# LBSC 796/INFM 718R: Information Retrieval Systems Spring, 2006

#### **Midterm Exam**

Name:

- Show sufficient work to demonstrate your understanding of the material. This is mostly for your benefit, as it will allow partial credit to be awarded.
- It might be a good idea to read over the entire exam first to get a sense of how to manage your time.
- You will have until 3:45pm to complete the exam.
- Good luck!

Question	Points	Total
1		16
2		8
3		28
4		12
5		18
6		18
Total		100

#### **Question 1. Evaluation (16 points)**

The following is the uninterpolated precision-recall graph of an IR system on a particular topic. You know that 20 hits were retrieved, and that there are 16 relevant documents for this topic (not all of which are retrieved).



**A.** (4 points) What does the interpolated graph look like? Draw neatly on the graph above.

**B.** (6 points) In the diagram below, each box represents a hit. Based on the above precision-recall graph, which hits are relevant? Write an "R" on the relevant hits. Leave the non-relevant hits alone.



Page 2 of 15

Page 1 of 15

**C.** (2 points) What is the Mean Average Precision (MAP)?

## **Question 2. TREC Collections (8 points)**

TREC-style evaluations with reusable test collections assume a model of the user and information-seeking behavior that is oversimplified and unrealistic. *Briefly* list two of these oversimplifications.

D. (2 points) What is the R-precision?

E. (2 points) What is Precision at 10?

Page 3 of 15

Page 4 of 15

#### Question 3. Boolean and Vector Space Retrieval (28 points)

Assume the following fragments comprise your document collection:

Doc 1: banking on banks to raise the interest rate Doc 2: jogging along the river bank to look at the sailboats Doc 3: jogging to the bank to look at the interest rate Doc 4: buzzer-beating shot banked in! Doc 5: scenic outlooks on the banks of the Potomac River

Assume that you drop stopwords. Assume that you stem.

**A.** (5 points) Construct the term-document matrix for the above documents that can be used in Boolean retrieval. The index terms have already been arranged for you alphabetically in the following table (assume terms not mentioned to be stopwords):

Term	Doc 1	Doc 2	Doc 3	Doc 4	Doc 5
bank					
buzzer-beating					
interest					
jog					
look					
outlook					
Potomac					
raise					
rate					
river					
sailboat					
scenic					
shot					

**B.** (2 points each) What documents would be returned in response to the following queries?

bank NOT interest

(interest AND rate) NOT jog

bank AND (scenic OR jog )

Page 5 of 15

Page 6 of 15

Doc 1: banking on banks to raise the interest rate Doc 2: jogging along the river bank to look at the sailboats Doc 3: jogging to the bank to look at the interest rate Doc 4: buzzer-beating shot banked in! Doc 5: scenic outlooks on the banks of the Potomac River

**C.** (10 points) Construct the vector space term-document matrix for the above documents (repeated from before) using *tf.idf* term weighting. Normalize your vectors. The following blank tables are provided for your convenience. You can use as many or as few of them as you wish. Clearly indicate your final answer.

Term	Doc 1	Doc 2	Doc 3	Doc 4	Doc 5
bank					
buzzer-beating					
interest					
jog					
look					
outlook					
Potomac					
raise					
rate					
river					
sailboat					
scenic					
shot					

Term	Doc 1	Doc 2	Doc 3	Doc 4	Doc 5
bank					
buzzer-beating					
interest					
jog					
look					
outlook					
Potomac					
raise					
rate					
river					
sailboat					
scenic					
shot					

Term	Doc 1	Doc 2	Doc 3	Doc 4	Doc 5
bank					
buzzer-beating					
interest					
jog					
look					
outlook					
Potomac					
raise					
rate					
river					
sailboat					
scenic					
shot					

Page 7 of 15

Page 8 of 15

**D.** (3 points each) Simulate the retrieval of documents in response to the following queries. Indicate the order in which documents will be retrieved, and the similarity score between the query and each document.

bank interest

#### Question 4. The dreaded disease... again (12 points)

In lecture 4, we talked about a hypothetical disease that's very rare, but there's a very accurate test for it. By Bayes' Rule, we computed the probability that you actually have the disease, given that you tested positive. This question explores this problem in more detail. Find the probability that you have the disease, given that you tested positive under the following two slightly different scenarios. The goal is for you to gain a better intuition of how probabilities interact.

**A.** (6 points) The incidence of the disease: 0.1% Accuracy of the test: 99%

bank river

**E.** (1 points) In two sentences or less, describe why indexing word senses does not yield higher retrieval performance.

Page 9 of 15

Page 10 of 15

**B.** (6 points) The incidence of the disease: 0.01% Accuracy of the test: 99.9%

### Question 5. The guessing game (18 points)

Assume we're playing a game in which I think of a number between 1 and 1000, and you have to guess that number in as few attempts as possible. With each guess, you can also ask me a yes/no question about the number and I'll answer truthfully.

To reiterate: each turn consists of you guessing a number, asking me a yes/no question, and receiving the answer. The goal is to guess the number I'm thinking of in as few turns as possible.

**A.** (6 points) Come up with a linear-time algorithm for guessing this number. By algorithm, I mean a strategy of what number you would guess first, guess next, etc., and what you would ask.

**B.** (6 points) Come up with a logarithmic-time algorithm for guessing this number. By algorithm, I mean a strategy of what number you would guess first, guess next, etc., and what you would ask.

Page 11 of 15

Page 12 of 15

**C.** (3 points) In the worst case, how many guesses would it take to figure out the number using the linear-time algorithm?

### Question 6. Language Processing and IR (18 points)

Three linguistic phenomena are listed below. For each,

- In ten words or less, describe what it is. (1 point)
- Give an original example, not presented in the lecture slides. (1 point)
- In two sentences or less, describe why it poses a problem for IR systems (2 points)
- Does it cause primarily precision or recall problems? (2 points)

### **Morphological Variation**

What is it?

**D.** (3 points) In the worst case, approximately how many times guesses would it take to figure out the number using the logarithmic-time algorithm?

Example:

Why is it a problem?

Precision or recall problem?

**Synonymy** What is it?

Example:

Why is it a problem?

Precision or recall problem?

Page 13 of 15

Page 14 of 15

Polysemy What is it?

Example:

Why is it a problem?

Precision or recall problem?

Page 15 of 15