9 and 10 Answer pg + 08 (disadvanlages & charlenges).

Chap #02 The -lerm vocabulary and posting lists

- \* Tokenization -> Process of chopping character streams into tokens. Linguistic proprocessing that deals with building equivalent class of tokens.
- \* Document processing pipeline
  - O Collect the documents to be indexed
  - @ Jokenize the text
  - 6) Do linguistic preprocessing of tokens
  - @ Index the documents that each term occurs in

\* Challenges in Document processing:

- O What characles sel is required for encoding eg: ASCII encoding for plain english text

  Mulli-byle or single byte encoding scheme like UTF-8.
- 1 What language the document is written in?
- 3 What formal is it in? eg: pdf/html/word

entres to consider for a certain query term. eg car search krenge to automobile ki posting list ka intersection bhi lega to caler synonyms.

@ Perform the expansion during index construction. eg when a document contains 'automobile' , we inder it under car' as well.

Mainlaining relations of unnormalized lerms is less efficient because there are more postings to store 2 merge These expansions are asymetric.

- Normalizing Totens to remove diacritics.

- lase folding: Feduce all letters to lower case.

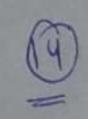
\* Morphological Analysis: (Forms of things).

1 Inglections: Adding suffix to a mord, that doesn't change its grammatical category eq: Jenses in verbs Plural in nouns.

Desivations: Adding sulfin to a word that changes its grammatical certegory.

eg: Noun > adjective > verb.

Dictionary	Thesaurus
O Corrlains an alphabélical list	1 List of words with synonyme
of words including meanings,	and antonyms, arranged
elymology and pronunciation.	and antonyms, arranged either alphabetically or themakally
@ 9t desines a word.	19 9t gives a list of words with similar meanings.
	with similar meanings.
3 Some dictionaries use the	3 There are various therauri some
International Phonetic Alphabet	of which may give additional
	into like phrases, sentences
	and clouses.
Stemming	Lemmalization.
1) Refers to crude heurestic	1) Refers to chopping properly
process that chops off	with the use of vocabulary
the ends of words.	7 morphological analysis of
	words.
@ Operates on a word	@ Removes andlection andings
the context.	2 Removes and rection andings.
the context.	



Date\_\_\_\_

## \* Normalization:

- We need to normalize words in documents as well as in queries, so that the lokens in the query exactly match the tokens in the tokens in the tokens in the

- Token normalization is a Tcanonicalizing tokens so that matches occur despite superficial differences in the characles sequences of the tokens.

- The most standard way is to implicitly create and equivalence classes of terms

eq: deleting periods to form a term

U.S.A - USA

eg: deleting hyphens anti-discriminatory, articliscriminatory

Normalizing loters to remove diacritics

Case folding: Reduce all letters to lower case

This can be done in two ways;

O Index unnormalized tokens and maintain

a query expansion list of multiple vocabulary

nice

10 Choosing a document unit . Selecting too large or too small document affect precision and recall.

1 Totenization is challenging

- Difficult to caler hyphenalion, apostrophe

- Compound inords

- Abbreviation, numbers, phone-number ele

- language identification.

eg: splitting on basis of white spaces is easier for english language only Challenges arise while dealing with other languages.

- Splitting on white space can also split what should be regarded as single token

eg: San Francisco.

\* Stop mords :=

- Extremely common words are called stop words

- We usually drop stop words entirely from dictionary since they have little importance and

- But modern lechnique do not autour lo throw

any thing from documents.

- Stop words play important tole in some special queies like phrase queies, relational queries, little of songs ele

3 stemming is faster, easier to implement and reduces accuracy.

3) Lemmalization is slowler, but accurale.

(4) Example: Operational Result: Oper, O, Op elc. Result: Operale.

(4) Example : Operational

Conclusion: The goal of both stemming of temmalization is to reduce inflection forms of some times desivationally related forms of word to a common base form.

