

a). Ex 6.10 from FOS.

Let n be the number of flips needed to see the 1st head.

$$P(n=1) = P, P(n=2) = (1-P)P, P(n=3) = (1-P)^2P \dots$$

Expectation:

$$E(n) = \sum_{k=1}^{\infty} k \cdot P(n=k)$$

$$= \frac{P}{1-P} \sum_{k=1}^{\infty} k \cdot (1-P)^{k-1} P$$

$$= \frac{P}{1-P} \sum_{k=1}^{\infty} k \cdot (1-P)^k, \text{ let } (1-P)=r$$

$$= \frac{P}{1-P} r \sum_{k=0}^{\infty} k \cdot \cancel{(1-P)} r^{k-1}, \text{ we have } \frac{d}{dr} r^k = k \cdot r^{k-1}$$

$$= \frac{P}{1-P} r \frac{d}{dr} \left(\sum_{k=0}^{\infty} r^k \right), \text{ since } |r| < 1$$

$$= \frac{P}{1-P} r \frac{d}{dr} \left(\frac{1}{1-r} \right)$$

$$= \frac{P}{1-P} r \frac{1}{(1-r)^2}$$

$$= \frac{P}{1-P} \cdot (1-P) \cdot \frac{1}{P^2}$$

$$= \frac{1}{P}$$

b) Ex. 2.3.1 MNDS

(a) The largest number.

Map: ~~Map~~ map all numbers to same key, $\langle \text{max}, n_i \rangle$, n_i is the number

Reduce: then we have only one key-value pair: $\langle \text{max}, [n_1, n_2, \dots, n_i] \rangle$, reduce it by selecting the maximum of value list. The output is $\langle \text{max}, n_{\text{max}} \rangle$, n_{max} is the max for current machine.

Map: map all the outputs of last stage to themselves (identity).

Reduce: we have $\langle \text{max}, [n_{\text{max}1}, n_{\text{max}2}, \dots] \rangle$, where $n_{\text{max}i}$ represents the maximum of different distributed machine, reduce it by selecting the maximum of value list. The output is

$\langle \text{max}, n_{\text{MAX}} \rangle$, n_{MAX} is the largest number.

(b) The mean of data set.

Map: same as last question: $\langle \text{mean}, n_i \rangle$

Reduce: compute the mean of value list: $\langle \text{mean}, \frac{n_1 + n_2 + \dots + n_i}{i} \rangle$

Map: Identity

Reduce: $\langle \text{mean}, [m_1, m_2, \dots, m_k] \rangle \rightarrow \langle \text{mean}, \frac{m_1 + m_2 + \dots + m_k}{k} \rangle$

(c) The set of integer.

Map: map each number n_i to a key-value pair ~~$\langle \text{set}, 1 \rangle$~~ $\langle \text{set}, n_i \rangle$

Reduce: reduce $\langle \text{set}, [n_1, n_2, \dots] \rangle$ to $\langle \text{set}, \text{set}([n_1, n_2, \dots]) \rangle$

Map: Identity

Reduce:

~~(d) The count of distinct numbers.~~ $\langle \text{set}, \text{set}([set_1, set_2, \dots]) \rangle \rightarrow \langle \text{set}, \text{set}(\text{set}([set_1, set_2, \dots])) \rangle$

~~Map: map each number to $\langle \text{set}, 1 \rangle$~~

~~Reduce: $\langle \text{number}, [1, 1, \dots] \rangle \rightarrow \langle \text{number}, \# \rangle$~~

(d) The count of distinct number.

~~We just need to take the result of (c), problem~~

Map: map each number n_i to $\langle n_i, 1 \rangle$

Reduce: remove duplicated 1's for each key, so that each distinct number will have a key-value pair $\langle \text{num}, 1 \rangle$

Map: identity

Reduce: remove duplicated 1's for each key.

finally, we count the number of keys, that's exactly the same number of distinct number in data set.

(c). we have $R_{ij} = \begin{cases} +1, & \text{w.p. } \frac{1}{2} \\ -1, & \text{w.p. } \frac{1}{2} \end{cases}$

$$\text{so } E[R_{ij}^2] = \frac{1}{2}(1) + \frac{1}{2}(1) = 1$$

Want to show that $E[\|v\|^2] = \|u\|^2$

$$\text{Proof: we have } E[\|v\|^2] = E\left[\sum_{i=1}^k v_i^2\right]$$

$$= E\left[\sum_{i=1}^k \left(\frac{1}{k} \sum_{j=1}^d R_{ij} u_j\right)^2\right]$$

$$= E\left[\frac{1}{k} \sum_{i=1}^k \left(\sum_{j=1}^d R_{ij} u_j\right)^2\right]$$

$$= \frac{1}{k} \sum_{i=1}^k E\left[\left(\sum_{j=1}^d R_{ij} u_j\right)^2\right]$$

$$= \frac{1}{k} \sum_{i=1}^k \sum_{j=1}^d \left(E[R_{ij}^2]\ u_j^2\right) \quad \begin{matrix} \text{assume independent} \\ R_{ij} \end{matrix}$$

$$= \frac{1}{k} \sum_{i=1}^k \sum_{j=1}^d u_j^2$$

$$= \frac{1}{k} \cdot k \cdot \sum_{j=1}^d u_j^2$$

$$= \sum_{j=1}^d u_j^2$$

$$= \|u\|^2 \blacksquare$$

Part D

(i)

The probability that an edge of graph $G = (V, E)$ is present is $p_e = \frac{M}{C_n^2} = \frac{2M}{N(N-1)}$.

