

## Assignment 4 (100 points)

Submit to TRACS

### 1. The following questions are related to chapter 1 and chapter 2.

1.1 (5 points) Draw the inverted index that would be built for the following document collection. (See Figure 1.3, Page 6 of the textbook, for an example.)

Doc 1 new home sales top forecasts

Doc 2 home sales rise in july

Doc 3 increase in home sales in july

Doc 4 july new home sales rise

1.2 (5 points) Consider the following fragment of a positional index with the format:

word: document: <position, position, . . .>; document: < position, . . .>

...

Gates: 1: <3>; 2: <6>; 3: <2,17>; 4: <1>;

IBM: 4: <3>; 7: <14>;

Microsoft: 1: <1>; 2: <1,21>; 3: <3>; 5: <16,22,51>;

The /k operator, word1 /k word2 finds occurrences of word1 within k words of word2 (on either side), where k is a positive integer argument. Thus k = 1 demands that word1 be adjacent to word2.

Describe the set of documents that satisfy the query Gates /2 Microsoft.

2. The following questions are related to chapter 6.

2.1 (10 points) Consider the table (left) of term frequencies for 3 documents denoted Doc1, Doc2, Doc3. Compute the tf-idf weights for the terms car, auto, insurance, best, for each document, using the idf values from the right table.

|           | Doc1 | Doc2 | Doc3 |
|-----------|------|------|------|
| car       | 27   | 4    | 24   |
| auto      | 3    | 33   | 0    |
| insurance | 0    | 33   | 29   |
| best      | 14   | 0    | 17   |

| term      | $df_t$ | $idf_t$ |
|-----------|--------|---------|
| car       | 18,165 | 1.65    |
| auto      | 6723   | 2.08    |
| insurance | 19,241 | 1.62    |
| best      | 25,235 | 1.5     |

2.2 (10 points) Compute the vector space similarity between the query “digital cameras” and the document “digital cameras and video cameras” by filling out the empty columns in the following table. Assume  $N = 10,000,000$ , logarithmic term weighting (wf columns) for query and document, idf weighting for the query only and cosine normalization for the document only. Treat **and** as a stop word. Enter term counts in the tf columns. What is the final similarity score?

| word    | query |    |         |                            | document |    |                              |                 |
|---------|-------|----|---------|----------------------------|----------|----|------------------------------|-----------------|
|         | tf    | wf | df      | $idf$ $q_i = wf \cdot idf$ | tf       | wf | $d_i = \text{normalized wf}$ | $q_i \cdot d_i$ |
| digital |       |    | 10,000  |                            |          |    |                              |                 |
| video   |       |    | 100,000 |                            |          |    |                              |                 |
| cameras |       |    | 50,000  |                            |          |    |                              |                 |

**3. The following questions are related to chapter 8.**

3.1 (10 points) An IR system returns 8 relevant documents, and 10 nonrelevant documents. There are a total of 20 relevant documents in the collection. What is the precision of the system for this search, what is its recall? what is the balanced  $F$  measure?

3.2 (10 points) Consider an information need for which there are 4 relevant documents in the collection. Compare two systems that run on this collection. Their top 10 results are judged for relevance as follows (the leftmost item is the top ranked search result):

System 1: R N R N N N N R R

System 2: N R N N R R R N N N

- a. What is the MAP of each system? Which has a higher MAP?
- b. What is the R-precision of each system? Does it rank the systems the same as MAP?

3.3 (20 points) The following list of R's and N's represents relevant (R) and nonrelevant (N) documents in a ranked list of 20 documents in response to a query from a collection of 10,000 documents. The leftmost item is the top ranked search result. This list shows 6 relevant documents. Assume that there are 8 relevant documents in total in the collection.

R R N N N N N R N R N N N R N N N R

- a. What is the precision of the system on the top 20?
- b. What is the  $F1$  (balanced  $F$  measure) on the top 20?
- c. What is the uninterpolated precision of the system at 25% recall?
- d. What is the interpolated precision at 33% recall?
- e. Assume that these 20 documents are the complete result set of the system. What is the MAP for the query?

Assume, now, instead, that the system returned the entire 10,000 documents in a ranked list, and these are the first 20 results returned.

- f. What is the largest possible MAP that this system could have?
- g. What is the smallest possible MAP that this system could have?

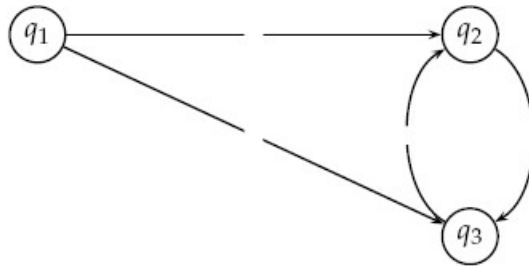
**4. The following questions are related to chapter 21.**

4.1 (15 points) Consider a web graph with three nodes 1, 2 and 3. The links are as follows: 1  $\rightarrow$  2, 3  $\rightarrow$  2, 2  $\rightarrow$  1, 2  $\rightarrow$  3. Write down the transition probability matrices for the random surfer's walk with teleporting, for the following three values of the teleport probability: (a)  $\alpha = 0$ ; (b)  $\alpha = 0.5$  and (c)  $\alpha = 1$ .

4.2 (15 points) For the web graph shown below, compute PageRank, hub and authority scores for each of the three pages. Also give the relative ordering of the 3 nodes for each of these scores.

PageRank: Assume that at each step of the PageRank random walk, we teleport to a random page with probability 0.1, with a uniform distribution over which particular page we teleport to.

Hubs/Authorities: Normalize the hub (authority) scores so that the maximum hub (authority) score is 1.



**Submission:** Type in Word (or write CLEARLY on paper and scan it) and submit electronically to TRACS.