AS A GRAPH

Link analysis of web (as a graph) is based on following 2 assumptions:

① Anchor text pointing to pg B is a good description of Pg B

② Hyperlink from A to B represent an endorsement of Pg B by creator of Pg A.

(→ This is not always the case) i.e. many ~~links~~ links amongst pgs within a single website stem from the use of a common template.

↓ endorsement can be for a single domain.

CHALLENGES

- Pg might not provide an accurate description of itself.

- There is a gap b/w terms in a web pg & how a web user would describe web pg.
(terms don't define content of web pg)

- Web pg is a composition of text, graphics & img. Standard IR approach doesn't support searching with these rich contents.

- Anchor text along with extended info from anchor extended gives a good idea of web pgs.

* Anchor text means the text along with a link (description of a link of pg)

& extended anchor text means a little content from above & below anchor text.

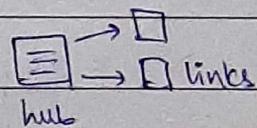
HITS (Hyperlink-Induced Topic Search)

→ Also known as hubs & authorities is a link analysis algo that rates pgs on web pages.

• To calculate/decide the score of each pg; hub score & authority score.

hub → pg that mainly contains outlinks (pgs it is referring)

authority → the web pg is the best link of a domain



→ [] ← most of referred link. Other pgs have used it as an inlink. Other pgs (hub) are referring to this pg.

Dazzle

Day / Date: _____

- A good hub represents a pg that pointed to many other pgs
- A good authority represents a pg that was linked by many different hubs.

* The algo assigns two scores for each pg ; authority & hub value.

Authority → estimates the value of the content of the pg.

Hub → estimates the value of its links to other pgs

HITS Algorithm

- ① Given a search query Q, collect a subset of 200 top web pgs as input (seed) that contain highest frequency of query Q.
- ② Add the collection of web pgs that points or are pointed by these top 200 web pgs.
Create adjacency matrix A among these web pgs.
- ③ Initialize the hub & authority column vectors U & V with values 1.
- ④ For a set k no. of iterations, do the following:
 - (i) Update the authority scores through authority matrix V.
 - (ii) Update the hub scores through hub matrix U.
 - (iii) Normalize the hub matrix and authority matrix U & V.
- ⑤ Rank the web pgs acc. to the authority score as reflected through authority matrix V.

Example (in video)

a → authority score

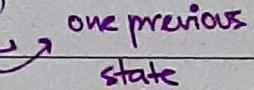
, h → hub score

$$a^{(1)} = A^T \cdot h^{(0)}$$

$h^{(0)}$ → hub score at iteration 0.

$$h^{(1)} = A \cdot (a)^0$$

$$h^{(0)} = a^{(0)} = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} \quad \text{or} \quad \begin{bmatrix} 1/n \\ 1/n \\ 1/n \\ 1/n \end{bmatrix}$$

for eg. $a^2 = A^T \cdot h^{(1)}$ 

Day / Date: _____

$$\begin{bmatrix} 1 \\ 3 \\ 5 \\ 1 \\ 2 \\ 1 \end{bmatrix} = \begin{bmatrix} 1/(1^2 + 3^2 + 5^2 + 1^2 + 2^2 + 1^2) \\ 3/(1^2 + 3^2 + \dots) \\ \vdots \end{bmatrix}$$

normalized

→ Authority score zyada tou most authority pgs.

→ If we update auth(p) to be summation :

$$\text{auth}(p) = \sum_{i=1}^n \text{hub}(i) \quad ; n \rightarrow \text{total no. of pgs connected to } p$$

$i \rightarrow \text{a pg connected to } p$

$$\text{hub}(p) = \sum_{i=1}^n \text{auth}(i)$$

matrix multiplication

ISSUES OF HITS

- ① expensive in sense of computing time.
 - ② Query dependent
 - ③ executed at query time, not at indexing time.
 - ④ Not commonly used by search engines
 - ⑤ It computes 2 scores, not 1.
 - ⑥ It is processed on a small subset of relevant documents (a focused subgraph / baseline), not all docs as was the case in pg rank algo.
- * With the help of query expansion, we can find other relevant & related docs to the query term.

Note: See numerical at end of slides.

$h^0 \text{ or } a^0 [1] [1,1]^0$ $a' = A^T \cdot h^0 \rightarrow$ first iteration of authority score
all 1's $\bullet h' = A \cdot a^0 \rightarrow$ $\dots \quad \dots \quad \dots$ hub score.
Day / Date: _____
eg. $A^2 = A^T \cdot h'$

Query extension is done so that

it includes the pages that are pointing towards page that are ranked by query.

In a lot of cases anchor text is not very clear so if query term is not present, we want to include those pages too; those pages which are tightly linked to the pages retrieved (either in a authority way or a hub way)

query pg filter kya that those pages that are have high freq. of query term.

PAGE RANK

→ no. of web pages → N

→ Anchor text → describes link & have hyperlink URL
① forward link $p_1 \xrightarrow{\text{to}} p_2 \rightarrow$ Outlink
② backward link $p_1 \xleftarrow{\text{from}} p_2 \rightarrow$ Inlink

Page Rank is also a link analysis algo

assigns a numerical weighting to each element of a hyperlinked set of documents → that pg. ki importance bataa hei pg. rank.

Damping factor → aik pg. pc pochan jayga without any link-probs. that a user opens a new web pg. to begin random walk.

Day / Date:

Summery

$\text{PR}(P) \rightarrow$ pg rank of pg. P

someday.

$\deg(P)^- \rightarrow \square^-$

$\deg(P)^+ \rightarrow \square^+$

$N(P)^- \rightarrow 1$

$N(P)^+ \rightarrow 2$

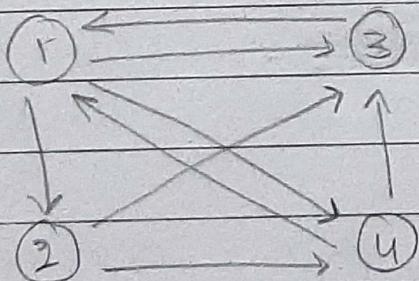
similar.

mornings.

$W \rightarrow$ hyperlink matrix representing network whose entries constitute the fractional page rank contributions.

$w \rightarrow$ eigenvector containing ranks for each vertex in the network.

d is kept from $0.1 - 0.2$. And usually $d = 0.85$



$$\text{PR}(1) = \frac{2}{3}$$

P

$$\gamma_a = \frac{\gamma_4}{2}$$

$$\gamma_y = \frac{\gamma_4 + \gamma_a}{2}$$

$$= \frac{\gamma_4}{2} + \frac{\gamma_4}{2} \times \frac{1}{2}$$

$$\gamma_y = \frac{2\gamma_4}{2} \times \frac{1}{2}$$

$$\gamma_y = \frac{\gamma_y}{2} + \frac{\gamma_a}{2}$$

$$\gamma_y = \frac{\gamma_y}{2} + \frac{\gamma_4}{2}$$

$$\frac{\gamma_4}{2} = \frac{\gamma_4}{2} \times \frac{1}{2}$$

$$\frac{\gamma_4}{2} = \frac{\gamma_4}{4}$$

Dazzle