# Information Organization Midterm Examination (November 28th, 2012)

**Notice**

- Write your name and student number to all the answer sheets.
- You can open printed materials (slides, notes, memos, etc).
- You can use a calculator. But other electronic devices such as laptop PC, tablet, mobile phone, electronic dictionary, are not allowed.
- You can write answers either in English or Japanese.

**Question 1.**

Suppose that the following document IDs are in a postings list.

$$22, 24, 31, 38, 101, 132$$

Show a compact representation of the above postings list, using bits as few as possible. How many bits are necessary?

**Question 2.**

In information retrieval, it is common to retrieve top-$k$ relevant documents.

1. Explain an algorithm that uses min-heap and finds top-$k$ documents from a collection of $N$ documents.
2. Illustrate how min-heap changes when the following relevance scores of $N = 5$ documents are processed in this order. Here, assume that $k = 3$, and a higher score should be ranked higher.

$$0.21, 0.33, 0.15, 0.82, 0.64$$

**Question 3.**

Suppose that there is a collection of $N = 100000$ documents, with document frequencies (df) of terms

orange, apple, melon, mango

are 10000, 50000, 1000, 100, respectively. There are two documents $d_1$ and $d_2$:

$d_1$:  orange melon mango
$d_2$:  apple melon melon

Now consider the following query $q$:

$q$:  melon mango

1. Compute the tf-idf score between query $q$ and document $d_1$. For tf-idf score, use ltn.lnc, namely logarithmic term frequency $(1 + \log(tf_{t,d}))$, idf $(\log N/df_t)$ document frequency, no normalization for the query, and logarithmic term frequency, no document frequency, and cosine normalization for the document. You can use $\log 2 = 0.30$ for approximation.
2. Compute the tf-idf score between query $q$ and document $d_2$, in the same manner as (1). Then answer which document is more relevant to $q$.
3. Compute the Jaccard coefficient (a) between $q$ and $d_1$, and (b) between $q$ and $d_2$. Which one this result is indicating as more relevant, (a) or (b)?
4. As a relevance measure, what are advantages of tf-idf score over Jaccard coefficient?

**Question 4.**

1. Suppose that System A and System B for information retrieval were evaluated on the same document collection, over query to find "ramen recipe". The results were as follows:

   ```
   System A:  Retrieved  80 documents, in which  20 documents were about
   ramen recipe.  All the documents about ramen recipe in the collection were retrieved.
   ```

   ```
   System B:  Retrieved 20 documents, in which 15 documents were about
   ramen recipe.
   ```

   Now compare qualities of System A and System B using appropriate measures.

2. In web search, potentially all the web pages in the Internet can be the document collection. Why recall is difficult to measure in this situation? What is a good measure to evaluate quality of retrieved results in web search?

| 科目/Subject | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|

Introduction to Information Retrieval

担当教員/Lecturer

氏名/Name : XU, Shan (ジョ サン)

| 4 | 4 | 1 | 2 | 1 | 5 | 4 | 0 | - | 7 |
|---|---|---|---|---|---|---|---|---|---|

Score **91**

20l2年 11月 28 日提出

---

Question l: There are maybe two ways possible to represente using gap: 22, 2, 7, 13, 67, 31

Variable byte code : 00010110, 00000010, 00000111, 00001011, 01011111, 00011111

Gamma codes : 01111001l0, 100, 11011, 11011, 1111101111, 1111101111   **48 bits**

It's obvious that using Gamma codes is better, the representation is shown below:

11110110, 100, 11011, 11011, 1111101111, 1111101111   **42 bits**.

42 bits are necessary   **20**

**47**
**44**

Question >: 1. Use a binary min heap. Takes O(N log k) operations to construct then read off k winners
in O(k logk) steps. A binary min heap is a binary tree in which each node's value is less than the values
of its children. To process a new document d' with score s' following this steps:   **41**

① Get current minimun hm of heap (O(1))

② If s' < hm skip to next document

③ if s' > hm heap-delete-root 10(log k))   **10**

④ Heap-add d/s' (O(log k))   **34**

2. As described in 1, firstly find current minimum $h_m = 0.21$.

then read $S' = 0.33 > 0.21$, So add to this tree ($0.21 \rightarrow 0.33$)

then read $S' = 0.15 <$ minimum $h_m = 0.21$, $S' < h_m$, So skip to next one.

then read $S' = 0.82$, minimum $h_m = 0.21$, $S' > h_m$, So add to this tree ($0.21 \rightarrow 0.82$)

then read $S' = 0.64$, minimum $h_m = 0.21$, $S' > h_m$, and turn to the children node of 0.33, as in the children node, $h_m = 0.33$, $S' > h_m$, So add to the children node of 0.33,

it is shown as

Question3: As the figure show:

| term | query | | | | document | | | |
|---|---|---|---|---|---|---|---|---|
| | tf | wt | idf | Wt,q | tf | wt | Wt,d | n'lized | product |
| orange | 0 | 0 | 1000 | 0 | 1 | 1 | 0.577 | 0 |
| melon | 1 | 1 | 1000 | 2 | 1 | 1 | 0.577 | 1.15 |
| mango | 1 | 1 | 100 | 3 | 1 | 1 | 0.577 | 1.73 |

