

University of Virginia  
Department of Computer Science

**CS 4780/6501: Information Retrieval**  
**Spring 2020**

**Tuesday 3:30pm, March 31<sup>st</sup> to Thursday 3:30pm, April 2<sup>nd</sup>**

Name:
ComputingID:

- This is a **closed book** and **closed notes** exam. No electronic aids or cheat sheets or discussing the questions with anyone else are allowed.
- You are expected to finish this exam within 75 minutes.
- There are 7 pages, 4 parts of questions (the last part is for bonus questions), and 115 total points in this exam.
- Please carefully read the instructions and questions before you answer them.
- Please directly fill in your answers to the editable area below each question. And make sure to ***save*** your answers before you close the PDF file. (***NO*** argument about this after the grading is done).
- If you need any clarification of the exam questions, please directly send an email to the instructor within the exam period.
- Try to keep your answers as concise as possible; our grading is *NOT* by keyword matching.

Total	/100+15
-------	---------

## Academic Integrity Agreement

I, the undersigned, have neither witnessed nor received any external help while taking this exam. I understand that doing so (and not reporting) is a violation of the University's academic integrity policies, and may result in academic sanctions.

Signature: \_\_\_\_\_

**Your exam will not be graded unless the above agreement is signed.**

# 1 True/False Questions (12×3 pts)

Please choose either True or False for each of the following statements. For the statement you believe it is False, please give your brief explanation of it (you do not need to explain when you believe it is True). Three point for each question. *Note: the credit can only be granted if your explanation for the false statement is correct.*

1. Cosine similarity is proved to be equivalent to Euclidean distance in a vector space.  
*True/False, and Explain:*
2. The time complexity of searching for query  $q$  in an inverted index is  $O(|q| \times |L|)$ , where  $|q|$  is the length of query  $q$  and  $|L|$  is the average length of posting lists in the index.  
*True/False, and Explain:*
3. Browsing mode is preferred for document retrieval when the user has a clear mind of what he/she is looking for.  
*True/False, and Explain:*
4. Vector Space Model is equivalent to Bag-of-Word model.  
*True/False, and Explain:*
5. We can easily get the number of unique terms in a particular document from an inverted index.  
*True/False, and Explain:*
6. kappa statistic is lower bounded by -1.  
*True/False, and Explain:*
7. Stemming helps improve recall of a Boolean retrieval model.  
*True/False, and Explain:*
8.  $P(D, Q|R) = P(Q|R)P(D|Q, R)$  is the key step used in deriving the query generation model, a.k.a, the language model, where  $D$  denotes a candidate document,  $Q$  denotes a user query, and  $R$  is the relevance label.  
*True/False, and Explain:*
9. Relevance quality of a returned document is judged against all the query terms in the given query.  
*True/False, and Explain:*

10. Mean Average Precision prefers a system to return as many relevant documents as possible.  
*True/False, and Explain:*
11. We seldom use a database system to solve web search problems mainly because of the efficiency concerns.  
*True/False, and Explain:*
12. Maximum likelihood estimator is problematic when one does not have good knowledge about how to set up the prior.  
*True/False, and Explain:*

## 2 Multiple Choice Questions (6×4 pts pts)

Please choose all the answers that you believe are correct under each question.

1. Good “*basic concepts*” in a vector space model should be:
  - (a) orthogonal to each other;
  - (b) based on linguistic study;
  - (c) able to automatically compute the weights in each document;
  - (d) understandable by human.
2. Which of the following answer choices describes the correct order of the steps necessary to build an inverted index? Assume each document is represented by the tuple.
  - (a) sort by document ID, divide documents among different machines, sort by term ID, merge information about each term;
  - (b) sort by term ID, divide documents among different machines, sort by document ID, merge information about each term;
  - (c) divide documents among different machines, sort by document ID, sort by term ID, merge information about each term;
  - (d) sort by term ID, merge information about each term, divide documents among different machines, sort by document ID.
3. Zipf’s law tells us:
  - (a) head words take major portion in English vocabulary;
  - (b) in a given French corpus, if the most frequent word’s frequency is 1, then the second frequent word’s frequency is around 0.5;
  - (c) comparing to tail words, removing head words helps more to reduce the storage of documents represented by a vector space model when using a dense matrix data structure;
  - (d) sublinear TF scaling is necessary.

4. What are the key assumptions behind the derivation of the RSJ ranking model:
  - (a) binary relevance;
  - (b) uniform document prior;
  - (c) available relevance annotations;
  - (d) attributes not in the query are equally likely to occur in relevant and non-relevant documents.
5. Classical retrieval evaluation typically makes the following assumptions:
  - (a) independent result relevance evaluation;
  - (b) top down sequential browsing;
  - (c) query is a proxy of a user's information need;
  - (d) precision is preferred over recall.
6. In which of the following situations a ranking system's MRR performance would be equal to its MAP performance:
  - (a) the first returned document is relevant;
  - (b) all relevant documents are returned;
  - (c) there is only one relevant document;
  - (d) there is no relevant document.

### 3 Short Answer Questions (40 pts)

Most of the following questions can be answered by one or two sentences. Please make your answer concise and to the point.

1. Name five essential components of a search engine. (5pts)
  - 
  - 
  - 
  - 
  -
2. Describe the three stages of sorting-based inverted index construction (aka map-reduce). (3pts)
  - 
  - 
  -
3. Why are N-grams not an ideal way to account for phrases in a query? (3pts)

4. In class, we analyzed how Yahoo utilized a browsing method for its home page, while Google currently utilizes a querying method for its home page/search engine design. Provide two key differences between these two methods. (4pts)
- -
5. How does Zipf's law ensure effective inverted index compression? (4pts)
6. What are the three key elements of classical IR evaluation? (6pts)
- - 
  -
7. What is NDCG? How is it better than P@k, MAP and MRR? (3pts)
8. Suppose that we know there are 4 relevant files C, D, E, F. The result is returned in the following position: K, D, A, F, G, C. When calculating the mean average precision, after adding up the precision at each relevant position, what should we divide the sum by? (2pts)
9. Describe a circumstance in which one would want to use P@K over another ranking approximation like Mean Reciprocal Rank? And how about vice-versa (i.e. use MRR over P@K)? (4pts)
- -
10. What are the three key heuristics commonly shared by vector space models and BM25? (6pts)
- - 
  -

## 4 Bonus Questions (15 pts)

The first question is supposed to be open ended. Your answers have to be very specific to convince the instructor that you deserve the bonus (generally mention some broad concepts will not count).

1. Microsoft Bing would like to become a *real-time personalized* web search engine, i.e., its search results can be optimized with respect every action the user has taken in the system. Where do you think in the classical search engine architecture that Bing needs to renovate, and how exactly? (12pts)
2. List your favorite (labeled as +) and disliked (labeled as -) aspects of this class. (3pts)