

## SUBLINEAR TIME AVERAGE DEGREE

Given  $G = (V, E)$ , where  $|V| = n$ , and connected can we estimate the average degree in  $G$  by inspecting less than all of the vertices?

**Theorem.** *Let  $d$  be the average degree in  $G$  and  $d_S$  be the average degree of a set  $S$  of  $s$  vertices chosen uniformly<sup>1</sup> from  $V$ . Then if  $s \geq \beta\sqrt{n}/\varepsilon^{O(1)}$ ,  $\varepsilon^{O(1)}$  simply means we can choose any constant that makes the proof workout, for some  $\beta$  to be fixed later*

$$(0.1) \quad P\left(d_S \geq \left(\frac{1}{2} - \varepsilon\right) \cdot d\right) \geq 1 - \frac{\varepsilon}{64}$$

For the moment assume the theorem and note that by Markov's inequality it's immediate that

$$(0.2) \quad P(d_S \leq (1 + \varepsilon) \cdot d) \geq 1 - \frac{1}{1 + \varepsilon} \geq \frac{\varepsilon}{2}$$

Why? Markov's inequality says

$$P(X \geq a) \leq \frac{E(X)}{a}$$

and by complement

$$1 - P(X < a) \leq \frac{E(X)}{a} \iff P(X < a) \geq 1 - \frac{E(X)}{a}$$

Then  $X = d_S$  and  $E(X) = d$  and  $a = (1 + \varepsilon) \cdot d$  gives

$$P(d_S < (1 + \varepsilon) \cdot d) \geq 1 - \frac{d}{(1 + \varepsilon) \cdot d} = 1 - \frac{1}{1 + \varepsilon} = \frac{\varepsilon}{1 + \varepsilon} \geq \frac{\varepsilon}{2}$$

What's the point? We can exploit these two inequalities to construct a  $(2 + \varepsilon)$ -approximation algorithm<sup>2</sup>: pick  $k = \frac{8}{\varepsilon}$  sets  $S_i$  of size  $s$  and output the set with smallest average degree. Using eqn. 0.2 (why eqn. 0.2 and not more obviously eqn. 0.1? To get the direction of the inequality to work) the probability that all of the sets  $S_i$  overestimate average degree is

$$\begin{aligned} (P(d_S \geq (1 + \varepsilon) \cdot d))^k &= (1 - P(d_S < (1 + \varepsilon) \cdot d))^k \\ &= \left(1 - \frac{\varepsilon}{2}\right)^k = \left(1 - \frac{\varepsilon}{2}\right)^{8/\varepsilon} \\ &= \left(\left(1 - \frac{1}{2(\frac{1}{\varepsilon})}\right)^{2(\frac{1}{\varepsilon})}\right)^4 \\ &\leq (e^{-1})^4 \\ &\leq \frac{1}{8} \end{aligned}$$

The probability that any one of the  $k$  sets underestimates the average degree is

$$\begin{aligned} k \cdot P\left(d_S < \left(\frac{1}{2} - \varepsilon\right) \cdot d\right) &= k \cdot \left(1 - P\left(d_S \geq \left(\frac{1}{2} - \varepsilon\right) \cdot d\right)\right) \\ &\leq k \cdot \left(1 - \left(1 - \frac{\varepsilon}{64}\right)\right) \\ &\leq k \cdot \frac{\varepsilon}{64} = \frac{8}{\varepsilon} \cdot \frac{\varepsilon}{64} = \frac{1}{8} \end{aligned}$$

---

<sup>1</sup> $d_S$  is hence a random variable.

<sup>2</sup> $d_S \leq d \leq d_S(2 + \varepsilon)$ .

Hence with probability  $1 - \frac{1}{8} - \frac{1}{8} = \frac{3}{4}$  both eqn. 0.1 and eqn. 0.2 will be satisfied. Finally  $\varepsilon \rightarrow \varepsilon/2$  yields a  $(2 + \varepsilon)$ -approximation algorithm.

