

MS Data Science: Information Retrieval and Text Mining (Fall 2020)

Mid-Term Exam

Instructor: Dr. Saeed Ul Hassan

Name Qazi Danish Ayub

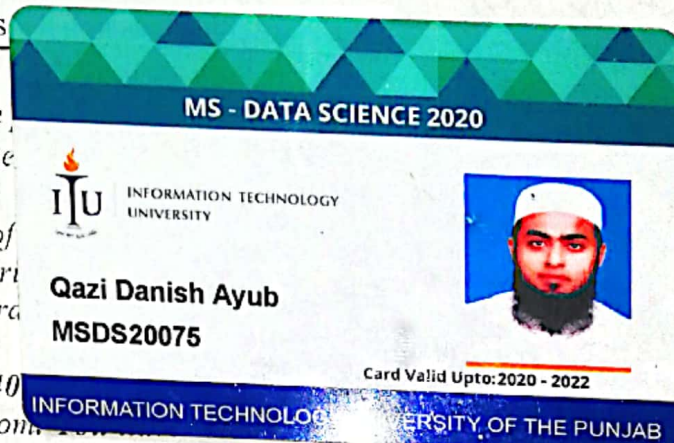
Roll No. MSDS 20075

Date 08-Dec-2020

Time Allowed: 2 Hours

Total Marks: 60

- Use plan A4 size
- Please follow the sheet.
- The first page of students must write if student ID card are visible
- An additional 40 google classroom following proper naming convention i.e. [rollnumber].pdf (e.g. MSDS19016.pdf).
- Please note no extra time will be given and wrong submissions will not be considered.
- Best of luck!



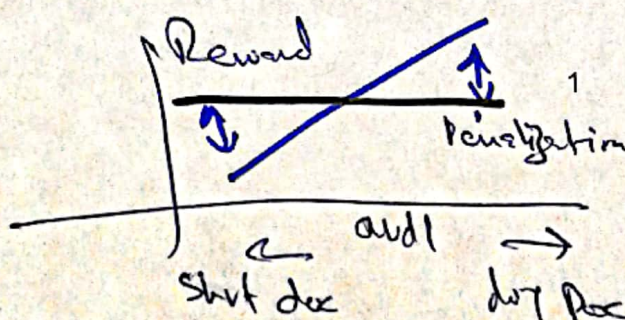
Q No 1

(4)

Long Doc — better chance to match any Query
Need avoid over Penalization

So Pivot length normalizer: take average doc length as "Pivot"

Norm = 1 if $|d| = \text{pivot}$ (and 1)



be $\{0, 1\}$

Q No 1

(4)

Answer:

Pseudo Relevance Feedback ..

A method for automatic local analysis on manual part of relevance feedback such that the user gets improved retrieval performance without extra interaction. It is also known as blind relevance feedback. → A unsupervised technique to some extent.

Relevance feedback ..

The involvement of user in the process for improving the result such that the user give feedback and the system computer better representation in that feedback is called Relevance feedback. → A pure supervised technique.

(d)

Ans.:

→ Stemming increases the total recall

Reason.:

According to the concept of stemming User entered term more documents are matched and as alternately words forms for a user entered terms are matched as well thus which produces the increase in recall - Although also reduces the precision -

(3)

Ans.:

Index Searcher:- Is a Gateway to Searching an index - All searches come through an Index Searcher instance using any of the several overloaded methods -

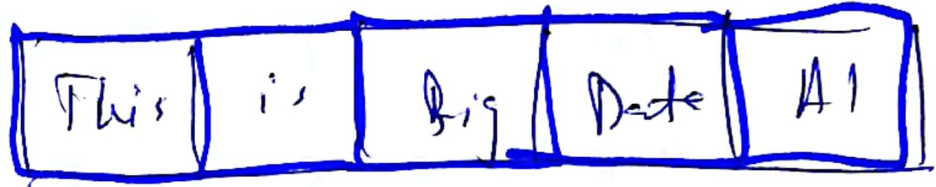
QueryParser:- Processes a human-entered (and readable) expression into a concrete Query object -

