

Q1 i) Given two sets A and B we hash them to get equivalent sets which contain numbers. Once we have these sets, we define 't' different random permutations T_1, T_2, \dots, T_t .

Now, we calculate the min wise signature for each of our set. To calculate the signature we simply check which number in the random permutation T_i matches first in our set.

Once we have the minwise signature for each set Jaccard similarity can be estimated by the following

Approximate Jaccard similarity = $\frac{\text{Number of Matching positions}}{\text{Total positions}}$

$$= \frac{|S \cap T|}{|S \cup T|}$$

i) Size of Min-Hash computing

Consider that we have 2 sets S_1 & S_2 .
We would like to estimate Jaccard Similarity.

Define indicator random variable X_i which is 1 if the i th min-hash match.

$$X = \sum_{i=1}^t X_i = r$$

For both S_1 & S_2 , the i th min-hash will be same.

Since all $S_1 \cup S_2$ elements are equally likely to come first in the permutation. Then chance that $|S_1 \cap S_2|$ are same is exactly $\frac{|S_1 \cap S_2|}{|S_1 \cup S_2|}$ which is

same as Jaccard similarity of S_1 & S_2 .

$$\Rightarrow \Pr[X_i=1] = \text{Jaccard similarity } (S_1, S_2)$$

$$\therefore E[X] = \sum_{i=1}^t E[X_i] = \sum_{i=1}^t \Pr[X_i=1]$$

$$= t \cdot \text{Jaccard Similarity } (S_1, S_2)$$

By the Chernoff bound

$$\Pr [X \in E[X](1 \pm \varepsilon)] \leq 2e^{-\frac{t \text{JaccardSim}(s_1, s_2) \varepsilon^2}{3}}$$

If we want that bad probability to be less than δ , we get

$$e^{-\frac{t^2}{3}} \leq \delta$$

(or) $t \geq \frac{\ln \frac{2}{\delta}}{\text{JaccardSimilarity}(s_1, s_2) \varepsilon^2}$

i.e. Suppose we want to measure Jaccard similarity to an accuracy $\pm \varepsilon$ only when similarity is at least $\frac{1}{2}$.

- Size of hash =

$$t \geq \frac{30 \ln 2}{\varepsilon^2}$$

ii) Similarity ≥ 0.8 document is selected

$$1 - 0.8 \leq 0.2 = K \text{ and now}$$

$$P_1 \geq 1 + K = 1 + 0.2 = 0.8$$

$$P_2 \leq 1 + e^{(a-b)} \quad \text{and now}$$

$$q = \frac{\log P_1}{\log P_2} = \frac{\log 1/0.8}{\log 1/p_2} = \frac{\log 10/8}{\log 10^2} = P_2 = \frac{8^2}{10^2}$$

$$P = 0.64$$

Similarity (at least) ≥ 0.8

$$\Rightarrow \text{and } c \leq C'K \leq 0.36$$

$$\Rightarrow \text{Similarity} \geq 0.64 \quad (c * 0.8)$$

=) Maximum possible value of c in our scheme

$$\Rightarrow c \leq 0.8$$

$$(0.8)^{0.8} \approx 0.73$$

2. Proof: algorithm frequent

for each $s \in \{1, 2, \dots, m\}$, counter on the list π atleast the number of occurrences of s in the stream minus $m/(k+1)$.

Suppose, in particular, if some s does not occur on the list, the counter = 0.

Also it says that it occurs fewer than $m/(k+1)$ times in the stream.

Consider each decrement counter step as eliminating some item.

An item is eliminated if it is currently being read and there are already k symbols different from it and k other items are simultaneously eliminated.

Thus, the elimination of each occurrence of an $s \in \{1, 2, \dots, m\}$ is elimination of $k+1$ items.

\therefore No more than $m/(k+1)$ occurrences of any symbol can be eliminated.

Item is not eliminated. (Ans)

\Rightarrow It will be on the list at the end

Any item that has frequency $> \frac{m}{k+1}$ times occur in the list.

