

CS49/249 (Randomized Algorithms), Spring 2021 : Lecture 4

Topic: Estimation Algorithms, Chebyshev, Chernoff-Hoeffding Bound

Disclaimer: These notes have not gone through scrutiny and in all probability contain errors.

Please discuss in Piazza/email errors to deeparnab@dartmouth.edu

- One extremely important application of randomized algorithms is in *estimation*. Imagine there is a statistic stat you are interested in. For concreteness, let us assume stat is the fraction of Hanover's population who are vegetarians. To calculate stat exactly we must ask everyone in Hanover of their dietary preference. However, the exact statistic is often not important, and what one really needs is an *estimate* est of the statistic. And polling or *random sampling* is the way to obtain such estimates. In this lecture, we will set the formal definitions of what a "good" estimate is, and what one needs to obtain one.
- For any statistic stat , an estimate est is a *random variable* often obtained via a randomized algorithm. We say est is an *unbiased estimate* if

$$\mathbf{Exp}[\text{est}] = \text{stat} \quad (\text{Unbiased Estimate})$$

For the statistic fraction of vegetarians in Hanover, here is an unbiased estimate: sample an individual from the Hanover population *uniformly at random*¹ and ask them about their dietary preference. If they say vegetarian, set $\text{est} = 1$, otherwise set $\text{est} = 0$. If we denote the population of Hanover as n , then there are $\text{stat} \cdot n$ many vegetarians, and $(1 - \text{stat})n$ many non-vegetarians. Thus, we get

$$\mathbf{Exp}[\text{est}] = \frac{\text{stat} \cdot n}{n} \cdot 1 + \frac{(1 - \text{stat})n}{n} \cdot 0 = \text{stat}$$

- Although the above estimate was unbiased, all of us really can feel that it was a pretty bad estimate. Indeed, it was always 0 or 1. The following definition gives a much more nuanced notion of the quality of an estimate.

Definition 1 ((ϵ, δ) -estimate.). An estimate est is a multiplicative (ϵ, δ) -estimate of a statistic stat if

$$\Pr[\text{est} \notin (1 \pm \epsilon) \cdot \text{stat}] \leq \delta$$

That is, with probability at least $(1 - \delta)$ (which is high if δ is small), the estimate satisfies $(1 - \epsilon)\text{stat} \leq \text{est} \leq (1 + \epsilon)\text{stat}$.

There is another notion of (ϵ, δ) -additive estimate often used when the statistic stat is bounded, as in the case of fraction of vegetarians where $\text{stat} \in [0, 1]$.

Definition 2 ((ϵ, δ) -additive estimate.). An estimate est is an additive (ϵ, δ) -estimate of a statistic stat if

$$\Pr[\text{est} \notin \text{stat} \pm \epsilon] \leq \delta$$

That is, with probability at least $(1 - \delta)$ (which is high if δ is small), the estimate satisfies $(\text{stat} - \epsilon) \leq \text{est} \leq (\text{stat} + \epsilon)$.

¹We are completely ignoring **how** one samples an individual uniformly at random from the population. We are also assuming that they answer our poll. All these are super-important considerations and are real bottlenecks to good polling. Unfortunately, it is not in the scope of this course.

- **Variance of an estimator.** The main result we want to establish today is to show how to obtain (ε, δ) -estimates from unbiased estimates. It is an extremely important fact which will be used many times in this course (and life).

Theorem 1. [Boosting Theorem or the Median-of-Means Theorem.]

Let $\widehat{\text{est}}$ be an *unbiased* estimator of some statistic stat . Then, one can obtain an (ε, δ) -multiplicative estimate of stat using K independent samples of $\widehat{\text{est}}$, where

$$K = \frac{C \text{Var}[\widehat{\text{est}}]}{(\text{Exp}[\widehat{\text{est}}])^2} \cdot \frac{1}{\varepsilon^2} \cdot \ln\left(\frac{2}{\delta}\right)$$

where C is some constant. Consequently, one can obtain an (ε, δ) -additive estimate of stat using K' independent samples of $\widehat{\text{est}}$, where $K' = \frac{C \text{Var}[\widehat{\text{est}}]}{\varepsilon^2} \cdot \ln\left(\frac{2}{\delta}\right)$

Thus, what really determines the quality of an unbiased estimate is its variance. If the variance, or more precisely the variance-to-mean-square ratio (which is the square of the standard-deviation-to-mean ratio). If that ratio is small, then one can obtain a good multiplicative estimate with not “too many” runs of the unbiased estimator.

