



Autumn Examinations 2021-2022

Course Instance Code(s)	4BCT, 1CSD1, 1CSD2, 4BS2
Exam(s)	4 th B.Sc. Computer Science and IT M.Sc. Computer Science (Data Analytics) M.Sc. Computer Science (Artificial Intelligence) B.Sc. (Hons)
Module Code(s)	CT4100
Module(s)	Information Retrieval
Paper No.	1
External Examiner(s)	Dr. Ramona Trestian, Dr. John Woodward
Internal Examiner(s)	Professor M. Madden *Dr. C. O’Riordan

Instructions: Answer any 3 questions. All questions carry equal marks

Duration	2 hours
No. of Pages	3
Discipline(s)	Computer Science
Course Co-ordinator(s)	Dr. Colm O’Riordan, Dr. Frank Glavin, Dr. Mattias Nickels

Requirements:

Release in Exam Venue	Yes
-----------------------	-----

CT4100 Information Retrieval

Question 1 (25 marks)

- (a) Recommender systems are used to generate recommendations for users on unseen items. Collaborative filtering is one such approach. Explain, in your own words, the main stages of collaborative filtering and illustrate how this approach can be used to generate a recommendation. (10)
- (b) Suggest approaches to evaluate the usefulness of a collaborative filtering system. Discuss any limitations of these approaches (8)
- (c) Outline a suitable approach to incorporating content information about items into a recommender system. (7)

Question 2 (25 marks)

- (a) Precision and recall are often used to evaluate the performance of an IR system. Explain these terms. Given a ranked answer set and evaluation judgements, explain how you would generate a precision-recall graph. (8)
- (b) Outline any limitations associated with using precision and recall and suggest alternative evaluation measures. (7)
- (c) Learning mechanisms has been used successfully in information retrieval. Using an approach of our choice, suggest a learning mechanism to identify good weights to apply to terms in documents. Discuss any limitations of your approach (10)

1/A 3 main steps in collaborative filtering:

① Find how similar each user to every other user-

ratings user	I ₁	I ₂	I ₃	I ₄
U ₁	0	10	5	0
U ₂	1	9	3	8
U ₃	0	0	1	2

→ U₁ & U₂ is similar

★ with cosine similarity / pearson correlation

② Form clusters of users who are similar



★ with correlation thresholding / vetting correlation.

③ In each cluster, make recommendations based on what other users in the group have rated which not yet rated by active user.

→ recommend I₄ to U₁ as U₂ has given high ratings to I₄.

★ with computed weighted average with user ratings.
compute weighted mean with neighborhood ratings.

1b) Evaluation metrics for CF.

- ① Coverage \rightarrow ability of system to provide a recommendation on a given item.
- ② accuracy \rightarrow correctness of the recommendations generated by the system.

Approaches to evaluate:

- MAE, RMSE to measure accuracy.
 - \swarrow error as percentage
 - \searrow contributions of the absolute errors between the prediction & the true value.

downside: The accuracy measurement might not be useful to users as they might ask for movies that are unfamiliar with but not old favorites they do not likely want again.

1c) for each of the item, extract key properties such as brand, type, functionality, genre and etc.

Then based on the ratings of each user we could suggest new items that have similar characteristics at high ratings to user (increase coverage).

for new user, we could suggest popular items based on the ratings of all users along with item characteristics matching

for users with no similar ratings as others, same method applies.

②*) precision \rightarrow the correctness in returned set $\frac{TP}{TP+FP}$ $\frac{|RA|}{|A|} \rightarrow$ answer

recall \rightarrow the returned percentage across entire relevant set.

$$\frac{TP}{TP+FN} \quad \frac{|RA|}{|R|} \rightarrow \text{relevant set.}$$

precision-recall graph.

\rightarrow should only be generated from ranked list.