*Proof.* Let  $H$  be the set of vertices with highest degree in  $G$  and  $L = V \setminus H$ . Fix  $\varepsilon$  such that  $|H| = \sqrt{\varepsilon n}$ . Let  $v_d$  be the degree of a vertex and firstly

$$(0.3) \quad \sum_{v \in L} v_d \geq \left(\frac{1}{2} - \varepsilon\right) \sum_{v \in V} v_d$$

i.e.  $\sum_{v \in L} v_d$  is at least some non-vanishing proportion of total degree of  $G$ . Why is this the case? Edges incident on a vertex in  $L$  could contribute 2 to the sum: if  $\{u, v\} \subset L$  then  $(u, v)$  is counted in  $u_d$  and  $v_d$ . But if either  $u$  or  $v$  is in  $H$  then  $(u, v)$  will be counted in either  $u_d$  or  $v_d$  but not both. This isn't so bad because it means each  $(u, v)$ , as counted by  $\sum_{v \in L} v_d$ , is at least half as much as counted by  $\sum_{v \in V} v_d$ . If that were the whole story you'd have immediately

$$\sum_{v \in L} v_d \geq \frac{1}{2} \sum_{v \in V} v_d$$

Unfortunately edges  $(u, v)$  for which  $\{u, v\} \subset H$  might be a supermajority of the edges in  $E$  (thereby making  $\sum_{v \in L} v_d$  small again). But this can't happen because of how we chose  $\varepsilon$ : since  $|H| = \sqrt{\varepsilon n}$  there can only be  $(\sqrt{\varepsilon n})^2 = \varepsilon n$  such edges<sup>3</sup>, which is small relative  $|E|^4$ . So

$$\sum_{v \in L} v_d \geq \frac{1}{2} \sum_{v \in V} v_d - 2\varepsilon n \geq \frac{1}{2} \sum_{v \in V} v_d - \varepsilon \sum_{v \in V} v_d \geq \left(\frac{1}{2} - \varepsilon\right) \sum_{v \in V} v_d$$

Full disclosure: I'm not certain here. Okay don't forget the goal here is to produce a lower bound on the average degree of the vertices that end up in our sample. Therefore without loss of generality we can assume all vertices are sampled<sup>5</sup> from  $L$ . Let  $S$  be that sample (with  $|S| = s$ ) and let  $X_i$  be degree of the  $i$ th vertex<sup>6</sup> in  $S$ . With  $d_H$  being the vertex of lowest degree in  $H$  we have that for all  $v \in L$  it's the case that  $1 \leq v_d \leq d_H$  (since by definition all vertices in  $H$  have higher degree than all vertices in  $L$ ) and hence  $1 \leq X_i \leq d_H$ . Then (*supposedly*) by Hoeffding's inequality

$$(0.4) \quad \begin{aligned} P\left(\sum_{i=1}^s X_i \leq (1 - \varepsilon) E\left(\sum_{i=1}^s X_i\right)\right) &= P\left(\sum_{i=1}^s X_i - E\left(\sum_{i=1}^s X_i\right) \leq -\varepsilon E\left(\sum_{i=1}^s X_i\right)\right) \\ &\leq \exp\left(-\frac{2\varepsilon^2 (E(\sum_{i=1}^s X_i))^2}{\sum_{i=1}^s (d_H - 1)^2}\right) \\ &\dots \\ &\exp\left(-\frac{\varepsilon^2 (E(\sum_{i=1}^s X_i))^2}{d_H}\right) \end{aligned}$$

I have no idea how they got the right hand side. For Bernoulli random  $X_i$  you use the Anghuin/Valiant<sup>7</sup> inequality to get

$$P\left(\sum_{i=1}^s X_i \leq (1 - \varepsilon) E\left(\sum_{i=1}^s X_i\right)\right) \leq \exp\left(-\frac{\varepsilon^2 E(\sum_{i=1}^s X_i)}{2}\right)$$

<sup>3</sup>If  $H$  is completely connected, i.e. a clique, then every vertex in  $H$  is connected to every other, i.e.  $\sqrt{\varepsilon n}$  edges for everyone of the  $\sqrt{\varepsilon n}$  vertices.

<sup>4</sup>At least  $n - 1$  since the graph is connected (a minimum spanning tree on  $n$  vertices has  $n - 1$  edges)

<sup>5</sup>Since all vertices in  $H$  have higher degree, any lower bound produced by sampling  $L$  will also be a lower bound for average degree of vertices in  $H$ .

<sup>6</sup>A random variable since  $S$  is a sample.

<sup>7</sup>D. Angluin and L.G. Valiant, "Fast probabilistic algorithms for Hamiltonian circuits and matchings", *Journal of Computer and System Sciences*, No. 19, 1979, pp. 155–193.