Let’s go back to the example of stat being the fraction of vegetarians in a population. We saw an unbiased estimate $\widehat{\text{est}}$. What is its variance? A calculation gives that

$$\text{Var}[\widehat{\text{est}}] = \text{stat}(1 - \text{stat})$$

Thus, if one obtains $K = O\left(\frac{(1-\text{stat}) \ln(1/\delta)}{\text{stat} \cdot \varepsilon^2}\right)$ many independent samples, or in other words polls that many people independently uniformly at random, then one can come with a very good estimate. As you can see, the number grows as stat becomes smaller: this is to be expected, if the number of vegetarians are small, we will need to sample more to detect them. On the other hand, if we only want to estimate the fraction of vegetarians to within *additive* 1%, say (so $\varepsilon = 0.01$), then the number of samples we need is $O(\text{stat}(1 - \text{stat}) \cdot \ln(1/\delta))$. Note, and this is something often counter-intuitive to many, this *doesn’t* depend on the population size. The same *number* of people need be sampled² even from New York City.

- **First obtain an $(\varepsilon, \frac{1}{3})$ -estimate using Chebyshev.** We will prove [Theorem 1](#) in two steps. Both steps will also introduce tools which are much more important than the theorem itself.

The first idea is to take a bunch of unbiased estimates and creating a more refined estimate by taking the mean. Here is the algorithm.

1: **procedure** BETTER-ESTIMATOR(s): \triangleright *Assumes access to an unbiased estimate $\widehat{\text{est}}$.*
 2: Sample s independent unbiased estimates $\widehat{\text{est}}_0, \widehat{\text{est}}_1, \dots, \widehat{\text{est}}_s$.
 3: Return $\text{est}' := \frac{1}{s} \sum_{i=1}^s \widehat{\text{est}}_i$.

Claim 1. Let $s \geq \frac{3 \text{Var}[\widehat{\text{est}}]}{\varepsilon^2 (\text{Exp}[\widehat{\text{est}}])^2}$. Then, est' returned by BETTER-ESTIMATOR(s) satisfies

$$\Pr[\text{est}' \notin (1 \pm \varepsilon)\text{stat}] \leq \frac{1}{3}$$

²assuming both cities have the similar fraction of vegetarians.

To prove this, we will use *Chebyshev inequality* which is the most general purpose deviation inequality.

Theorem 2. (Chebyshev’s Inequality)

Let X be any random variable. Then for any $t > 0$, we have

$$\Pr[|X - \mathbf{Exp}[X]| \geq t] \leq \frac{\mathbf{Var}[X]}{t^2}$$

Proof. We first note that

$$\Pr[|X - \mathbf{Exp}[X]| \geq t] = \Pr[(X - \mathbf{Exp}[X])^2 \geq t^2]$$

Then we notice that $D := (X - \mathbf{Exp}[X])^2$ is a non-negative random variable, and therefore we can apply Markov’s inequality on it to get

$$\Pr[|X - \mathbf{Exp}[X]| \geq t] = \Pr[D \geq t^2] \leq \frac{\mathbf{Exp}[D]}{t^2} = \frac{\mathbf{Var}[X]}{t^2}$$

□

Proof of Claim 1.

- We first calculate the variance of est' . We use two simple facts: $\mathbf{Var}[cX] = c^2 \mathbf{Var}[X]$, and the linearity of variance³ for *independent* random variables X_1, \dots, X_t , $\mathbf{Var}[\sum_{i=1}^t X_i] = \sum_{i=1}^t \mathbf{Var}[X_i]$. Using these, we obtain

$$\mathbf{Var}[\text{est}'] = \mathbf{Var}\left[\frac{1}{s} \sum_{i=1}^s \widehat{\text{est}}_i\right] = \frac{\mathbf{Var}[\widehat{\text{est}}]}{s}$$

Also, recall that $\mathbf{Exp}[\text{est}'] = \frac{1}{s} \sum_{i=1}^s \mathbf{Exp}[\widehat{\text{est}}_i] = \text{stat}$ since each $\widehat{\text{est}}_i$ is an unbiased estimate of stat . Thus, taking the mean of s independent samples is still unbiased but the variance decays. Useful nugget of information.

