

Lecture 3:

Sublinear time algorithms II

Course: Algorithms for Big Data

Instructor: Hossein Jowhari

Department of Computer Science and Statistics
Faculty of Mathematics
K. N. Toosi University of Technology

Spring 2021

Outline

- ▶ Deviation bounds: Markov, Chebyshev, Chernoff
- ▶ Estimating the average degree (continued)

Recap from previous lecture

We have a sequence of n (unknown) integers between 1 and $n - 1$ (these are degrees of a graph on n nodes.)

$$d_1, d_2, d_3, d_4, \dots, d_n$$

Want to estimate $d = \frac{d_1 + d_2 + \dots + d_n}{n}$

We sample s integers (with replacement) and output the average.

Let X_i be a random variable associated with the i -th sample.

$$\text{Algorithm's output: } X = \frac{1}{s}(X_1 + \dots + X_s)$$

$$\text{We know: } E[X] = d$$

Deviation from Expectation

We want to know how often X deviates from $E[X]$ by a considerable degree.

In other words, we want to bound this probability ($\epsilon \geq 0$)

$$\Pr\left(\underbrace{|X - E[X]|}_{\text{the amount of deviation}} \geq \epsilon E[X] \right)$$

We have some useful inequalities for this.

Deviation Bounds

Markov Inequality: For any non-negative random variable X ,

$$\Pr(X \geq t) \leq \frac{E[X]}{t} \quad \Rightarrow \quad \Pr(X \geq tE[X]) \leq \frac{1}{t}$$

Chebyshev Inequality: For any random variable X and $t > 0$,

$$\Pr(|X - E[X]| \geq t) \leq \frac{\text{Var}[X]}{t^2}$$

Specially (when $t = \epsilon E[X]$),

$$\Pr(|X - E[X]| \geq \epsilon E[X]) \leq \frac{\text{Var}[X]}{\epsilon^2 E^2[X]}$$

Proof: Apply Markov inequality to the random variable $Y = (X - E[X])^2$.

Applying Chebyshev

We need an upper bound on $Var[X]$.

Since X_i 's are independent,

$$Var[X] = Var\left[\frac{1}{s}(X_1 + \dots + X_s)\right] = \frac{1}{s^2}(Var[X_1] + \dots + Var[X_s])$$

Since X_i 's are identical, $Var[X] = \frac{1}{s^2}sVar[X_i] = \frac{1}{s}Var[X_i]$

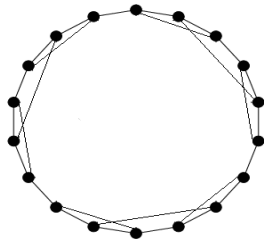
$$Var[X_i] = E[X_i^2] - E^2[X_i] = \left(\frac{d_1^2}{n} + \dots + \frac{d_n^2}{n}\right) - d^2$$

$$\begin{aligned} Pr(|X - E[X]| \geq \epsilon E[X]) &\leq \frac{\frac{1}{s}\left(\frac{d_1^2 + \dots + d_n^2}{n} - d^2\right)}{\epsilon^2 d^2} \\ &= \frac{1}{\epsilon^2 s} \left(n \frac{d_1^2 + \dots + d_n^2}{(d_1 + \dots + d_n)^2} - 1 \right) \end{aligned}$$

How large the term $D = n \frac{d_1^2 + \dots + d_n^2}{(d_1 + \dots + d_n)^2}$ can be?

Lets consider two cases:

3, 3, 3, 3, ..., 3

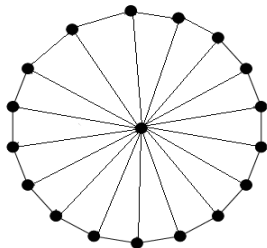


3-regular graph

$$d = 3 \quad D = 1$$

$\Rightarrow s = 1$ is enough

3, 3, 3, ..., 3, $n-1$, 3, ..., 3



wheel graph

$$d \approx 4 \quad D = O(n)$$

$$\Rightarrow s = O\left(\frac{n}{\epsilon^2}\right)$$

- ▶ It can be shown that $D \leq \frac{d_{max}}{d}$ when $d_{max} = \max\{d_i\}$. It suggests $s = O(\frac{d_{max}}{\epsilon^2 d})$ is enough.
- ▶ The above cases tell us we need random $\Omega(n)$ degree queries to distinguish between $d = 3$ and $d \approx 4$.
- ▶ This shows $\frac{3}{4} + \epsilon$ approximation is not possible using $o(n)$ degree queries.
- ▶ Uriel Feige showed that $O(\frac{\sqrt{n}}{\epsilon})$ random degree queries is enough to get a $\frac{1}{2} - \epsilon$ approximation of d .
- ▶ Note that even if we set $\epsilon = \frac{2}{3}$ (for $\frac{1}{3}$ approximation), Chebyshev inequality needs $s = \Omega(n)$.