So the tf-idf score is $0 + 1.15 + 1.73 = 2.88$

次ページへつづく

早稲田大学大学院情報生産システム研究科

前ページのつづき

2　term

| term | query tf | wf | df | idf | weight | document tf | wf | wtd | n'lized | product |
|------|----|----|-----|-----|--------|----|----|-----|---------|---------|
| apple | 0 | 0 | 50000 | 0.3 | 0 | 1 | 1 | 1 | 0.61 | 0 |
| melon | 1 | 1 | 1000 | 2 | 0 | 2 | 1.3 | 1.3 | 0.79 | 1.58 |
| mango | 1 | 1 | 100 | 3 | 3 | 0 | 0 | 0 | 0 | 0 |

So tf-idf score is 1.58

As 2.88 > 1.88, So document 1 is more relevan to q.

3. (a) Jaccard $(q, d_1)$ = $\frac{2}{3}$

(b) Jaccard $(q, d_2)$ = $\frac{1}{3}$

As Jaccard is an measure of overlap of two sets. So, (a) is more relevant.

4. Jaccard coefficient doesn't consider term frequency. And Rare terms are more informative than frequent terms. Jaccard doesn't consider this information. However, tf-idf score both take the term frequency and rare terms more informative into consideration. At the same time, it normalizes for the length of a document.

Question 4: 1. For the document collection, we can calculate precision, recall, and $F_1$ to compare.

| For System A | relevant | not relevant | for system B | relevant | not relevant |
|---|---|---|---|---|---|
| retrieved | 20 | 60 | retrieved | 15 | 5 | 20 |
| not retrieved | 0 | 20 | not retrieved | 5 | 20 |

For system A, precision $P_1 = \frac{20}{80} = \frac{1}{4}$, recall $R_1 = \frac{20}{20} = 1$, $F_1 = \frac{2}{5}$

For system B, precision $P_2 = \frac{15}{20} = \frac{3}{4}$, recall $R_2 = \frac{15}{20} = \frac{3}{4}$, $F_1 = \frac{3}{4}$

We can know that B has a higher precision, and lower recall than A, because it returns less results. and B has a better harmonic mean than A.

2. Recall (R) is the fraction of relevant documents that are retrieved.

$$Recall = \frac{\#(relevant\ items\ retrieved)}{\#(relevant\ items)}$$

Recall is a non-decreasing function of the number docs retrieved. And users care about the high-ranking documents, like the top 4, 15 of web search is better to use precision at top k or use measures that rewards user more for getting rank 1 right than for getting rank 10 right.

# Information Organization Final Examination (January 23th, 2013)

**Notice**

- Write your name and student number to all the answer sheets.

- You can bring printed materials (slides, notes, memos, etc). But NO electronic device, such as cell phone, laptop, electronic dictionary, or calculator is allowed.

- You can write answers either in English or Japanese.

**Question 1.** Alice wants to send a document $D$ to Bob securely over the Internet. But Bob wants that $D$ must be checked by Charlie. So Charlie receives $D$ from Alice first. After he checks $D$, he sends $D$ and his confirmation message to Bob. Bob wants that copies of $D$ received from Alice and Charlie are authentic (truly from them) and identical. Alice wants a message that Bob correctly received $D$.

Describe a sequence of communications over the Internet which uses public-key encryption and satisfies these requirements.

**Question 2.** Consider the following transactions T1, T2 and T3:

| | |
|---|---|
| T1: | R(A), R(B), W(A), W(C) |
| T2: | R(B), R(A), W(B) |
| T3: | R(C), W(C) |

Show schedules that demonstrate the following notions (1)-(4), by interlacing one or more transactions from T1, T2, and T3. You can add S(), X(), U(), commit, and abort, if necessary. Here, S(A) means a shared-lock on A, X(A) means an exclusive-lock on A, and U(A) means unlocking A. A schedule should be written like:

T3: S(A), T3:X(C), T3: R(C), T1: S(A), T1:R(A), T3:X(C), T3: W(C), T3:U(C), T3: U(A), T3: commit

Also you need to explain why each schedule satisfies its notion.