Q3a The expected storage requirement of the algorithm will depend on the number of nodes selected and the count of their neighbours.

$$E[\text{storage requirement}] =$$

$$= E[\sum]$$

First let us define an indicator random variable

$$X_u = \begin{cases} 1 & \text{if node } u \text{ is sampled} \\ 0 & \text{otherwise} \end{cases}$$

$$E[\text{storage requirement}] = E\left[\sum_{u \in V_p} X_u + \sum_{u \in V_p} X_u |N(u)|\right]$$

Using linearity of expectation

$$= E\left[\sum_{u \in V_p} X_u\right] + E\left[\sum_{u \in V_p} X_u |N(u)|\right] 2M$$

$$\geq \sum E[X_u] + \sum E[X_u |N(u)|] 2M$$

$$= \sum P X_u + \sum P X_u |N(u)| 2M$$

$$= NP + 2NP |N(u)| M$$

$$= \frac{N}{\sqrt{M}} \cdot 60 + 2 \cdot \frac{N}{\sqrt{M}} \cdot \frac{60}{|N(u)|} M$$

Q 3 b: In order to compute the expected storage requirement as a function of M , consider

Lower Bound

For the lower bound, consider a graph which has minimum possible edges

$$\text{ie } M = N - 1$$

$$\Rightarrow N = M + 1$$

substituting this in our previous result we get

$$\text{Ans: } (M+1) \cdot \frac{60}{\sqrt{M}} + (M+1) \cdot \frac{60}{\sqrt{M}} \cdot \text{Ans}(\alpha) \cdot 2M$$

Upper Bound

For the upper bound, consider a fully connected graph

$$\text{ie } M = N_C = \frac{N!}{2! \times (N-2)!}$$

$$= \frac{N(N-1)}{2}$$

$$2M = N^2 - N$$

$$N^2 - N - 2M = 0$$

$$N = 1 \pm \sqrt{\frac{1+8M}{2}}$$

substituting in the previous result

$$\text{Ans: } \frac{1 \pm \sqrt{1+8M}}{2} \cdot \frac{60}{\sqrt{M}} \left(1 + \frac{2M}{\text{Ans}(\alpha)} \right)$$

Q3c

i) Define an indicator random variable

$$X_u = \begin{cases} 1 & \text{if node } u \text{ is sampled} \\ 0 & \text{otherwise} \end{cases}$$

$$E[\hat{D}] = E\left[\frac{1}{P} \cdot D_p\right] = E\left[\frac{1}{P} \cdot \frac{1}{|V|} \sum_{u \in V} \deg G_p(u)\right]$$

$$= E\left[\frac{1}{P} \cdot \frac{1}{|V|} \sum_{u \in V} X_u \deg G(u)\right]$$

$$= \frac{1}{P} \cdot \frac{1}{|V|} \sum_{u \in V} E[X_u \deg G(u)]$$

$$= \frac{1}{P} \cdot \frac{1}{|V|} P \cdot \deg G(u)$$

$$= \frac{1}{|V|} \deg G(u)$$

$$E[\hat{D}] = D$$

QED

ii) $\text{Var}[\hat{D}] = ?$

$$\text{Var}[\hat{D}] = E[\hat{D}^2] - (E[\hat{D}])^2$$

$$= E\left[\left(\frac{1}{P} \cdot \frac{1}{|V|} \sum_{u \in V} \deg G_p(u)\right)^2\right] - (E[\hat{D}])^2$$

Let $X_u = 1$ if node u is sampled
0 otherwise

$$\text{Var}[\delta] = \frac{1}{p^2 |V|^2} \mathbb{E} \left[\left(\sum_{j=0}^{|V|} X_{uj} \deg_G(u_j) \right)^2 \right] - D^2$$

$$= \frac{1}{p^2 |V|^2} \mathbb{E} \left[\sum_{j=0}^{|V|} (X_{uj} \deg_G(u_j))^2 + \sum_{i \neq j}^{|V|} X_{ui} \deg_G(u_i) \cdot X_{uj} \deg_G(u_j) \right]$$

$$\therefore (\sum a_i)^2 = \sum a_i^2 + \sum_{i \neq j} a_i a_j$$

$$= \frac{1}{p^2 |V|^2} \sum_{j=0}^{|V|} \mathbb{E} \left[(X_{uj} \deg_G(u_j))^2 \right] + \sum_{i \neq j}^{|V|} \mathbb{E} \left[X_{ui} \deg_G(u_i) \cdot X_{uj} \deg_G(u_j) \right]$$