Each vertex $v \in V$ could potentially connect $(N-1)$ other vertex, pair-wisely. Then the expected number of neighbors of a vertex of the graph is $(N-1)\frac{2M}{N(N-1)} = \frac{2M}{N}$.

Since each vertex is sampled identically and independently, their sum of expectation equals to their expectation of sum.

The expected number of vertex being sampled is $E(|V_p|) = N\frac{60}{\sqrt{M}}$, and the expected storage requirement is $E(|V_p| + \sum_{u \in V_p} |N(u)|) = N\frac{60}{\sqrt{M}} + N\frac{60}{\sqrt{M}}\frac{2M}{N} = \frac{60N}{\sqrt{M}} + 120\sqrt{M}$.

(ii)

The probability that an edge of graph $G = (V, E)$ is present is between 0 and 1. i.e.

$$0 \leq \frac{2M}{N(N-1)} \leq 1. \text{ So, } 2M \leq N(N-1) \leq N^2. \sqrt{2M} \leq N.$$

$$\text{Thus, } \frac{60N}{\sqrt{M}} + 120\sqrt{M} \geq \frac{60\sqrt{2M}}{\sqrt{M}} + 120\sqrt{M} = 60\sqrt{2} + 120\sqrt{M} \text{ (lower bound).}$$

A connected graph with N vertex, has at least $(N-1)$ edges. It has $(N-1)$ when all the vertices are connected one by one and there is no circle in that graph. Reducing edge from the graph above will make the graph no longer connected. So, $M \geq N-1$. $M+1 \geq N$.

$$\text{Thus, } \frac{60N}{\sqrt{M}} + 120\sqrt{M} \leq \frac{60(M+1)}{\sqrt{M}} + 120\sqrt{M} = \frac{60}{\sqrt{M}} + 180\sqrt{M} \text{ (upper bound).}$$

$$\text{The upper bound is } \frac{60}{\sqrt{M}} + 180\sqrt{M} \text{ and the lower bound is } 60\sqrt{2} + 120\sqrt{M}.$$

(iii)

(A)

Assume that the number of neighbors of a vertex $v \in V$ of the graph $G = (V, E)$ is k . As calculated above, the expected number of neighbors of a vertex of the graph $G = (V, E)$ is

$$E(k) = \frac{2M}{N}.$$

Since a vertex $u \in V_p$ in graph $G_p = (V, E_p)$ is sampled from graph $G = (V, E)$, and for every $u \in V_p$, its entire neighbors are stored, the expected number of neighbors of a vertex $u \in V_p$ of

$$\text{the graph } G_p = (V, E_p) \text{ is also } E[\deg_{G_p}(u)] = (N-1)\frac{2M}{N(N-1)} = \frac{2M}{N}.$$

$$E[\hat{D}] = E\left[\frac{1}{p} \frac{1}{|V|} \sum_{u \in V_p} \deg_{G_p}(u)\right] = \frac{1}{p} \frac{1}{|V|} E\left[\sum_{u \in V_p} \deg_{G_p}(u)\right] = \frac{1}{p} \frac{1}{|V|} \sum_{u \in V_p} E[\deg_{G_p}(u)] = \frac{1}{p} \frac{1}{|V|} \frac{60N}{\sqrt{M}} \frac{2M}{N} = \frac{2M}{N}$$

(B)

For each vertex $u \in V_p$ in graph $G_p = (V, E)$, there are potentially $(N - 1)$ neighbors and the probability that a neighbor is connected with the vertex is $p_e = \frac{M}{C_n^2} = \frac{2M}{N(N - 1)}$. Since presence of edge between the vertex and a potential neighbor is independent and the degree of $u \in V_p$ in graph $G_p = (V, E_p)$ is sum of the presence of edges (1 indicates presence; 0 indicates absence), the distribution of degree of u , $\deg_{G_p}(u)$, is binomial.