But that's still not correct because it's missing  $d_H$ . Anyway moving on. This has to be interpreted now in terms of  $d_S$  and  $d$ . Firstly  $d_S = \sum_{i=1}^s X_i$  and  $d = E(d_S) = E(\sum_{i=1}^s X_i)$ . So

$$P(d_S \leq (1 - \varepsilon) d) \leq 1 - \exp\left(-\frac{\varepsilon^2 E(\sum_{i=1}^s X_i)}{d_H}\right)$$

Since any vertex in  $H$  has degree at least  $d_H$  the average degree has to be at least<sup>8</sup>  $d_H \cdot \left(\frac{|H|}{n}\right)$ . Hence by eqn. 0.3 we have that

$$E(X_i) \approx \frac{1}{|S|} \sum_{v \in S} v_d \geq \frac{1}{|L|} \sum_{v \in L} v_d \geq \frac{1}{|L|} \left(\frac{1}{2} - \varepsilon\right) \sum_{v \in V} v_d \geq \frac{1}{|V|} \left(\frac{1}{2} - \varepsilon\right) \sum_{v \in V} v_d \geq \left(\frac{1}{2} - \varepsilon\right) d_H \left(\frac{|H|}{n}\right)$$

The  $\approx$  comes from the law of large numbers<sup>9</sup> I think and the first inequality is because of ???. By linearity of expectation we have that

$$E(\sum_{i=1}^s X_i) \geq s \left(\frac{1}{2} - \varepsilon\right) d_H \left(\frac{|H|}{n}\right)$$

Given our choice of  $s \geq \beta\sqrt{n}/\varepsilon^{O(1)}$

$$\begin{aligned} E(\sum_{i=1}^s X_i) &\geq s \left(\frac{1}{2} - \varepsilon\right) d_H \left(\frac{|H|}{n}\right) \\ &\geq \left(\frac{1}{2} - \varepsilon\right) d_H \left(\frac{|H|}{n}\right) \frac{\beta\sqrt{n}}{\varepsilon^{O(1)}} \\ &= \left(\frac{1}{2} - \varepsilon\right) d_H \left(\frac{|H|}{\sqrt{n}\varepsilon^{O(1)}}\right) \beta \\ &= \left(\frac{1}{2} - \varepsilon\right) d_H \left(\frac{|H|}{\sqrt{n}\varepsilon\varepsilon^{O(1)}}\right) \beta \\ &= \left(\frac{1}{2} - \varepsilon\right) \left(\frac{d_H\beta}{\varepsilon^{O(1)}}\right) \end{aligned}$$

Note between lines 4 and 5 in the above the constant in  $\varepsilon^{O(1)}$  changes. Therefore

$$\begin{aligned} P(d_S \leq (1 - \varepsilon) d) &\leq \exp\left(-\frac{\varepsilon^2 E(\sum_{i=1}^s X_i)}{d_H}\right) \\ &\leq \exp\left(-\frac{\varepsilon^2 \left(\frac{1}{2} - \varepsilon\right) \left(\frac{d_H\beta}{\varepsilon^{O(1)}}\right)}{d_H}\right) \\ &\leq \exp\left(-\beta \left(\frac{1}{2} - \varepsilon\right)\right) \end{aligned}$$

Choose  $\beta$  such that

$$\exp\left(-\beta \left(\frac{1}{2} - \varepsilon\right)\right) \leq \frac{\varepsilon}{64}$$

and taking the compliment we get

$$P(d_S \geq (1 - \varepsilon) d) \geq 1 - \frac{\varepsilon}{64}$$

Finally since  $B \Leftarrow A$  implies  $P(B) \geq P(A)$  we have that

$$P\left(d_S \geq \left(\frac{1}{2} - \varepsilon\right) d\right) \geq P(d_S \geq (1 - \varepsilon) d) \geq 1 - \frac{\varepsilon}{64}$$

□

<sup>8</sup> $d = \frac{1}{n} \sum_{i=1}^n v_{d,i} = \frac{1}{n} \sum_{v \in L} v_d + \frac{1}{n} \sum_{v \in H} v_d \geq \frac{1}{n} \sum_{v \in L} v_d + \frac{1}{n} |H| d_H$

<sup>9</sup>The law of large numbers: for all  $\epsilon$  and random variables  $X_i$   $\lim_{n \rightarrow \infty} (|\bar{X}_n - E(\bar{X}_n)| > \epsilon) = 0$