Top 10 ranked list =

$\hookrightarrow \underline{d_1} \underline{d_2} d_3 d_4 \underline{d_5} \underline{d_6} d_7 \underline{d_8} d_9 \underline{d_{10}}$

underlined = relevant doc.

document set = 20

relevant in document = 10

Then we can calculate the precision-recall pairs one by one

1. considering 1st document only:

$$\text{precision} = \frac{1}{1}, \text{recall} = \frac{1}{10}$$

2. considering first 2 doc:

$$\text{precision} = \frac{2}{2}, \text{recall} = \frac{2}{10}$$

3. considering first 3 docs

$$p = \frac{2}{3}, \text{ recall} = \frac{2}{10}$$

4. considering first 4 docs.

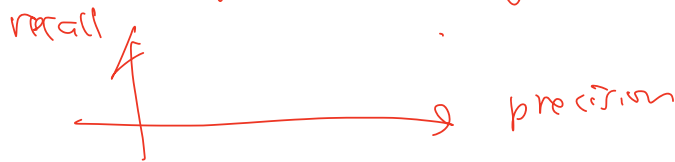
$$p = \frac{2}{4}, r = \frac{2}{10}$$

5. first 5:

$$p = \frac{3}{5}, r = \frac{3}{10}$$

then continue until the last answer.

to get 10 (x,y) points to plot.



2.1 - easily achieved recall 100% by returning everything

- precision-recall graph - impractical grid would be a large 1.74.

- might not be relevant as not measuring the true performance of system but only a single query

- alternative (1) average precision of n queries.
(2) then take the mean of the average precision for MAP as a single value measurement for

the system.

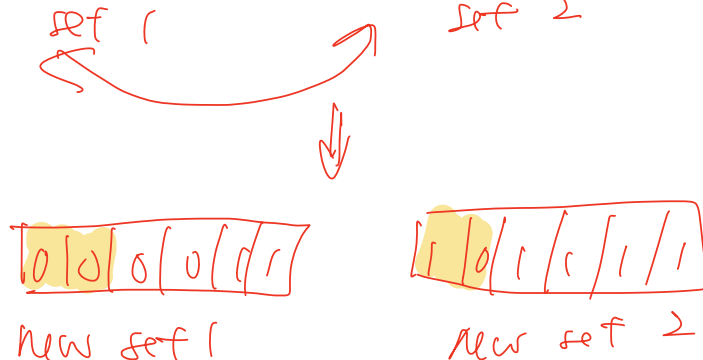
(3) use f -measure which combines precision & recall

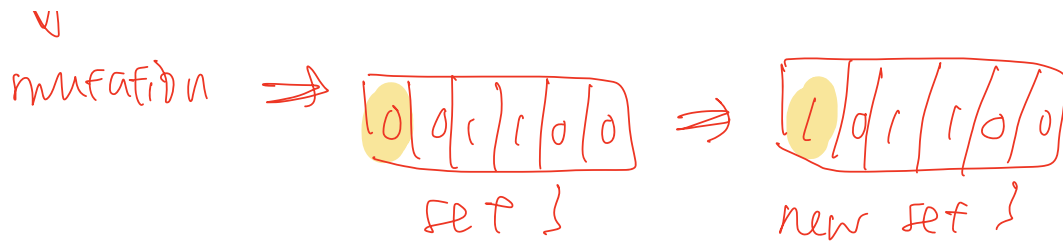
(2c) LR - genetic algorithm

steps :
 with length = size of lexicon
 + random rates

1. Create random population (genotype)
2. calculate fitness with evaluation function to result in map
3. select good population
4. perform cross over / mutation
5. repeat to step 2-

Crossover \Rightarrow





Limitation : (1) choice of fitness function

(2) choice of representation for encoding

(3) parameters setting : I - size of population

II - number of generations

III - probability of crossover

IV - probability of mutation

(4) definition of selection scheme

Question 3 (25 marks)

- (a) Describe, in your own words, with reference to any well-known term weighting scheme, the main constituents of a good weighting scheme. *Automatic* (9)
- (b) Local analysis of the returned answer set is often used to select a number of terms to add to the current query. Describe an approach that may be used to select candidate terms to add to the query. Discuss the advantages and potential limitations of your approach. *association / graph metric*
correlation context. distance $d(t_i, r)$ (8)
- (c) Suggest an approach to extend local analysis in order to ensure a diverse set of candidate terms. *Assignment \rightarrow clustering*
select N % from each cluster. (8)