since each vertex is sampled independently

$$= \frac{1}{p^2 |V|^2} NP \cdot (\deg_G(u_j))^2 \sum_{i=0}^{|V|} \sum_{\substack{j=0 \\ i \neq j}}^{|V|} \mathbb{E} \left[X_{ui} \deg_G(u_i) \cdot X_{uj} \deg_G(u_j) \right]$$

$$= \frac{1}{p^2 |V|^2} NP (\deg_G(u_j))^2 \cdot NP \deg_G(u_i) \cdot NP \deg_G(u_j)$$

$$= \frac{1}{p^2 |V|^2} N^3 P^3 (\deg_G(u))^4$$

$$= NP (\deg_G(u))^4 \quad (\because |V| = N)$$

Q 3c
(iii)

$$\text{Prob} \left(|D - \bar{D}| \geq \frac{D}{2} \right)$$

$$= \text{Prob} \left(|D - E(D)| \geq \frac{E(D)}{2} \right)$$

$$\leq \frac{\text{Var}(D)}{(E(D))^2}$$

$$\leq \frac{4 \cdot \text{Var}(D)}{D^2}$$

$$\leq \frac{4 \cdot N \cdot P \cdot (\deg G(u))}{D^2}$$

$$\leq \frac{4 \cdot N \cdot P \cdot |E|^2}{D^2}$$

$$\leq \frac{4 \cdot N \cdot P \cdot |E|^2 \cdot N^2}{|E|} = 4PN^3M$$

$$\text{Prob} \left(|D - \bar{D}| \geq \frac{D}{2} \right) \leq 4PN^3M$$

3 d) Let G be the entire graph which is going to be sampled.

Given the \hat{T} - number of distinct pair of vertices that are reachable via path of length $= 2$ from G_p

G_p - sampled graph with probability p .

Given that $E[\hat{T}] = T$ proving that \hat{T} is no. of distinct pairs of vertices that are reachable via length $= 2$

For every edge $e \in G(G_1, G_2)$ define an indicator random variable

$$X_e = \begin{cases} 1 & \text{if } e \text{ is sampled} \\ 0 & \text{otherwise} \end{cases}$$

$$\Rightarrow X = C_{G_p}(S, \bar{S})$$

$$= \sum X_e$$

$$e \in G(S, \bar{S})$$

$$E[X] = \sum_{e \in G(S, \bar{S})} E[X_e] = \sum_{e \in G(S, \bar{S})} E[X_e]$$

$$= P \cdot C_G(S, \bar{S})$$

$$E\left[\frac{X}{P}\right] = C_G(S, \bar{S})$$

The cut (sampling) value of the sampled graph is same as the original graph.

using Chernoff bound \Rightarrow

$$P\left[\left|\frac{X}{P} - C_G(S, \bar{S})\right| > \epsilon \cdot C_G(S, \bar{S})\right]$$

for length between 1/2 edges of 1 and atmost 2 hop.

~~atmost 2 hop~~ 2 hop

$$\overline{T} = \left| \sum_{\substack{u, v \\ u \in S, v \in \bar{S}}} (u, v) \right| = \sum_{u \in S} \sum_{v \in \bar{S}} (u, v)$$

$$P\left[\left|\frac{\overline{T}}{P} - C_G(S, \bar{S})\right| > \epsilon \cdot C_G(S, \bar{S})\right]$$

$$= P\left[|X - E[X]| > t \cdot \sigma[X]\right]$$

$$\leq e^{-\frac{t^2 \sigma^2}{2}}$$

$$\leq e^{-\frac{P \epsilon^2}{2}}$$

. - probability value \Rightarrow we need to reduce the deviation.

$$P \in \frac{2d \log n}{\epsilon^2} \Rightarrow$$

$$P \left[\frac{|X - T|}{P} \cdot \epsilon \cdot (1 + \epsilon) \cdot c \right] > \left(1 - \frac{1}{n^d}\right)$$

. - Hence $E[T'] = T$ from the above statement.

. - Hence I justified that sampling at 2 hops is G_2 is equal to sampling G_p .