1. Read-write conflict.

2. Deadlock.

3. Schedule that is conflict-serializable, but causing cascading aborts.

4. Schedule that is view-serializable, but not conflict-serializable.
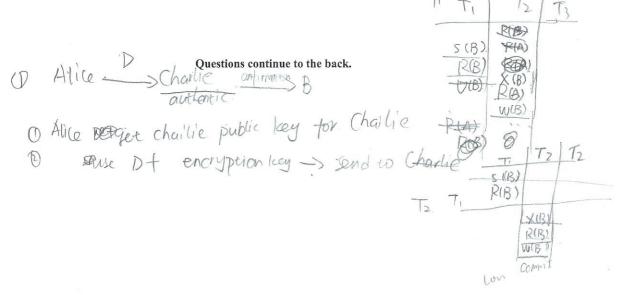
# Question 3.

Answer following questions.

1. In what situations SIX lock mode is used? Is IX lock compatible with SIX lock? Justify your answers.

2. Suppose that table emp(EID, salary) holds records of employee IDs (eid) and his/her salary (Salary). Now consider the following transaction T:

   Compute the average salary (AvgSalary) of all the employees. Then for each employee, if he/she has salary below AvgSalary, then increase his/her salary by (AvgSalary - Salary)*0.1.

In the questions below, justify your answers.

**(a)** Describe a lock schedule for T that guarantees conflict serializability.

**(b)** Describe a lock schedule for T that may not be conflict-serializable, may produce an inaccurate result, but has more concurrency than (a).

**Questions continue to the back.**

**Question 4.** Consider the following XML document $d_1$:

```
<cars>
    <car> <maker>Toyota</maker>
        <year>2012</year>
        <model>Voxy <special>Tourist</special></model>
        <option>Air bag</option>
        <color>blue</color>
    </car>
    <car> <maker>Honda</maker>
        <year>2008</year>
        <model grade=''G''>Insight</model>
        <color>red</color>
        <special>navi
            <option>radar</option>
        </special>
    </car>
</cars>
```

*car[ //model ][*

1. Write the answers to the following XPath queries applied on $d_1$.

   **(1)** /*/car[color][option]

   **(2)** //[color[text() ="blue"]/ancestor:://*[text()="radar"]]

   **(3)** //option/../../[//model][following:*]

2. Show a DTD $e$ such that the above document $d_1$ is valid against $e$.

3. Consider the following XPath queries (4) and (5).

   **(4)** //a/following::b//a

   **(5)** //a/ancestor::*/following::b/*/a

   (a) Show a document $d_2$ that returns an empty result to query (4) and returns a non-empty result to query (5).

   (b) Show a document $d_3$ that returns a non-empty result to query (4) and returns an empty result to query (5).

   (c) The product (6) of queries (4) and (5) is an XPath query such that i) if query (4) and query (5) return the same result $r$ on a document, then (6) also returns $r$, and ii) otherwise (6) returns an empty result. Show such an XPath query (6).

4. Concisely describe advantages and disadvantages of using XML.

Score **64**

Information Organization

氏名/Name  XU, Shan

**Question1:** Public-key Encryption: User's public encryption key is known to all, decryption key is only known by user.

① Bob generates a public encryption key and sends it to Charlie. Charlie issues a certificate to Bob.

② This certificate is stored in encrypted form, encrypted with Charlie's private key, known only to Charlie ...

③ Charlie's public key is know to all users, including A, which can encrypt the certificate and obtain Bob's public key.

first with session key, then ...

⑤ Alice encrypted the document with Bob's public key to Charlie, Charlie check it, and send the confirmation message to Bob.

④ Bob ask Alice to send copies of D.

⑥ Bob received D, and decrypted it with the decryption key.

⑦ Bob send a message to Alice to inform that he has received D.

**Question2:**

1. $T_1 = S(B)$, $T_2 = R(B)$, $T_2 = S(B)$, $T_2 = R(B)$, $T_2 = X(B)$, $T_2 = W(B)$, $T_2 = U(B)$, $T_1 = R(B)$, $T_1 = commit$.

It is obvious that $T_2$ can not repeatable read ...

*(red handwritten margin notes: 24, 40, 20, 64)*

2: T1: S(A), R(A),    S(B)
   T2:        X(B) W(B)
   T3:               S(C), R(C),    X(C)

There is a cycle in the dependency graph of the schedule.

3. T1: X(A), W(A), U(A)
   T2:              X(A), R(A), U(A),    X(A), W(A),    commit
   T3:                                            abort.

Because T2 read the written by T1, so if T1 abort, T2 will abort.

The result of the schedule is serializable to conflict equivalence of the committed transactions

4. T1:
   T2:

Question3: 1. If T1 like S and 2X at the same time, it will use SIX lock mode.

For example, T1 scan R, and updates a few tuples.

T1 get an SIX lock on R, then repeatedly get an S lock on tuples of R, and occasionally updates to X on tuples.

SIX lock is : compatible with SIX lock. As mentioned SIX is S and 2X, @ S and 2X is compatible, so SIX and 2X lock is compatible.

IX lock is : compatible with SIX lock. As mentioned SIX is S and 2X, @ S and 2X is compatible.

早稲田大学大学院情報生産システム研究科                    次ページへつづく