CON-1

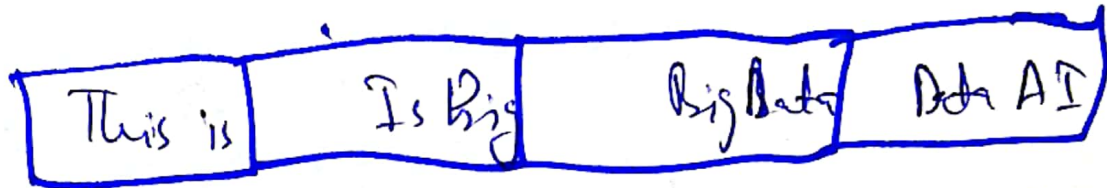
(5)

This is big Data AI

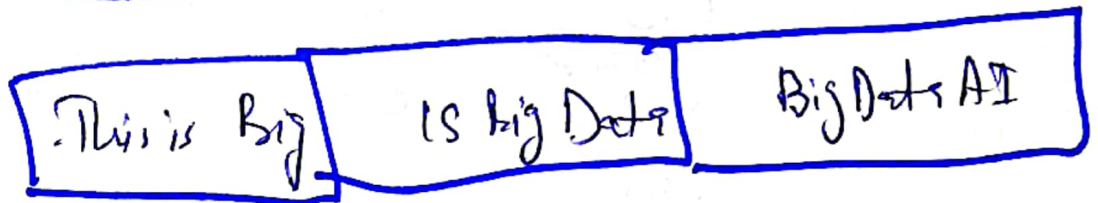
UniGram



Bi-Gram



Tri-Gram



Q2: Long Questions (50 Marks)

1. Provide answers for a, b, c and d parts

(10 marks)

a) Fill the recall and precision values for the below given documents w.r.t relevant documents found for each query? Draw similar table on your paper for both Query 1 and Query 2 separately. (3)



= Relevant Documents for Query 1

Ranking 1:



Recall:

$\frac{1}{5}$	$\frac{1}{5}$	$\frac{1}{5}$	$\frac{2}{5}$	$\frac{3}{5}$	$\frac{3}{5}$	$\frac{4}{5}$	$\frac{4}{5}$	$\frac{4}{5}$	$\frac{5}{5}$
---------------	---------------	---------------	---------------	---------------	---------------	---------------	---------------	---------------	---------------

Precision:

$\frac{1}{1}$	$\frac{1}{2}$	$\frac{1}{3}$	$\frac{2}{4}$	$\frac{3}{5}$	$\frac{3}{6}$	$\frac{4}{7}$	$\frac{4}{8}$	$\frac{4}{9}$	$\frac{5}{10}$
---------------	---------------	---------------	---------------	---------------	---------------	---------------	---------------	---------------	----------------



= Relevant Documents for Query 2

Ranking 2:



Recall:

$\frac{0}{3}$	$\frac{0}{3}$	$\frac{1}{3}$	$\frac{1}{3}$	$\frac{1}{3}$	$\frac{1}{3}$	$\frac{2}{3}$	$\frac{2}{3}$	$\frac{3}{3}$	$\frac{3}{3}$
---------------	---------------	---------------	---------------	---------------	---------------	---------------	---------------	---------------	---------------

Precision:

$\frac{0}{1}$	$\frac{0}{2}$	$\frac{1}{3}$	$\frac{1}{4}$	$\frac{1}{5}$	$\frac{1}{6}$	$\frac{2}{7}$	$\frac{2}{8}$	$\frac{3}{9}$	$\frac{3}{10}$
---------------	---------------	---------------	---------------	---------------	---------------	---------------	---------------	---------------	----------------

b) Find the Average Precision of Query 1

(2.5)

c) Find the Average Precision for Query 2

(2.5)

$$(b) \quad \frac{\frac{1}{1} + \frac{2}{4} + \frac{3}{5} + \frac{4}{7} + \frac{5}{10}}{5} = \frac{3.17}{5} = 0.6342$$

$$(c) \quad \frac{\frac{1}{3} + \frac{2}{7} + \frac{3}{9}}{3} = \frac{0.952}{3} = 0.317$$

$$MAP = \frac{0.6342 + 0.317}{2} = 0.4756$$

Long
Q No 2

(2)

(a)

Boolean Models

Deals with the
the Document Selection and thus Does not
contribute in Ranking as much [0, 1]

Vector Space Model

Deals with Term Frequency
Inverse Document Frequency and thus
also contribute is Document Ranking also
values is not boolean.

