

Q1-a) Use fingerprints to remove all but one copy of identical documents.

Remove common HTML tags and integers from the single computation, to eliminate shingles that occur very commonly in documents without telling us anything about duplication.

Use union-find algorithms to create clusters that contain documents that are similar.

Calculate  $d_i$  and  $d_j$  to check similarity.

Compute # of shingles in common for any pair of documents whose sketches have any members in common.

For list  $\langle x, i, d_i \rangle$  sorted by  $(x, i)$  pairs, generate  $j$  for all  $i$  whose  $(x, i)$  is present in both their sketches. For each pair  $i, j$ , get a non-zero sketch overlap, a count of # of  $(x, i)$  values.

Use a threshold to find pairs  $(i, j)$  that have overlapping sketches.

Run union find to group documents into near duplicate 'syntactic clusters'.

In this way, we can use single-link clustering algorithm.



Q5. a) No, we cannot say that documents are duplicates as jaccard similarity is bag of words based approach and doesn't consider word order.

eg. " I live in Delhi but- I often visit Gurgaon. "  
" I live in Gurgaon but- I often visit Delhi. "

have jaccard similarity = 1 but are different in meaning

(b) Log function allows damping effect.  
Any function that is monotonic and grows slowly compared to a linear function can be used.

eg.  $x^{1/k}$  where  $k=1$

(c) If we normalise  $w_a, w_b, w_c$  before computing  $(w_b - w_a + w_c)$  then the result will be normalised.

Adding  $\rightarrow$  Normalising

$$p = i$$

$$q = j$$

$$p+q = \frac{i+j}{\sqrt{2}}$$

$$= \frac{i}{\sqrt{2}} + \frac{j}{\sqrt{2}}$$

$$|p+q| = 1$$

Normalising  $\rightarrow$  Adding Individually

$$\frac{i}{\sqrt{2}} \quad \frac{j}{\sqrt{2}}$$

$$\hat{p} + \hat{q} = \frac{i+j}{\sqrt{2}}$$

$$|\hat{p} + \hat{q}| = 1$$

Radhika Patwari

18CS10062

Q3. a) Heap's law =  $(M = kT^b)$   
= estimating  
vocabulary size given collection size

$M$  = vocabulary size

$T$  = # of tokens in  
collection

$b \approx 0.5$

$30 \leq k \leq 100$

Zipf's law = the  $i$ th most frequent term has frequency  
proportional to  $1/i$ .  
= estimating collection frequency (cf) given  $cf_i \propto \frac{1}{i}$   
for term rank  $i$

Both laws are power laws but heap law talks about  
size of vocabulary and zipf's law talks about  
distribution of the vocabulary.

Q3. c) (i)  $\frac{1}{8} \left( \frac{1}{1} + \frac{1}{1} + \frac{3}{9} + \frac{4}{11} + \frac{5}{15} + \frac{6}{20} + \frac{7}{21} + \frac{8}{22} \right)$

↳ Max MAP

↳ common part -

assuming relevant docs are at rank 21, 22)

↓  
we also disregard relative relevance score &

just consider relevant / non relevant.

min MAP (relevant docs are at 9999 and 10,000)

$$\frac{1}{8} \left( 2 + \frac{3}{9} + \frac{4}{11} + \frac{1}{3} + \frac{3}{10} \right) + \frac{1}{8} \left( \frac{7}{9999} + \frac{8}{10000} \right)$$

Q3. c)(ii) Only NDCG can be computed as we can just determine ranking of documents and not the absolute relevance.

To compute it we can merge pairwise rankings (if possible, i.e. common elements exist) to recover total ordering of docs and compute NDCG.

Q3. (b) Backup approach

Keep a backup copy of every auxiliary index  $Z_i$  in the disk as  $Z_i$ -backup.

Now while retrieving, we will check if  $Z_i$  is present. If not present ( $Z_i$  is deleted), we can use  $Z_i$ -backup. If it is not present, then  $Z_i$  does not exist. If  $Z_i$ -backup is present, then it implies  $Z_i$  is deleted. We use  $Z_i$ -backup.

This method has search overhead but requires 2 if condition only.

18CS10062

$$\begin{aligned}
 Q4.(a) \quad O(R|q, \vec{x}) &= \frac{P(R=1|q, \vec{x})}{P(R=0, q, \vec{x})} \cdot \cancel{P(R=1|q, \vec{x})} \\
 &= \frac{P(R=1|q) \cdot P(\vec{x}|R=1, q)}{P(R=0|q) \cdot P(\vec{x}|R=0, q)}
 \end{aligned}$$

Now if probability of term  $x_i$  is dependent only on previous term  $x_{i-1}$ , then this is a Markov assumption.

$\therefore$  Under Markov assumption,

$$\begin{aligned}
 O(R|q, \vec{x}) &= O(R|q) \cdot \frac{P(x_1|R=1, q)}{P(x_1|R=0, q)} \\
 &= \frac{\prod_{i=2}^n P(x_i|R=1, q, x_{i-1})}{\prod_{i=2}^n P(x_i|R=0, q, x_{i-1})}
 \end{aligned}$$

Q4.(b) Let vocabulary size =  $|V|$

$\Rightarrow$  BIM model

Every token / word  $x_i$  has 2 parameters  $r_i$  and  $p_i$

$\therefore$  Total # of parameters =  $2 \cdot |V|$

$\rightarrow$  Present Model

$P(x_i|x_{i-1})$  and  $r(x_i|x_{i-1})$  are required for each

token pair.

$\therefore$  Total # of parameters =  $2 \cdot |V|^2$



Radhika Patwari

18CS10062

⑥

Q2.(b) Semantic relations are captured by word pairs.  
Sharing common offset vectors (semantic relations are captured by certain directions in the vector space)

Another implicit assumption is that a word isn't ambiguous (no polysemy)

eg. Bank - money (place - product analogy)  
might not match with

BMW - car

because BMW is clearly an automobile company but bank can be a place or riverbank.