$$Pr\left(X \leq (1 - \epsilon)E[X] \right) \leq \frac{1}{\epsilon^2 s} (D - 1)$$

In our analysis using Chebyshev inequality, we did not use the fact that d_1, \dots, d_n is the degree sequence of an undirected graph. In other words, the numbers d_1, \dots, d_n have a certain relation. Consider the following sequence:

$$\underbrace{1, 1, 1, \dots, 1, 1}_{n-t}, \underbrace{n, \dots, n}_t \quad \text{average} \approx t$$

Here we need $s = \Omega(n/t)$ samples for a $\alpha > 1/t$ approximation of the average. ☹

But wait! These numbers cannot be degree sequence of a graph.

Union bound

Let E_1, \dots, E_k be a collection of events. Then

$$Pr(E_1 \cup E_2 \cup \dots \cup E_k) \leq \sum_{i=1}^k Pr(E_i)$$

Example: In the above sampling task, let E_i be the event that we sample an n in the i -th round.

$$E = E_1 \cup E_2 \cup \dots \cup E_t \Rightarrow \text{at least one } n \text{ sampled}$$

If we sample s times, by the union bound, we have

$$Pr(E) \leq \sum_{i=1}^s Pr(E_i) = \sum_{i=1}^s \frac{t}{n} = \frac{st}{n}$$

$$\text{If } s < \frac{n}{2t} \Rightarrow Pr(E) < 1/2$$

Lets try a different tool: Chernoff bound

Chernoff Bound: Let $0 \leq \epsilon \leq 1$. Suppose Y_1, \dots, Y_t are independent random variables taking values in the interval $[0, 1]$. Let $Y = \sum_{i=1}^t Y_i$. Then

$$\Pr\left(Y \leq (1 - \epsilon)E[Y]\right) \leq e^{-\frac{\epsilon^2 E[Y]}{2}}$$

$$\Pr\left(Y \geq (1 + \epsilon)E[Y]\right) \leq e^{-\frac{\epsilon^2 E[Y]}{3}}$$

$$\Pr\left(|Y - E[Y]| \geq \epsilon E[Y]\right) \leq 2e^{-\frac{\epsilon^2 E[Y]}{3}}$$

Proof: We present an incomplete proof shortly.

Lets apply Chernoff inequality to our problem.

Recall that $X = \frac{1}{s}(X_1 + \dots + X_s)$ where $X_i \in \{1, \dots, d_{max}\}$.

We define $Y_i = \frac{X_i}{d_{max}} \Rightarrow Y_i \in [0, 1]$.

$$Y = Y_1 + \dots + Y_s \Rightarrow Y = \frac{s}{d_{max}}X$$

$$E[Y] = \frac{s}{d_{max}}d$$

$$\begin{aligned} Pr(|X - E[X]| \geq \epsilon E[X]) &= Pr\left(\left|\frac{sX}{d_{max}} - E\left[\frac{sX}{d_{max}}\right]\right| \geq \epsilon E\left[\frac{sX}{d_{max}}\right]\right) \\ &= Pr(|Y - E[Y]| \geq \epsilon E[Y]) \\ &\leq 2e^{-\frac{\epsilon^2 E[Y]}{3}} = 2e^{-\frac{\epsilon^2 s}{3} \frac{d}{d_{max}}} \end{aligned}$$

A direct application of Chernoff bound suggest $s = O\left(\frac{d_{max}}{d\epsilon^2}\right)$.

This is the same bound that we obtained using Chebyshev!

This is again what we expected because we have not yet used the fact that the numbers d_1, \dots, d_n are the degree sequence of a graph.

In comparison with Chebyshev inequality:

- ▶ Chernoff does not need a knowledge of the variance. It only needs the expectation.
- ▶ Chernoff gives a much higher probability of concentration.

Comparing Chebyshev and Chernoff

Suppose we want to have error probability $\delta < 0$.

