

**ComS 535x:
Homework #2**

Due on Mar 3, 2015

Instructor: Professor Pavankumar Aduri

Chenguang He

Contents

1	Question 1	3
2	Question 2	4
3	Question 3	5
4	Question 4	6
5	Question 5	7
6	Question 6	8
7	Question 7	9

Question 1

Consider the following documents **D1 = 1, 4, 6, 7, 8** and **D2 = 2, 3, 9, 4, 7**

(a) What are the binary term-frequency vectors of D1 and D2?

Answer:

$$T = (1, 4, 6, 7, 8, 2, 3, 9)$$

$$Tf_{11} = \langle 1, 1, 1, 1, 1, 0, 0, 0 \rangle, Tf_{12} = \langle 0, 1, 0, 1, 0, 1, 1, 1 \rangle$$

(b) What is the Jaccard Similarity of D1 and D2 (with respect to binary term-frequency vectors)

Answer:

$$Jac(Tf_{11}, Tf_{12}) = \frac{|Tf_{11} \cap Tf_{12}|}{|Tf_{11} \cup Tf_{12}|} = \frac{Tf_{11} \cdot Tf_{12}}{Tf_{11}^2 + Tf_{12}^2 - Tf_{11} \cdot Tf_{12}} = 0.25$$

(c) What is the cosine similarity of D1 and D2 (with respect to binary term-frequency vectors)

Answer:

$$Cos(Tf_{11}, Tf_{12}) = \frac{Tf_{11} \cdot Tf_{12}}{\|Tf_{11}\| \|Tf_{12}\|} = 0.4$$

Question 2

Let D_1 and D_2 be two documents. Let C be the cosine similarity of the documents with respect to binary term-frequency vectors and J be the jacquard similarity with respect to binary term-frequency vectors. Show that

(a) $C^2 \leq J$

Answer:

Let $\frac{C^2}{J} = \frac{|D_1 \cap D_2|^2}{|D_1||D_2|} \times \frac{|D_1 \cup D_2|}{|D_1 \cap D_2|} = \frac{|D_1 \cap D_2| \cdot |D_1 \cup D_2|}{|D_1||D_2|}$, when $D_1 = D_2$, $C^2 = J$, otherwise, $|D_1 \cap D_2|$ will go smaller and the inequality will less than 1.

(b) $J \leq \frac{C}{2-C}$

Answer:

Let $D_1 = X$ and $D_2 = Y$, and i from 1 to n .

Since, $J = \frac{\sum X_i Y_i}{\sum X_i^2 + \sum Y_i^2 - \sum X_i Y_i}$ and $C = \frac{\sum X_i Y_i}{\sqrt{\sum X_i^2} + \sqrt{\sum Y_i^2}}$

We have $\frac{C}{J} = \frac{\sum X_i^2 + \sum Y_i^2 - \sum X_i Y_i}{\sqrt{\sum X_i^2} + \sqrt{\sum Y_i^2}} = \sqrt{\frac{\sum X_i^2}{\sum Y_i^2}} + \sqrt{\frac{\sum Y_i^2}{\sum X_i^2}} - C$

Thus, $J = \frac{C}{\sqrt{\frac{\sum X_i^2}{\sum Y_i^2}} + \sqrt{\frac{\sum Y_i^2}{\sum X_i^2}} - C}$

Because of X and Y are two vectors, we can write the relationship between X and Y as: $|X| = c|Y|$ where $c > 0$. then

when $c = 1$, $|X| = |Y|$, we have $J = \frac{C}{2-C}$, we have $J \leq \frac{C}{2-C}$, since $\sqrt{\frac{\sum X_i^2}{\sum Y_i^2}} + \sqrt{\frac{\sum Y_i^2}{\sum X_i^2}}$ has minimum value,

when $X = Y$. Otherwise, it $J < \frac{C}{2-C}$ when $X \neq Y$

Question 3

3. Let D_1 and D_2 be two documents such that $|D_1 \cup D_2| = 1, \dots, n$. Show that the Jaccard similarity of D_1 and D_2 can be computed exactly in time $O(n \log n)$.

Answer:

There are two parts for Jaccard similarity, the union of two sets and the intersection of two sets. When we calculate the union of two sets, we can simply count all elements in two sets, it takes $O(n)$. When we calculate the intersection of two sets, we simply sort two sets, it takes $O(n \log n)$ times and for each element in set D_1 , we simply do binary search in D_2 , since D_2 is sorted. It takes $O(\log n)$ to find an element in D_2 , and there are n elements in D_1 . Therefore, the entire program takes $O(n \log n)$ time.

