



Beijing-Dublin International College



SEMESTER 2 FINAL EXAMINATION - 2022/2023

School of Computer Science

COMP3009J Information Retrieval

Dr. Robert Ross
Assoc. Prof. Neil Hurley
Dr. David Lillis *

Time Allowed: 120 minutes

Instructions for Candidates

Answer Question 1 and any two other questions. Question 1 has 30 marks available. All other questions have 35 marks available.

BJUT Student ID: _____ **UCD Student ID:** _____

I have read and clearly understand the Examination Rules of both Beijing University of Technology and University College Dublin. I am aware of the Punishment for Violating the Rules of Beijing University of Technology and/or University College Dublin. I hereby promise to abide by the relevant rules and regulations by not giving or receiving any help during the exam. If caught violating the rules, I accept the punishment thereof.

Honesty Pledge: _____ **Le Liu** _____
(Signature)

Instructions for Invigilators

Candidates are allowed to use non-programmable calculators during this examination.

Question 1:

- (a) Describe what is meant by “*information need*” in the context of Information Retrieval.

What are the different types of information need?

User only use IR when there is some information that they interesting.

Information need is motivation for using IR

Visceral Need, Conscious Need, Formalised Need, Compromised Need

[6 marks]

- (b) Below is part of a positional index relating to the term “same”. In creating this index, stopword removal and stemming have not been used. Postings lists begin at 1 for the first term in each document. Which document(s) could contain the phrase “the same day at the same time”? Explain your answer.

<same: 41825;

1: 120, 124, 167;

2: 9, 10, 13;

3: 121, 162;

4: 4, 101, 105, 106;

5: 1, 5, 88, 888, 889;

...>

1,,4could contain the phrase

Because they each have two occurances of "same" that are exactly 4 positions apart. While, In to must have "same same". In 5 same must at begain.

[6 marks]

- (c) A *modern Information Retrieval pipeline* may include Boolean searches, simple ranking and reranking based on machine learning. Explain why these are all useful to make an effective Information Retrieval system.

[6 marks]

- (d) The *BM25* method of Information Retrieval is based on the belief that a good term weighting scheme is based on three principles. Briefly describe each of these principles.

Inverse Document Frequency(IDF): Terms that appear in fewer documents are more informative.

Term Frequency(TF): The higher frequency, more important

Document Length: Long documents are penalized to avoid unfair advantage from term repetition.

[6 marks]

- (e) Compare and contrast the preprocessing steps of *stemming* and *lemmatisation*. In particular, what are the advantages and disadvantages of each?

Stemming is the process reduce the words to a common rot, often by suffix stripping.

[6 marks]

Adv: Fast, Simple, Reduces vocabulary size

Disadvan: Overstemming and Homographs

Lemmatisation is a NLP technique for converting word into lemmas

Adv: More accurate, return real word

Dis: slower, need more analysis

[Total 30 marks]

Question 2:

- (a) The *Boolean Model* makes use of the query operators *AND*, *OR* and *NOT*. Explain how these work and how they affect the number of documents returned by an Information Retrieval system. Also show how each of these can be implemented by using operations from Set Theory. 01b

AND is used to narrow a search. More AND, Fewer records

[6 marks]

OR is used to broaden a search, Documents contains any term specified will be return

NOT is used to specifically exclude a term from search. More NOT, Fewer records.

- (b) Briefly describe **two** ways in which the process of running Boolean queries can be optimised so that they can be processed more efficiently. 01c

1. Process in order of increasing of frequency: start with smallest, if reach the end of shorter list, we can stop. [6 marks]
2. Skip Pointer: it can reach later parts of list without interating through every element.

- (c) The *probabilistic model* of Information Retrieval makes use of two probabilities relating to query terms. These are $P(k_i|R)$ (the probability that a relevant document will contain the term k_i) and $P(k_i|\bar{R})$ (the probability that a non-relevant document will contain the term k_i). However, these probabilities cannot be calculated directly and must be estimated.

- (i) Briefly describe how initial values for these probabilities may be generated.
- (ii) Explain how these initial estimates can be improved with user feedback.

[8 marks]

- (d) Below is a small document collection, containing three documents. Answer the questions that follow.

1. Lowercasing all word.
2. Remove stopwords
3. Tokenisation

Stopwords: and, be, is, it, to, will

Document 1: It is going to rain and rain and rain today.

Document 2: Today I will be playing sport.

Document 3: I am going to watch the play.

- (i) Describe the preprocessing steps you would use when creating an index for these documents.
- (ii) Calculate a vector to represent each document, using the TF-IDF weighting system. You should use the stopwords list provided, but do not perform stemming.
- (iii) Calculate the cosine similarity for each vector using the query "going to play football". and show the final ranked list of documents for this query.

football", and show the final ranked list of documents for this query.