Using Chebyshev we should have:

$$\Pr(|X - E[X]| \geq \epsilon E[X]) \leq \frac{1}{\epsilon^2 s} (D - 1) < \frac{1}{\epsilon^2 s} \left(\frac{d_{max}}{d} \right) \leq \delta$$

$$s > \frac{1}{\delta} \frac{d_{max}}{\epsilon^2 d}$$

Using Chernoff we should have:

$$\Pr(|X - E[X]| \geq \epsilon E[X]) \leq 2e^{-\frac{\epsilon^2 s}{3} \frac{d}{d_{max}}} \leq \delta$$

$$s \geq 3 \ln\left(\frac{1}{2\delta}\right) \frac{d_{max}}{\epsilon^2 d}$$

An (incomplete) proof of Chernoff bound

Claim: Let $Y = Y_1 + \dots + Y_t$ where Y_i 's are independent random variables taking values in the interval $[0, 1]$. Let $\mu = E[Y]$. Then

$$\Pr(Y \geq (1 + \epsilon)\mu) \leq \left(\frac{e^\epsilon}{(\epsilon + 1)^{\epsilon+1}}\right)^\mu$$

Proof: Fix $\theta > 0$.

$$\Pr(Y \geq (1 + \epsilon)\mu) = \Pr(e^{\theta Y} \geq e^{\theta(1+\epsilon)\mu})$$

because $f(x) = e^x$ is a monotone function.

$$\Pr(e^{\theta Y} \geq e^{\theta(1+\epsilon)\mu}) \leq \frac{E[e^{\theta Y}]}{e^{\theta(1+\epsilon)\mu}} = \frac{E[e^{\theta(Y_1 + \dots + Y_t)}]}{e^{\theta(1+\epsilon)\mu}}$$

by Markov inequality.

Since Y_i 's are independent, ($E[XY] = E[X]E[Y]$ when Y and X are independent.)

$$\frac{E[e^{\theta(Y_1+\dots+Y_t)}]}{e^{\theta(1+\epsilon)\mu}} = \frac{E[e^{\theta Y_1}] \times \dots \times E[e^{\theta Y_t}]}{e^{\theta(1+\epsilon)\mu}}$$

We show $E[e^{\theta Y_i}] \leq e^{(e^\theta - 1)E[Y_i]}$. Since $Y_i \in [0, 1]$,

$$E[e^{\theta Y_i}] \leq E[1 + (e^\theta - 1)Y_i] = 1 + (e^\theta - 1)E[Y_i] \leq e^{(e^\theta - 1)E[Y_i]}$$

Because for all $x \in [0, 1]$ and $\theta > 0$, we have

$$e^{\theta x} \leq 1 + (e^\theta - 1)x \leq e^{(e^\theta - 1)x}$$

$$\frac{\prod_{i=1}^t E[e^{\theta Y_i}]}{e^{\theta(1+\epsilon)\mu}} \leq \frac{\prod_{i=1}^t e^{(e^\theta - 1)E[Y_i]}}{e^{\theta(1+\epsilon)\mu}} = e^{(e^\theta - 1)\mu - \theta(1+\epsilon)\mu} = e^{((e^\theta - 1) - \theta(1+\epsilon))\mu}$$

The claim follows after setting $\theta = \ln(1 + \epsilon)$.

An application of Chernoff bound

Amplifying the success probability

Suppose we have a randomized algorithm A that processes the input data D and approximate some $f(D)$ where

$$|A(D) - f(D)| \leq \epsilon f(D) \text{ with probability at least } 3/4.$$

How to amplify the success probability of A ?

We want to have a randomized algorithm A' with error probability $\delta \ll 1/4$.

Idea: Run A on input data D , $O(\ln(\frac{1}{\delta}))$ times and output the **median** of the outcomes.

Each (independent) repetition of A succeeds with probability $3/4$. Suppose a_i is the outcome of i -th repetition. We have

$$\Pr(|a - f(D)| \geq \epsilon f(A)) \leq 1/4.$$

We define $X_i = 1$ if i -th repetition is good (its error is less than $\epsilon f(A)$), otherwise we let $X_i = 0$.

$X = X_1 + \dots + X_t$ is the number of good outcomes in t repetitions.

The median of $\{a_1, \dots, a_t\}$ is bad \Rightarrow Less than $t/2$ repetitions are good. In other words, $X < t/2$.

By Chernoff bound, we have

$$\Pr(\text{median is bad}) \leq \Pr(X < t/2) \leq e^{O(-t)} \leq \delta \Rightarrow t = (\ln(\frac{1}{\delta}))$$