

Q-1 Descriptive Questions

Q-1 State the tasks involved in preprocessing text data with suitable examples -

Ans - After we have gathered good amount of data, we need to do preprocessing before inverted index construction. Following steps can be used -

① Tokenization → we need to convert documents into stream of tokens and then proceed with index creation. These tokens can then be used in query processing.

② Normalization → Remove punctuations can help save a lot of memory and vocab list can be made efficient.

③ ~~Re~~ Lowercasing/Uppercasing → for efficient query processing, it is always better to convert terms to a particular case (either upper or lower)

④ Removing Stop Words - The most frequent occurring words (example - a, the, in, of, at) are called stop words. They can be removed before creating index being less significant

⑤ Stemming - Transforming the word into its root word (example - mostly to most), removing s/es or any stemmer can help make better index

Name - Sayam Kumar
Roll No - 20180010158
(Page No 2)

⑥ Lemmatization → we can also create mappings from completely different words to their root form.

Eg → ate → eat, went → go

Necessity → We always look for building a scalable Information Retrieval system. Therefore in order to efficiently store and access indices, it's necessary to complete these steps. We cannot afford user waiting considerable long time for their queries.

Q-2

Assuming 1-based indexing, no stemming, no removal of stop words

Shipment →

d1: 0 1
d3: 1

of → d1: 2 | d2: 2 | d3: 2

platinum → d1: 3 | d4: 1

damaged → d1: 4 |

in → d1: 5 | d2: 8 | d3: 7

Delhi → d1: 6 | d2: 6 | d4: 10

India → d1: 7 |

delivery → d2: 1 |

a → d2: 3, 9 | d3: 8

parcel → d2: 4 | d3: 5

for → d2: 5

arrived →

d2: 7 | d3: 6

new → d2: 10 | d3: 4 | d4: 9

truck → d2: 11 | d3: 9 |

my → d3: 3

and → d4: 2

parcels → d4: 3

are → d4: 4

the → d4: 5

shipment → d4: 6

today → d4: 7

from → d4: 8

Q-3 Ans → wildcard queries
contains * to match words

e.g. - fun* to match words
starting with fun

We can create B trees containing backwards addition
words. for queries with abc*, B trees are efficiently
faster and yield results b/w abc \leq w $<$ abd.

① For queries ending with *, (e.g. * abc)

we can look upon B tree and find cba \leq w $<$ cbb.

② It is very expensive operation to find
queries with (*) in the middle.

Permutation Index

① Create * ana at the end and store words in
backwards fashion.

for example word tiger, maintain dollar as ending
 \rightarrow tiger\$, regi\$t, reg\$ti, se\$tig, s\$tige,
\$tiger

Lookups

X*	\$X*
X	X\$
@ X *	X *

Comparison

Permutation index is very efficient in retrieving
queries by requires a lot more space (quadruple in
size). Instant time comparison take place in
permutation index

Name - Sayam Kumar
Roll No - S20180010158

(Page No 3)

Q-4 Biwords

In biwords, we index every consecutive pair of terms in text.

e.g. 'hello world people' can be indexed as -

- hello world
- world people

(*) We keep biwords as a dictionary term

Advantages

- ① Two word phase query processing is fast
- ② Even longer queries can be broken down to 2phase (biword) queries
- ③ Speed of retrieval is increased
- ④ More contextual information and relevance can be formed

Co, seeing ups and downs, these all are not part of standard solution.

Name-Sayam Kumar
RegNo-S20180010158

Page No-4

Disadvantage

- ① Space is increased
- ② False +res can be there
- ③ Not implemented as standard solutⁿ.
- ④ \rightarrow
'hello world' and
'world hello' are stored differently.

Sayam Kumar
 Roll No - S20180010158
 Page No 5

Q-5

Term	freq	$x w$	$P(x w)$	$P(w)$	$P(x w)P(w)$
actress	321	c/ct	0.000117	0.012	1.4×10^{-6}
cress	80	a/#	0.000144	0.003	4.3×10^{-7}
careess	121	ac/ca	0.000164	0.004	6.5×10^{-7}
access	527	r/c	0.000209	0.021	$4.3 \times 10 \times 10^{-6}$
across	803	e/o	0.00321 0.00093	0.031	2.8×10^{-5}
acres	412	es/e	0.00321	0.005	1.6×10^{-6}
acres	412	s/s	0.00342	0.005	1.7×10^{-6}

Best word = across (2.8×10^{-5})

Q.6 5 term weighting
example -

Name - Sayam Kumar
Roll No - S20180010156
Page No - 6

① It is important to have term weighting to measure relevance of a term in a document. This also leads to good results if done correctly.

① Raw TF factor → Just keep the no. of occurrences.
Ex - if standford occurs in all documents, then it is ranked high.

② Logarithmic Scale - computed as

$$1 + \log_e (tf)$$

This helps in normalizing out high freq. terms with lower ones

③ Binary TF factor - Whether or not a term is there in a document. Either 0/1, Count and relevance is not stored. Results in poor performance.

④ Augmented TF factor - Reduces the range of frequencies to (0.5, 1) providing a normalized view of all terms. calculated as

$$0.5 + 0.5 \frac{tf}{\max tf \text{ in document}}$$

⑤ Okapi's TF factor → Approximations of 2 Poisson Model

$$\frac{tf}{2+tf}$$

In all of the above, log scale is better.