- Next, we apply Chebyshev on est' to get

$$\begin{aligned} \Pr[\text{est}' \notin (1 \pm \varepsilon)\text{stat}] &= \Pr[|\text{est}' - \mathbf{Exp}[\text{est}']| \geq \varepsilon \mathbf{Exp}[\widehat{\text{est}}]] \\ &\leq \frac{\mathbf{Var}[\text{est}']}{\varepsilon^2 (\mathbf{Exp}[\widehat{\text{est}}])^2} = \frac{\mathbf{Var}[\widehat{\text{est}}]}{s\varepsilon^2 (\mathbf{Exp}[\widehat{\text{est}}])^2} \leq \frac{1}{3} \quad \square \end{aligned}$$

- **The Chernoff Bound.** Claim 1 is great, but the “error probability” is still $\frac{1}{3}$ which we would like to ramp down to δ . If you stare at the proof of the claim, you probably notice that if instead of using $s \geq 3 \cdot \left(\frac{\mathbf{Var}}{\varepsilon^2 \mathbf{Exp}^2}\right)$ samples, we used $\frac{1}{\delta} \cdot \left(\frac{\mathbf{Var}}{\varepsilon^2 \mathbf{Exp}^2}\right)$ many we would have our (ε, δ) -estimate. However, the next technique shows how we can actually have a muuuuuuuch better dependence on δ . This is really important if we want our confidence levels really high.

³I guess this is also an important fact which you should recall from CS30/Discrete Probability.

To explain this best, we introduce one of the *most used* theorems in TCS: the Chernoff bound, or the Azuma-Bernstein-Chernoff-Hoeffding-Stein... bound⁴. Indeed, we will dedicate the rest of this lecture to understanding this bound, and continue with the proof of [Theorem 1](#) in the next class.

Theorem 3 (Chernoff Bounds). Let X_1, X_2, \dots, X_n be *independent Bernoulli* random variables with each $X_i \in \{0, 1\}$. Let $X = \sum_{i=1}^n X_i$. Then, for any $\varepsilon \in (0, 1)$,

$$\Pr[X \leq (1 - \varepsilon) \mathbf{Exp}[X]] \leq e^{-\frac{\varepsilon^2 \mathbf{Exp}[X]}{2}} \quad (\text{LT})$$

and

$$\Pr[X \geq (1 + \varepsilon) \mathbf{Exp}[X]] \leq e^{-\frac{\varepsilon^2 \mathbf{Exp}[X]}{3}} \quad (\text{UT1})$$

For the “upper tail”, that is for “larger” deviations, we have when $1 \leq t \leq 4$, we have the following (changing ε to t so as to underscore that the deviation is big)

$$\Pr[X \geq (1 + t) \mathbf{Exp}[X]] \leq e^{-\frac{t^2 \mathbf{Exp}[X]}{4}} \quad (\text{UT2})$$

and for $t > 4$ (really large), we have

$$\Pr[X \geq (1 + t) \mathbf{Exp}[X]] \leq e^{-\frac{t \ln t \mathbf{Exp}[X]}{2}} \quad (\text{UT3})$$

Remark: Important: *Equations (UT1) to (UT3) hold with all $\mathbf{Exp}[X]$ occurrences replaced by any upper bound $\mathbf{Exp}[X] \leq U$.*

Remark: *Note the asymmetry in the denominators in the exponent. For most applications this is not important. What are qualitatively more important : (a) things are in the exponent, (b) the expectation of the sum shows up in the exponent irrespective of the number of terms, (c) the X_i 's are independent but not necessarily identical, and (d) the dependence on ε in the exponent is quadratic.*

It is instructive to compare with what Chebyshev gives us. To apply Chebyshev, set $t = \varepsilon \mathbf{Exp}[X]$ to get $\Pr[X \geq (1 + \varepsilon) \mathbf{Exp}[X]] \leq \frac{\mathbf{Var}[X]}{\varepsilon^2 \mathbf{Exp}^2[X]}$. Unless $\mathbf{Var}[X] \leq \varepsilon^2 \mathbf{Exp}^2[X]$, the RHS is a trivial bound. The Chernoff-bound always gives something non-trivial.