- (iv) What effect on the results would you see if you had used stemming for this corpus?

[15 marks]

[Total 35 marks]

doc1 : going, rain, rain, rain, today

doc2 : today, i, playing, sports

doc3 : i, am, going, watch, play

Vocabulary:

am, going, i, play, playing, rain, sport, today, watch

IDF	$\log \frac{3}{1}$	$\log \frac{3}{2}$	$\log \frac{3}{2}$	$\log \frac{3}{1}$	$\log \frac{3}{1}$	$\log \frac{3}{1}$	$\log \frac{3}{1}$	$\log \frac{3}{2}$	$\log \frac{3}{1}$
doc1	0	$\frac{1}{3}$	0	0	0	1	0	$\frac{1}{3}$	0
doc2	0	0	1	0	1	0	1	1	0
doc3	1	1	1	1	0	0	0	0	1
doc1 =	[0, 0.1949, 0, 0, 0, 1.584, 0, 0.1949, 0]								
doc2 =	[0, 0, 0.5849, 0, 1.584, 0, 1.584, 0.584, 0]								
doc3 =	[1.584, 0.5849, 0.5849, 1.584, 0, 0, 0, 0, 1.584]								

q = [0, 0.5849, 0, 1.5849, 0, 0, 0, 0]

$\text{len}(d_1) = 1.6038$

$\text{len}(d_2) = 2$

$\text{len}(d_3) = 2.236$

$\text{len}(q) = 1.6813$

$\text{cosSim}(d_1, q) = 0.042$

$\text{cosSim}(d_2, q) = 0$

$\text{cosSim}(d_3, q) = 0.7555$

Page 4 of 6

d_3, d_1, d_2

(iv): if stemming, "playing" and "play" to "play".

This would increase the similarity between Doc2. It results Doc2 more relevant in the ranking.

Question 3:

- (a) The link structure of some web pages is shown below. There are six web pages shown (d1, d2, d3, d4, d5 and d6), and the arrows show links between the pages (e.g. d4 contains links to d3 and d6).

① Initial start with 1

$$② R(d_1) = (1 - 0.8) + 0.8 \times 1 = 1$$

$$R(d_2) = (1 - 0.8) + 0.8 \times 0 = 0.2$$

$$R(d_3) = (1 - 0.8) + 0.8 \times \frac{1}{2} = 0.6$$

$$R(d_4) = (1 - 0.8) + 0.8 \times 0 = 0.2$$

$$R(d_5) = (1 - 0.8) + 0.8 \times 1 = 1$$

$$R(d_6) = (1 - 0.8) + 0.8 \times (1 + 1 + 0.5) = 2.2$$

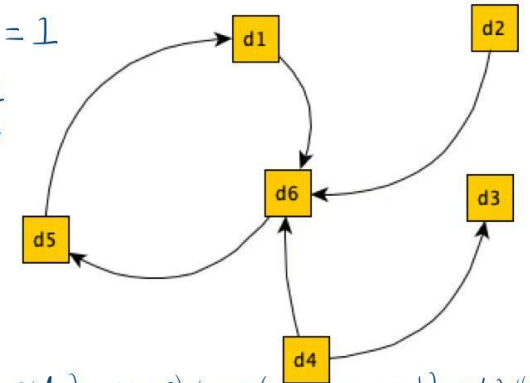
$$③ R(d_1) = (1 - 0.8) + 0.8 \times 1 = 1$$

$$R(d_2) = (1 - 0.8) + 0.8 \times 0 = 0.2$$

$$R(d_3) = (1 - 0.8) + 0.8 \times \left(\frac{2.2}{2}\right) = 0.28$$

$$R(d_4) = (1 - 0.8) + 0.8 \times 0 = 0.2$$

$$R(d_5) = (1 - 0.8) + 0.8 \times 1.22 = 1.96$$



$$R(d_6) = (1 - 0.8) + 0.8(1 + 0.2 + 0.1) = 1.24$$

Show a worked example of how PageRank scores are calculated for these documents.

Use a damping factor of 0.8 and show at least 3 iterations.

$$③ R(d_1) = (1 - 0.8) + 0.8 \times 1.96 = 1.768$$

$$R(d_2) = (1 - 0.8) + 0.8 \times 0 = 0.2$$

$$R(d_3) = (1 - 0.8) + 0.8 \times \left(\frac{2.2}{2}\right) = 0.28$$

$$R(d_4) = 0.2$$

$$R(d_5) = (1 - 0.8) + 0.8 \times 1.24 = 1.192$$

$$R(d_6) = (1 - 0.8) + 0.8 \times (1 + 0.2 + 0.2) = 1.24$$

[14 marks]

ts.