(b) Pseudo Code for removing Stop word punctuation
and Case Sensitiveness

① Dictionary - stopwords

② Dictionary.lower()

③ Dictionary.removePunctuation()

	TF-D ₁	TF-D ₂	TF-D ₃	doc _{he}	idf
representation	0	0	1	1	2.0
important	1	1	0	2	1.0
realizing	1	0	0	1	2.0
sub	0	0	1	1	2.0
semantic	1	1	0	2	1.0
knowledge	0	0	1	1	2.0
and	0	1	2	2	1.0
language	0	0	1	1	2.0
idf	0	1	0	1	2.0
technologies	1	1	0	2	1.0
several	1	0	1	2	1.0
language	0	0	1	1	2.0
web	1	1	0	2	1.0
amongst	0	1	0	1	2.0

c)

TF-IDF ₁	TF-IDF ₂	TF-IDF ₃
0	0	2.0
1.0	1.0	0
2.0	0	0
0	0	2
1.0	1.0	0
0.0	0	2
0.0	1.0	2
0.0	0.0	2
0.0	2.0	0
1.0	1.0	0
1.0	0	1.0
0.0	0	2.0
1.0	1.0	0
0.0	2.0	0

Q₁

D₃ → 2.0

D₂ → 1.0

D₁ → 0

Q₂

D₃ — 2.0

D₁ — 0

D₂ — 0

Long
Q No 3

$$v_{in} = \vec{\alpha v} + \frac{\beta}{|D_r|} \sum_{d_j \in D_r} d_j - \frac{\gamma}{|D_r|} \frac{\sum d_j}{\sum_{d_j \in D_r} 1}$$

Calculation Term Frequency

$$TF = \frac{\text{no. of rep of word}}{\text{no. of words in doc}}$$

$$D_1 = \{2, 1, 1, 1\}$$

the dog attached girl

$$D_2 = \{1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1\}$$

Using Frou

$$+ D_1 = \{0.4, 0.2, 0.2, 0.2\}$$

$$- D_2 = \{0.1, 0.1, 0.1, 0.1, 0.1, 0.1, 0.1, 0.1, 0.1, 0.1, 0.1\}$$

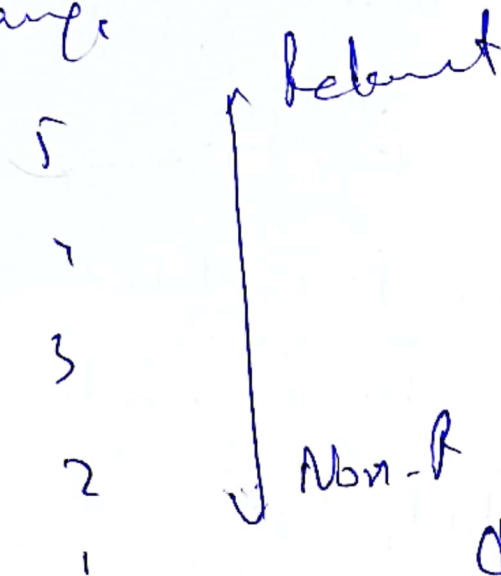
$$\alpha = \gamma = 1 \quad \beta = 0.5$$

$$\{1 + 0.2 \times 0.5 = 0.1, 1 + 0.2 \times 0.5 = 0.1\}$$

$$\{1, 1\} \text{ Ans.}$$

Q No 7 Long.

Range



$$Q_m = (a \cdot \vec{D}_0) + b \cdot \frac{1}{|D_r|} \sum_{\vec{D}_i \in D_r} \vec{D}_i - \left(c \cdot \frac{1}{|D_{nr}|} \cdot \sum_{\vec{D}_k \in D_{nr}} \vec{D}_k \right)$$

4 = Non-Relevant \vec{D}_k

5 = Relevant \vec{D}_i