Question 4 (25 marks)

- (a) Discuss approaches to identify the difficulty of a user provided query. Your answer should include features that are used to measure the level of difficulty and you should distinguish between pre-retrieval and post-retrieval approaches. (8)
- (b) Given a set of scientific articles (which contain title, abstract, authors, key words, main body of the paper and a bibliography), suggest a suitable approach to measure the similarity between these documents in the collection into useful sub clusters that may be of use in user search tasks. *vector space representation* (8)
- (c) Given the collection described in (b) and using the similarity measure(s), you have defined, describe an approach to clustering this collection in sub collections that may be of use to users engaging in search tasks. (9)

②A - A good weighting scheme should obey the constraints introduced in the axiomatic approach.

↳ for example the BM25 - which:

$$\sum_{t \in Q \cap D} \left(\frac{\text{tf}_t^D}{\text{tf}_t^D + K_1} \cdot \frac{\log \left(\frac{N - \text{df}_t + 0.5}{\text{df}_t + 0.5} \right)}{(1-b) + b \cdot \frac{\text{dl}}{\text{dl}_{\text{avg}}}} \cdot \frac{\text{tf}_t^Q}{\text{Query}} \right)$$

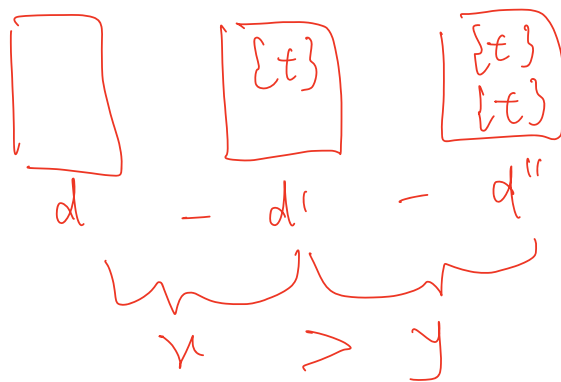
fulfill

1. adding query term increases score of document.
2. adding non-query term decreases score of document.
3. adding successive occurrence of term to a document increase less score than the addition of successive occurrence.

$$\text{sim}(d', q) > \text{sim}(d, q)$$

more than

$$\text{sim}(d'', q) > \text{sim}(d', q)$$



4. using a sublinear term-frequency factor component with k_r as controller. To ensure appearance of non-query term weighted less.

3b) - Association clustering is to find out the correlation of 2 terms by constructing a style

matrix $M :=$ and formula:

	t_1	t_2	t_3
t_1			
t_2			
t_3			

$$M_{ij} = \frac{\text{freq } t_i, j}{\text{freq } i + \text{freq } j - \text{freq } i, j}$$

then we could select Top N terms that correlates to t_{term_i} based on the matrix.

Pros: Cheaper computational cost.

Cons: Less evidence, no distance & no context.

3c) - clustering terms that selected from 3b.

then select N% from each / top 10 clusters.

4/11 - Pre-retrieval approaches :
estimate difficulty without running system.

① Linguistic approach:

I - To identify if query is ambiguous. (probabilistic parsing.)

② Statistical approach :

I - Identify the frequency distribution of the query terms in the collection (Idf/Isf).

II - Identify the non-specific terms in queries.

III - Identify term relatedness (mutual information / Jaccard coefficient)

IV ; Identify query scope -

↳ how many docs contain the query term -
(measure difference from collection language model & query language model)

Post retrieval approaches :

examine results from system of query

① Clarity: Compare language model induced from result set and 1 induced from collection set.

② robustness : measure by a deviation of :

I - query : generate a sub-query and measure overlap -

II, document : apply same query in modified version of collection

III = system: search for same query in same collection but different system

③ Scope analysis =

I = measure difficulty based on distribution of values in both results and collection sets

④ Measure similarity between articles:

Vector space model (cosine similarity)

↳ can use BM25 to generate weights for each article

④c) K-means clustering -

↳ use 4b to compute relatedness between each vector (euclidean distance)