[s]

- (b) Compare the *MAP*, *bpref* and *NDCG* evaluation metrics. In your answer, outline any advantages or disadvantages of each. For each metric, suggest a situation where it is more appropriate than the others.

MAP: Advantage: rewards ranking relevant docs early. Disadvantage: Assumes complete relevance judgments. Best for small or medium-sized collections with complete relevance labels.

bpref: Advantage: ignores judged docs. Works well with incomplete relevance judgements; Disadvantage: ???

Best for: collections with incomplete relevance data.

NDCG: Adv: Suitable for real-world ranking tasks; emphasizes top-ranked relevant items. Disadvantage: Requires graded relevance judgments. Best for: If graded relevance judgments are available

[9 marks]

- (c) Below is a set of results and relevance judgments for a query:

Retrieved = d13, d21, d19, d12, d6, d24, d11, d1, d3, d17, d9, d23, d10, d14

Relevant = {d2, d3, d7, d9, d12, d15, d17, d23} $r=8$

$$MAP = \frac{P@4 + P@9 + P@10 + P@11 + P@12}{5}$$

Calculate the following metrics:

(i) Mean Average Precision (MAP)

(ii) Recall

(iii) R-Precision

$$R\text{-Precision} = \frac{|\text{Relevant in top } r \text{ result}|}{r}$$

$$= \frac{1}{8} = 0.125$$

$$= \frac{0.25 + 0.2 + 0.3 + 0.363 + 0.416}{5}$$

$$= 0.191$$

$$\text{Recall} = \frac{|\text{Rel} \cap \text{Ret}|}{|\text{Ret}|} = \frac{5}{8} = 0.625$$

[12 marks]

[Total 35 marks]

[12 marks]

[Total 35 marks]

Question 4:

- (a) The *Rocchio Algorithm* uses a *modified query vector* to achieve relevance feedback. Explain why this is effective in improving the effectiveness of an Information Retrieval system, and how this modified query vector can be calculated.

It moves the query closer to relevant document and away from non-relevant ones in vector space. [10 marks]

This helps better capture the user's true information need.

Why works: This method is especially helpful when recall is important, and users can identify relevant results

- (b) The table below shows results from three search engines in response to the same query. Each set of results consists of a ranked list of unique document identifiers (DocID), along with the ranking score. Complete the following tasks, showing your workings for each.
- Calculate the ranking score that document D6 would have using *CombSum*.
 - Calculate the ranking score that document D10 would have using *CombMNZ*.
 - Calculate the ranking score that document D5 would have using *Borda Fuse*.

Engine A		Engine B		Engine C	
DocID	Score	DocID	Score	DocID	Score
D10	0.60	D5	971	D12	9.23
D9	0.57	D1	936	D1	9.00
D12	0.48	D11	860	D2	7.88
D11	0.46	D2	516	D6	6.69
D8	0.41	D8	414	D7	5.03
D1	0.37	D6	300	D8	4.22
D6	0.26	D4	153	D5	3.63
D7	0.19	D10	99		

[9 marks]

- (c) Different levels of *corpus overlap* can influence the design of a data fusion algorithm. Explain why this is the case.

Disjoint databases(Collection Fusion): No document appears in more than one result set. Chorus effect is not applicable. Fusion must treat each result set independently. [10 marks]

Identical databases (Data Fusion): Documents will frequently appear in multiple result sets. Repeated appearances are strong signals of relevance(Chorus Effect). Fusion should give higher scores to such doc

Overlapping databases: Documents may appear in multiple result sets. It hard to judge if absence means irrelevance or unawareness.

- (d) Briefly describe *three* sources of synonyms for use in *query expansion*.

[6 marks]

- Manually-created thesaurus, such as WordNet
- Automatic-created thesaurus from some external corpus like Wikipedia
- Word embeddings, where words are represented by vctors, e.g. woerd2vec

[Total 35 marks]

$$C=12$$

$$\text{system A: } (4+3+2+1)/4 = 2.5$$

$$\text{system B give } D_s: 12$$

$$\text{system C give } D_s: 6$$

$$\text{Borda } D_s = 12+6+2.5 = 20.5$$

Q4. b) i)

d_b (Normalized)

$$\frac{\text{unNormal} - \min}{\max - \min}$$

$$\text{sys A } 0.171$$

$$\text{sys B } 0.231$$

$$\text{sys C: } 0.546$$

$$\text{CombSum}(d_b) = 0.171 + 0.231 + 0.546 = 0.948$$

q11 D_{10}

	system A	system B	system C	Fused
d_{10}	1	0	0	<u>1</u>