

09/20/2024

## Big Data HW 1

Grothulath  
Thirumaran  
675086674

- 1) 1) Code has been uploaded on the gradescope.
- 2) After a spark session is initiated to read a text file, the extract\_friends function structures these relationships into direct & mutual connections. Which are then combined using the create\_map function. The get\_sort\_recommendation function aggregates and ranks the recommendation for each user.

3)

9550	9554, 9533, 9544, 9558, 153, 1220, 1421, 1436, 1951, 2413
8997	8998, 8987, 8992, 9001, 9003, 9009, 4849, 7174, 7279, 7364
4985	79, 577, 4839, 4984, 4986, 4987, 4988, 4989, 4990, 4991
4049	4871, 4875, 4889, 8492, 8685, 439, 660, 1100, 1137, 1156
2791	21185, 8783, 13280, 18359, 18363, 23667, 35740, 2204, 2786, 5996
3151	3161, 43162, 3230, 3450, 8692, 161, 2036, 3136, 3137, 3162
1724	1711, 1663, 1712, 1718, 1662, 1697, 1700, 1715, 1716, 1658
5060	5052, 5057, 5086, 14271, 98, 364, 575, 596, 611, 622
11	27552, 7785, 27573, 27574, 27589, 27590, 27600, 27617, 27620, 27667
8961	12241, 8973, 8965, 8963, 8966, 8967, 7174, 8969, 12243, 7177
739	732, 367, 381, 336, 21526, 28064, 677, 704, 728, 736

2) a)

$$\text{Conf}(A \rightarrow B) = \text{Pr}(B|A)$$

$$= \frac{\text{Pr}(A \cap B)}{\text{Pr}(A)}$$

Consider,

Basket 1  $\rightarrow$  A, B, E, F

Basket 2  $\rightarrow$  A, B, E

Basket 3  $\rightarrow$  A, B

Basket 4  $\rightarrow$  A, D

Basket 5  $\rightarrow$  B, C, E, F

Basket 6  $\rightarrow$  B, C, E

Basket 7  $\rightarrow$  B, C

Basket 8  $\rightarrow$  B, C, D

So,

$$\text{Pr}(A) = 4/8 = 0.5 \rightarrow \textcircled{1}$$

$$\text{Pr}(A \cap B) = 3/8 = 0.375 \rightarrow \textcircled{2}$$

using  $\textcircled{1}$  &  $\textcircled{2}$

$$\text{Conf}(A \rightarrow B) = \frac{\text{Pr}(A \cap B)}{\text{Pr}(A)} = \frac{0.375}{0.5} = \underline{\underline{0.75}}$$

The confidence might misrepresent the importance of an association, which is a major drawback. It concentrates on the popularity of A rather on B. If B is as popular as A, there is a high probability of basket containing A and also B, thus inflating confidence measure. From the data above

$$Pr(B) = 7/8 = 0.875$$

That infers that B appears in basket very often.

Lift measures how much "A + B" occur together than what if A + B are statistically independent.

If  $Lift > 1$ , Then B is more likely to be in the basket if A is there.

If  $Lift < 1$ , Then B is unlikely to be in the basket if A is there.

$$Lift(A \rightarrow B) = \frac{conf(A \rightarrow B)}{S(B)} = \frac{0.75}{0.875} < 1$$

Conviction compares the pr of A if B is not there and are independent, with actual freq of the appearance of A without B.

High conviction value  $\Rightarrow$  highly dependent on antecedent

Conviction = 1  $\Rightarrow$  items are independent.

$$\text{Conv}(A \rightarrow B) = \frac{1 - S(B)}{1 - \text{conv}(A \rightarrow B)} = \frac{1 - 0.875}{1 - 0.75} = 0.5$$

So, B is not highly depending on antecedent.

Hence lift and conviction do not suffer from the drawback.

b)

Basket 1  $\rightarrow$  A B C F

Basket 2  $\rightarrow$  A B E

Basket 3  $\rightarrow$  A B

Basket 4  $\rightarrow$  A D

Basket 5  $\rightarrow$  B C E F

Basket 6  $\rightarrow$  B C E

Basket 7  $\rightarrow$  B C

Basket 8  $\rightarrow$  B C D

From the table, confidence is not symmetric.

$$\text{Conf}(A \rightarrow B) = P_r(B|A)$$

$$= \frac{P_r(A \cap B)}{P_r(A)} = \frac{0.375}{0.5} = 0.75 \text{ (1)}$$



$$\begin{aligned}\text{conf}(B \rightarrow A) &= \frac{\text{Pr}(A \cap B)}{\text{Pr}(B)} \\ &= \frac{0.355}{0.875} = 0.428 - \textcircled{2}\end{aligned}$$

We can infer from ① + ② that the confidence is not symmetric

To prove:

Lift is symmetric

⇒

$$\begin{aligned}\text{Lift}(A \rightarrow B) &= \frac{\text{conf}(A \rightarrow B)}{S(B)} \\ &= \frac{\text{Pr}(B|A)}{\text{Pr}(B)} \\ &= \frac{\frac{\text{Pr}(A \cap B)}{\text{Pr}(A)}}{\text{Pr}(B)} \\ &= \frac{\frac{\text{Pr}(A \cap B)}{\text{Pr}(B)}}{\text{Pr}(A)} \\ &= \frac{\text{Pr}(A|B)}{\text{Pr}(A)} = \text{Lift}(B \rightarrow A).\end{aligned}$$

Thus proved.