$$\text{Prob}[\deg_{G_p}(u) = k] = C_{N-1}^k p_e^k (1 - p_e)^{N-1-k}$$

Then the variance of degree of u is

$$\begin{aligned} \text{Var}[\deg_{G_p}(u)] &= (N - 1)p_e(1 - p_e) = (N - 1)\frac{2M}{N(N - 1)}[1 - \frac{2M}{N(N - 1)}] \\ \text{Var}[\hat{D}] &= \text{Var}\left[\frac{1}{p} \frac{1}{|V|} \sum_{u \in V_p} \deg_{G_p}(u)\right] = \frac{1}{p^2} \frac{1}{|V|^2} \sum_{u \in V_p} \text{Var}[\deg_{G_p}(u)] \\ &= \frac{1}{p^2} \frac{1}{N^2} \frac{60N}{\sqrt{M}} (N - 1) \frac{2M}{N(N - 1)} [1 - \frac{2M}{N(N - 1)}] = \frac{M\sqrt{M}}{30N^2} [1 - \frac{2M}{N(N - 1)}] \end{aligned}$$

(C)

$$\begin{aligned} D &= \frac{1}{|V|} \sum_{v \in V} \deg_G(v); E(D) = \frac{1}{N} E\left[\sum_{v \in V} \deg_G(v)\right] = \frac{1}{N} \sum_{v \in V} E[\deg_G(v)] = \frac{1}{N} N \frac{2M}{N} = \frac{2M}{N}. \\ \hat{D} &= \frac{1}{p|V|} \sum_{u \in V_p} \deg_{G_p}(u); E[\hat{D}] = \frac{2M}{N}. \end{aligned}$$

Thus, $E[D] = E[\hat{D}]$, $E\left[\sum_{v \in V} \deg_G(v)\right] = \frac{1}{p} E\left[\sum_{u \in V_p} \deg_{G_p}(u)\right]$.

$\text{Prob}(|D - \hat{D}| \geq \frac{D}{2}) = \text{Prob}\left[|\frac{1}{|V|} \sum_{v \in V} \deg_G(v) - \frac{1}{p|V|} \sum_{u \in V_p} \deg_{G_p}(u)| \geq \frac{1}{2} \frac{1}{|V|} \sum_{v \in V} \deg_G(v)\right]$ can

be rewritten as $\text{Prob}(|S - \mu| \geq \frac{S}{2})$, where $S = \sum_{v \in V} \deg_G(v)$ and the mean of S is

$$\mu = E\left[\sum_{v \in V} \deg_G(v)\right] = \frac{1}{p} E\left[\sum_{u \in V_p} \deg_{G_p}(u)\right] = \frac{1}{p} \sum_{u \in V_p} \deg_{G_p}(u) = 2M.$$

($\frac{1}{p} E\left[\sum_{u \in V_p} \deg_{G_p}(u)\right] = \frac{1}{p} \sum_{u \in V_p} \deg_{G_p}(u)$ because it is estimation)

From Chernoff Bound,

$$Prob(S - \mu \geq \frac{S}{2}) = Prob(S \geq 2\mu) \leq e^{-\mu/3} = e^{-2M/3}$$

$$Prob(S - \mu \leq -\frac{S}{2}) = Prob(S \leq \frac{2\mu}{3}) \leq e^{-\mu(\frac{1}{3})^2/2} = e^{-M/9}$$

So,

$$Prob(|D - \hat{D}| \geq \frac{D}{2}) = Prob(|S - \mu| \geq \frac{S}{2}) = Prob(S - \mu \geq \frac{S}{2}) + Prob(S - \mu \leq -\frac{S}{2})$$

$$Prob(|D - \hat{D}| \geq \frac{D}{2}) \leq e^{-2M/3} + e^{-M/9}$$

(iv)

As described above, the probability an edge is present is $p_e = \frac{M}{C_n^2} = \frac{2M}{N(N-1)}$.

For $G = (V, E)$ or $G_p = (V, E_p)$, the expected number of neighbors of a vertex $v \in V$ is

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

For G , there are C_N^3 triples of vertices, each triple has probability P_e^3 of being a triangle. The

number of triangles is $C_N^3 P_e^3 = \frac{N(N-1)(N-2)}{6} \left[\frac{2M}{N(N-1)} \right]^3 = \frac{4M^3(N-2)}{3N^2(N-1)^2}$. The number of

distinct pairs of vertices that are reachable via paths of length 2 equals to the number of distinct triangles. Thus, $T = \frac{4M^3(N-2)}{3N^2(N-1)^2} \approx \frac{d^3}{6}$.

For G , assume that N_p vertices remain in G_p . Then remain edge probability for G_p is $\frac{d}{N_p}$. For G_p ,

there are $C_{N_p}^3$ triples of vertices, each triple has probability $(\frac{d}{N_p})^3$ of being a triangle. Let Δ_{ijk} be

the indicator variable for the triangle with vertices i, j, and k being present. The expected value of a sum of random variables is the sum of the expected values. Thus, the expected number of

triangles in G_p is $E\left(\sum_{ijk} \Delta_{ijk}\right) = \sum_{ijk} E(\Delta_{ijk}) = C_{N_p}^3 \left(\frac{d}{N_p}\right)^3 \approx \frac{d^3}{6}$. Then, $E[\hat{T}] = E\left(\sum_{ijk} \Delta_{ijk}\right) \approx \frac{d^3}{6}$.

Taken together, $E[\hat{T}] = T$.