-> Porter Stammer.

Consist of 5 phases of word reduction, applied sequentially Within each phase - here are various conventions to select rules.

See slides 15-18



## \* Faster postings list intersection via skip pointers:

- 91 the postings list lengths are mand n, the intersection lakes O(m+n) operations.
- Can we do intersection in sub linear lime?
- We can, if index is n't changing too fast

## Skip list:

- Augmenting postings list with skip pointers forms skip list.
- Skip pointers are shortcuts that allow us to avoid processing parts of the postings list that will not figure in the search result.
- Where to place skip pointers?
- More skips means; shorter skip span

lots of space storing skip pointers.

- Feinler skips means

long skip span - Penier comparisons

Ferrer opportunitées to skip.

- Skip pointers are easy to implement, is indices are static.
- Only use ful for 'AND conjunctive queries.

Unhy skip pointers are not useful for disjunctive queries (XOR 4)?

Because in queries of form "x OR Y", it is essential to visit every doc ID in the postings lists of either terms, thus no need for skip pointers

Phrase Queries

- We mant to be able to answer queries in the form of opposites phrases Like Stanford University.

- Simple postings lists require post processing to eliminate false positive results reducing elliciency.

\* Solution 1 - 10 Phrase Queries: Bimord indexes - Index every consecutive pair of terms in the lext

as phrase.

- Each of these bimords is now a dictionary term.

- Dino-inord phrase query processing is non immediate - longer phrases can be processed by breaking phrases into conjunctive bimord bookean query.

eg: S-lanford university palo alto "Stanford university" AND "University Palo" and "Palo Alto". Results in False positives.

nice

Extended bimords:

- Dokenize the text & perform parts of speech

-lagging.

- Group terms into nouns (N) and function mords (X)

- Any string of form NXXN called extended binord.

Issues:

- False positives

- Index blow up due lo bigger dictionary.

\* Solution 2 to Phrase quesies: Positional indexes - In the posternes, store for each term the position (s) in which tokens of it appear:

term | doc-frequency -> Di: Position 1, 2, 3 -> Dz Position 45

eg < (hello): (200), doc frequency

docid - 1: 9, 11, 14, ...; positions

2: 3, 4 . - - ;

3:6,7,8 ---;

- Can handle proximity queries.



eg to be or not to be

- Extract inverted index entries for each distinct term. ie: to , be, or, not
- Look for documents containing 'to' and be
- look position of 'to' en a document,
  then look if position of be occur one higher than position of 'to'. This means both lerms are adjacent.
- Process all terms in a similar manner.

\* Proximity Queries:

- Have k words in between wi; and wij
- /K mean "within K words of"
- Only positional indexes can be used for such queries

\* Positional index size :

- Positional index expands postings slovage
- Useful because phrase & proximity queries can be handled easily.
- Need an entry for each occurrence in a document.
- Index sixe depends on aug document sixe.

\* Combination Schemes:

- It is a combination of biword indexes and positional indexes.

- For queries like "Michael Jackson" it is ine ficient to keep on merging postings lists

- This -lechnique uses bimord indexes for certain queries & positional indexes for other phrase queries based on d'Oberent criteria.

- Gives speed up

Next word index - Another scherne

< Food for thoughts> chap # 02.

1 Answer pg # 12 3 13

3 Answer pg # 16 7 17

(4) (a) No. of comparisons = 11

Compares till we reach 47 in upper postings list

(b) No of comparisons = 6. list 2 -> Skip pointer of length = 116 = 4

3 AnsiNer pg # 20

**3** 

O-Positional index extracts inverted index for each distinct wi and wij.

each distinct wi and wij.

look for document d'containing vin; at index m and wij at index m+ 12, where k is equals to 2 for phrase query

eg Q = hello world

Output:

00000000

9

9 9

.

9

9

9

9

9

9

d, b/c hello at position 3 7 would at position 4 de b/c hello at position 4 to would at position 5.

① pg # 22