(6)

To prove: Conviction is not symmetric.

⇒

$$\text{Conv}(A \rightarrow B) = \frac{1 - S(B)}{1 - \text{conf}(A \rightarrow B)}$$

$$= \frac{1 - 0.875}{1 - 0.75} = \frac{0.125}{0.25} = 0.5 \rightarrow \textcircled{1}$$

$$\text{Conv}(B \rightarrow A) = \frac{1 - S(A)}{1 - \text{conf}(B \rightarrow A)}$$

$$= \frac{1 - 0.5}{1 - 0.428} = \frac{0.5}{0.571} = 0.875 \rightarrow \textcircled{2}$$

$$\textcircled{1} \neq \textcircled{2}$$

Thus proved.

x

c)

Basket 1 A B

Basket 2 A B

Basket 3 C D

Basket 4 C F

$$\text{conf}(A \rightarrow B) = \text{Pr}(B|A)$$

$$= \frac{\text{Pr}(A \cap B)}{\text{Pr}(A)} = 1$$

$$\begin{aligned}\text{Conv}(A \rightarrow B) &= \frac{1 - \text{Pr}(B)}{1 - \text{conf}(A \rightarrow B)} \\ &= \frac{1 - 0.5}{1 - 1} = \infty\end{aligned}$$

~~for~~

$$\begin{aligned}\text{Lift}(A \rightarrow B) &= \frac{\text{conf}(A \rightarrow B)}{\text{Pr}(B)} \\ &= \frac{1}{0.5} = 2\end{aligned}$$

$$\begin{aligned}\text{Lift}(C \rightarrow D) &= \frac{\text{conf}(C \rightarrow D)}{\text{Pr}(D)} \\ &= \frac{1}{0.25} = 4\end{aligned}$$

Thus  $A \rightarrow B$ ,  $C \rightarrow D$  are 100% of the time, they have different lift. value.

d)

Number of frequent item sets: 553  
 DAI93865  $\rightarrow$  FR040251 = 1.0000000000  
 GR085051  $\rightarrow$  FR040251 = 0.9991762768  
 DAI88079  $\rightarrow$  FR040251 = 0.9867256637  
 FR092469  $\rightarrow$  FR040251 = 0.9835100118  
 DAI43868  $\rightarrow$  SNA82528 = 0.9729729730

e)

8

DAI23334, ELE92920  $\rightarrow$  DAI62779 = 1.0000000000  
 DAI55911, GR085051  $\rightarrow$  FR040251 = 1.0000000000  
 DAI75645, GR085051  $\rightarrow$  FR040251 = 1.0000000000  
 ELE17451, GR085051  $\rightarrow$  FR040251 = 1.0000000000  
 ELE20847, FR092469  $\rightarrow$  FR040251 = 1.0000000000

\_\_\_\_\_ x \_\_\_\_\_

3) a) no. of columns with  $m$  1's out of  $u$

$$\Rightarrow \binom{u}{m} = \frac{u!}{m! (u-m)!}$$

no. of col. having no 1's for selected  $k$  rows

$$\Rightarrow \binom{u-k}{m} = \frac{(u-k)!}{m! (u-k-m)!}$$

$$\begin{aligned} \Pr(\text{don't know}) &= \frac{\binom{u-k}{m}}{\binom{u}{m}} = \frac{(u-k)! \cdot \cancel{m!} (u-m)!}{u! \cdot \cancel{m!} (u-k-m)!} \\ &= \frac{(u-k)! (u-m)!}{u! (u-k-m)!} \end{aligned}$$

$\Pr(\text{don't know}) \leq \left(\frac{u-k}{u}\right)^m$ , so probability is  
 at most  $\left(\frac{u-k}{u}\right)^m$



(9)