Question 4

Suppose we picked the following permutations (2x+1)the MinHash matrix.

$$(2x + 1) \% 5 = 0 \rightarrow 1, 1 \rightarrow 3, 2 \rightarrow 0, 3 \rightarrow 2, 4 \rightarrow 4 \text{ so } D1 = 0, D2 = 1, D3 = 2, D4 = 0$$

$$(3x + 4) \% 5 = 0 \rightarrow 4, 1 \rightarrow 2, 2 \rightarrow 0, 3 \rightarrow 3, 4 \rightarrow 1 \text{ so } D1 = 0, D2 = 2, D3 = 1, D4 = 0$$

$$(x + 3) \% 5 = 0 \rightarrow 3, 1 \rightarrow 4, 2 \rightarrow 0, 3 \rightarrow 1, 4 \rightarrow 2 \text{ so } D1 = 0, D2 = 3, D3 = 1, D4 = 0$$

0	1	2	0
0	2	1	0
0	3	1	0

Question 5

Suppose that we toss a biased coin (probability of head $1/4$) n times. Give a lower bound the probability that we see at least $\log n$ consecutive heads. In locality sensitive hashing, we showed that if two documents are s -similar, then the probability that they are mapped to the same bucket in some hash table is at least $1 - (1 - s^r)^b$. Do you see similarity between the two proofs?

Yes. Assume that, we toss the coin n times. For each time, we map it into $\lceil \frac{n}{\log n} \rceil$ blocks. Because the coin is biased, the probability of consecutive $\log n$ tosses is $(\frac{1}{4})^{\log n}$, and the probability of toss in each block of tail is at most $1 - (\frac{1}{4})^{\log n}$.

Therefore, in $\log n$ toss, the probability of each block get all head is $1 - (1 - (\frac{1}{4})^{\log n})^{\frac{n}{\log n}}$, where is similarity with $1 - (1 - s^r)^b$ for $s = \frac{1}{4}$, $r = \log n$, $b = \frac{n}{\log n}$

Question 6

Prove Claim 3 from Notes II

Define random variable X_i to be the number of i th permutation where MH_a and MH_b matches. So we have, $X_i = 1$ if $\min[\prod_i(D_a)] = \min[\prod_i(D_b)]$, $X_i = 0$. otherwise.

Since, for each X_i , it is independent event. We use chernoff's bound there:

$$Pr\left[\left|\frac{x}{k} - Jac(D_a, D_b)\right| \geq Jac(D_a, D_b) \cdot \delta\right] \leq 2 \cdot e^{-\frac{\delta^2 \cdot k \cdot Jac(D_a, D_b)}{2}}$$

$$\implies Pr\left[\left|\frac{x}{k} - Jac(D_a, D_b)\right| \leq Jac(D_a, D_b) \cdot \delta\right] \geq 1 - 2 \cdot e^{-\frac{\delta^2 \cdot k \cdot Jac(D_a, D_b)}{2}}$$

The output of algorithm A is $\frac{l}{k}$. Let $k = c \cdot \frac{1}{\varepsilon^2 \cdot \log \frac{1}{\delta}}$ and $\delta = \frac{\varepsilon}{Jac(D_a, D_b)}$, then we have:

$$Pr[|Output of A - Jac(D_a, D_b)| \leq \varepsilon] = 1 - 2e^{-\left(\frac{\varepsilon}{Jac(D_a, D_b)}\right)^2 \cdot \frac{c \cdot \frac{1}{\varepsilon} \cdot \log \frac{1}{\delta} \cdot Jac(D_a, D_b)}{2}} = 1 - 2 \cdot \delta^{\left(\frac{c}{2 \cdot Jac(D_a, D_b)}\right)}$$

Because when c increase to large enough, then $2 \cdot \delta^{\left(\frac{c}{2 \cdot Jac(D_a, D_b)}\right)} \leq \delta$

We add $(1 - \delta)$ at both side, then we have: $1 - 2 \cdot \delta^{\left(\frac{c}{2 \cdot Jac(D_a, D_b)}\right)} \geq 1 - \delta$

Therefor $Pr[|Output of A - Jac(D_a, D_b)| \leq \varepsilon] \geq 1 - \delta$

Question 7

Because the hash function is one to one mapping, by the Claim 1 in Note, we have $Let S = 1, \dots, m, T = 1, \dots, m+1$ and $Let h$ be a hash function. We want to prove that:

$$Pr[h(i) = j] = \frac{1}{m+1}$$

Proof:

$$Pr[h(i) = j] = \frac{m}{m+1} \times \frac{m-1}{m} \times \frac{m-2}{m-1} \dots \times \frac{m-i+1}{m-i+2} \times \frac{1}{m-i+1} = \frac{1}{m+1}$$

Now, by Claim 2, we want to prove that:

$$Pr[\min[h(D_a)] = \min[h(D_b)]] = Jac(D_a, D_b)$$

Proof:

$$Pr[\min[h(D_a)] = \min[h(D_b)]] = \frac{|D_a \cap D_b|}{m+1} = \frac{|D_a \cap D_b|}{|D_a \cup D_b|} = Jac(D_a, D_b)$$