# Exam IST 441 Spring 2014

- Last name: \_\_\_\_\_ First name: \_\_\_\_\_
- Student ID:
- I acknowledge and accept the University Policies and the Course Policies on Academic Integrity

This 100 point exam determines 30% of your grade.There are 45 point extra credit questions which graduate students have to answer but are optional for undergrads.

If you use the back of any paper, please note on the front of that page.

Name:

### Information Retrieval (15 pts)

- You are asked to design and implement an information retrieval / search engine system for enterprise search for at least 10,000,000 documents that are on an internal internet and is expected to grow.
  - (7pts) Explain to your clients what an enterprise search engine is by describing its basic components and what each does. What differences are there between enterprise search and a web search engine such as Google?

# Information Retrieval - Characteristics & Evaluation

- You are asked to design and implement an information retrieval / search engine system for enterprise search for at least 10,000,000 documents that are on an internal internet and is expected to grow.
  - (5pts) What are the important characteristics you must consider in the design including the size of storage needed? What is the storage for a mixture of documents and images?

(3pts) Define the measures that will help you evaluate its performance, including how you can design for scaling.

#### (5 nts)

Name:

Crawling (5 pts)

You are using a web crawler for your project.

• (5 pts) How does a crawler work? How does the crawler determine from a web site what and how much it can crawl? Give an example.

Crawling (10 pts)

For the following graph give the breath first and depth first paths starting from node A



Name:

Name:

# Document Analysis (35 pts)

• (5 pts) What is the document vector model, why is it important and how is it used by information retrieval and modern search engines? Describe what tokens are and how they are represented.

• Consider the following documents D and query Q for the following:

Stop words : {a, in, of}

- D1: A shipment of gold damaged in a fire.
- D2: Delivery of silver arrived in a silver truck.
- D3: A gold shipment arrived in a truck?
- Q1: gold silver truck
- (5 pts) Define for the above documents and queries the token vocabulary. Explain your tokenization.

Name:

#### **Document Vectors**

• (5 pts) Using the 3 documents and query, construct the document vector for each document and query with term or token frequency weights.

#### Inverted index

• (5 pts) Construct the dictionary file and posting file for the 3 documents

Name:

# Document Similarity and Ranking

• (5 pts) Define the inner product similarity metric between a document and a query.

• (5 pts) Construct the inner product for the term frequency document vectors for all documents with all queries. Make sure you show all the computation.

• (5 pts) Rank the documents for the query.

Name:

# Document Similarity and Ranking (extra credit – 10 pts)

• (4 pts - extra credit) Define the cosine similarity metric between documents and queries.

• (6 pts – extra credit) Calculate the cosine similarity between each document and other documents and query. Just show your work.

### Precision-Recall (20 pts)

• (5 pts) Define relevance; contrast it with importance.

• (5 pts) Define recall and precision in terms of relevant and irrelevant documents. Use set drawings to explain both.

• (3 pts) Which is more important for a search engine, recall or precision? Why?

Name:

#### **Precision-Recall Calculation**

- (7 pts) Consider the following universe of documents:
  - D1, D2, D3, D4, D5, D6, D7, D8, D9
  - For a particular query, documents D1, D3, D4, D8 are relevant. However our information retrieval system returns D3, D5, D6, D7, D8.
  - Calculate the Recall and Precision for this query. Be sure you show what documents belong in the ratio values for complete credit. Do not just put in a number.

Name:

#### Precision-Recall Calculation (extra credit – 5 pts)

• (5 pts) Your IR system works over a fixed set of documents. The number of returned documents is 100 and the number of those relevant is 10. Recall is 0.4 times greater than precision. Calculate precision and recall. Show your work.

Name:

# Precision-Recall (extra credit – 5 pts)

• Plot how in general precision and recall will change as the number of documents retrieved is increased. Make sure you give the numerical limits on precision and recall on your plots.

Name:

# Text processing (5 pts)

• What is Zipf's law and what does it mean?

Pagerank Calculation (10 pts & 5 pts extra credit)

Name:

- (10 pts) From ranked pages, calculate simple pagerank for the unranked pages, A, B, & C giving values for the links used. Also calculate the values on the outgoing links from pages A, B, & C. Assume no normalization.
- (5 pts extra credit) Normalize the outgoing values (just show your work).



Pagerank vs hubs and authorities (extra credit - 10 pts)

Name:

• (5 points extra credit) Give a definition of pagerank and what it means.

• (5 points extra credit) Define hubs and authorities; contrast with Pagerank.

Name:

## Queries and the Size of things (extra credit 5 pts)

You are trying to determine how many documents a search engine has indexed. You know the search engine is a full text indexer and has a complete Boolean query system.

1. (3 pts) What queries will give an estimate of the number of documents they have indexed? Which for the exact size? Why?

2. (2 pts) What is the size of Google? How do you determine it?

Name:

#### Complexity (extra credit 10 pts)

Give the Big O complexity for each term. Give the order of complexity and Label each as reasonable/unreasonable for scaling. The most complex requires the most work. All unspecified terms are undetermined constants

O(n)	Order of O(n)	Scaling

c.  $10 a^n + n^3$ 

1000

d. 100 k<sup>n</sup>

a.

b.

e.  $.06 n + \log n$ 

 $100 n^2 + 10^6 n$ 

f.  $20 + k \log n$