3 b) To prove:

$$\left(\frac{n-k}{n}\right)^n \leq e^{-10}$$

Inference:-

$$\left(\frac{n-k}{n}\right)^n \leq e^{-10}$$

$$\Rightarrow \left(1 - \frac{k}{n}\right)^n \leq e^{-10}$$

$$\Rightarrow \left[\left(1 - \frac{k}{n}\right)^{n/k}\right]^{\frac{nk}{n}} \leq e^{-10}$$

$$\Rightarrow e^{-\frac{nk}{n}} \leq e^{-10}$$

$$\Rightarrow -\frac{nk}{n} \leq -10$$

$$\Rightarrow k \geq \frac{10n}{m}$$

So, the smallest  $k$  is  $\frac{10n}{m}$ 

c)

$\pi_1$	$\pi_2$	$\pi_3$	$s_1$	$s_2$	$s_1$	$s_2$
1	3	2	0	0	2	2
2	1	3	1	1	1	1
3	2	1	0	1	3	1
Cyclic Permutation			input matrix		Signature matrix	

Succesed Similarity between  $s_1, s_2$  is  $\frac{1}{2}$ ,But for that the minhash value agree =  $\frac{2}{3}$

4) a) Given,

(6)

$G = H^k$  is  $(\lambda, c\lambda, P_1^k, P_2^k)$  sensitive.

$\Rightarrow$

So each  $1 \leq j \leq L$  and  $x \in T$   
the  $P_x [x \in T \cap \omega_j] \leq P_2^k$   $\left[ \begin{array}{l} \text{as } k = \log 1/P_2^k \\ \text{when } P_2^k = 1/n \end{array} \right.$

Therefore  $P_x [x \in T \cap \omega_j] \leq 1/n$

$$\Rightarrow E[|x \in T \cap \omega_j|] \leq 1.$$

By  $E[\sum_{j=1}^L |T \cap \omega_j|] \leq L$ . (linearity of expectation)

By Markov's inequality,

$$P_x [x \geq a] \leq \frac{E[x]}{a}.$$

$$\text{So, } P_x \left[ \sum_{j=1}^L |T \cap \omega_j| \right] \leq \frac{E[\sum_{j=1}^L |T \cap \omega_j|]}{3L} < \frac{L}{3L} = \frac{1}{3}$$

b) Given

$G_j$  is  $(\lambda, c\lambda, p_1^k, p_2^k)$  sensitive for  $1 \leq j \leq L$

$$P_\sigma[g_i(x^*) = g_i(z)] \geq p_1^k \quad - (1)$$

By defn,

$$L = n^p \quad - (2)$$

$$p = \frac{\log 1/p_1}{\log 1/p_2} \quad - (3)$$

$$k = \log 1/p_2^u \quad - (4)$$

from (4)

$$k = \log 1/p_2^u \Rightarrow \log 1/p_2 = \frac{\log n}{k} \quad - (5)$$

from (3) subs in (5)

$$p = \frac{\log 1/p_1}{\log 1/p_2} = \frac{\cancel{\log 1/p_1}}{\log n} = \frac{\log 1/p_1}{\log n}$$

$$\Rightarrow p \log n = k \log 1/p_1$$

$$\Rightarrow n^p = (1/p_1)^k$$

(6) in (2)

$$L = n^p = \left(\frac{1}{p_1}\right)^k \Rightarrow p_1^k = \frac{1}{L} \quad - (7)$$

From ②

$$P_r [g_i(x^*) \neq g_i(z)] \leq 1 - p_i^t$$

$$= 1 - 1/2$$

$$\text{Thus } P_r [\forall 1 \leq j \leq L, g_j(x^*) \neq g_j(z)] \leq (1 - 1/4)^L \leq \frac{1}{e}$$

c)  $U$  be the set of  $(c, \lambda)$ -ANN,

$$\text{Thus, } U = \{x \in A : d(x, z) \leq c\lambda\}$$

①  $\Rightarrow$  None of  $(c, \lambda)$  ANN points are hashed into the same bucket with  $z$ , that means  $\forall j$   
 $1 \leq j \leq L, \omega_j \cap U = \emptyset$ .

$$P_r [E_1] \leq P_r [x^* \notin \bigcup_{j=1}^L \omega_j] = P_r [\forall 1 \leq j \leq L, g_j(x^*) \neq g_j(z)] \leq \frac{1}{e}$$

⑤  $\Rightarrow$  Atleast one  $(c, \lambda)$  ANN point is  $z^0$  be hashed,  
 but are more than  $3L$  points at distance greater than  $c\lambda$  in the union bucket.  
 If are less than  $3L$  points, the probability less than  $1/3$



$$P' = P_0[E_1 \cup E_2] \leq P_0[E_1] + P_0[E_2] \leq \frac{1}{e} + \frac{1}{3}, \quad (13)$$

$$\text{Thus } \underline{P > 1 - \frac{1}{e} - \frac{1}{3}}$$

d) LSH  $\Rightarrow$  0.406 second per Query

Linear Search took average 0.503 second per query

Error decreases and increases at 16 for L

Error increases for  $k \leq 16$ :