- *Some consequences.*

- *Number of heads.* Suppose we toss n fair coins and let X be the number of heads we observe. We know $\mathbf{Exp}[X] = \frac{n}{2}$. But what really can we say about the “range” of X ? In particular, what is the probability see t more heads than expected, that is, $\Pr[X \geq \frac{n}{2} + t]$? Or rather for what value of t does this get below probability δ ?

⁴There are many such inequalities which bound how far the sum of random variables can deviate from its expectations. [Entire books have been written](#) on this (pdf copies of many of these are freely available).

This random variable (number of heads in n fair coin tosses) has a name : it is called the *Binomial Distribution* or *Binomial Random Variable*. And the exact answer to the above question is:

$$\Pr[X \geq \frac{n}{2} + t] = \frac{1}{2^n} \sum_{j \geq \frac{n}{2} + t} \binom{n}{j}$$

and you can sweat hard and figure it out. But it still doesn't answer after what t will the probability be at most δ . Is this $t = \Theta(1)$ or $\Theta(\log n)$ or $\Theta(n)$ or what? We now see that the Chernoff bound will lead to the answer.

To this end, define n independent random variables where $X_i = 1$ if the i th coin toss is head, and $X_i = 0$ otherwise. Next, observe $X = \sum_{i=1}^n X_i$ and that $\mathbf{Exp}[X] = \frac{n}{2}$. Apply [Theorem 3](#) with ε such that $\varepsilon \mathbf{Exp}[X] = t$, that is, $\varepsilon = \frac{2t}{n}$.

$$\Pr[X \geq \frac{n}{2} + t] \leq e^{-\frac{4t^2}{n^2} \cdot \frac{n}{6}} = e^{-2t^2/3n}$$

Therefore, if $t \geq \sqrt{\frac{3n}{2} \cdot \ln(\frac{1}{\delta})}$, the RHS probability is $\leq \delta$. The answer is $\Theta(\sqrt{n})$. To put in some numbers, the chances of observing $\frac{n}{2} + 20\sqrt{n}$ is less than 10^{-100} .

Exercise: Repeat a similar calculation to figure for what t we get $\Pr[X \leq \frac{n}{2} - t] \leq \delta$.

Remark: By the way, the variance of X is $\frac{n}{4}$ and thus the standard deviation of X is $\frac{\sqrt{n}}{2}$. Indeed, in this case the Chernoff bound tells us that the probability of seeing the number of heads more than c std-deviations away drops like $e^{-O(c^2)}$. Chebyshev on the other hand says that the probability drops like $\frac{1}{c^2}$.

- *The Fraction of Vegetarians.* For the simple estimation problem of fraction of vegetarians in Hanover, the Chernoff bound already solves the problem without needing [Theorem 1](#). Indeed, let us pick k samples $\widehat{\text{est}}_1, \dots, \widehat{\text{est}}_k$ and return $\text{est} := \frac{\sum_{i=1}^k \widehat{\text{est}}_i}{k}$. Let $X_i := \widehat{\text{est}}_i$ and let $X = \sum_{i=1}^k X_i$. We get $\mathbf{Exp}[X] = \text{stat} \cdot k$, and thus,

$$\Pr[\text{est} \geq (1 + \varepsilon)\text{stat}] = \Pr[X \geq (1 + \varepsilon)\text{stat} \cdot k] \leq e^{-\frac{\varepsilon^2}{3} \cdot \text{stat} \cdot k}$$

Therefore, if we choose k such that $k \cdot \text{stat} \cdot \frac{\varepsilon^2}{3} \geq \ln(2/\delta)$, we would get that $\Pr[\text{est} \geq (1 + \varepsilon)\text{stat}] \leq \frac{\delta}{2}$. The application of the other Chernoff bound would give $\Pr[\text{est} \leq (1 - \varepsilon)\text{stat}] \leq \frac{\delta}{2}$. Which would imply $\Pr[\text{est} \notin (1 \pm \varepsilon)\text{stat}] \leq \delta$, that is, est is an (ε, δ) -multiplicative estimate.

This gives us $k \geq \frac{3 \ln(2/\delta)}{\text{stat} \cdot \varepsilon^2}$ samples are sufficient (as we had found before, up to constants).

Learning Tidbits:

- Algorithm Design: When you design an unbiased estimate $\widehat{\text{est}}$, it is the $\mathbf{Var}[\widehat{\text{est}}] / \mathbf{Exp}^2[\widehat{\text{est}}]$ which defines the quality. If this is "small", then with that many samples times $\frac{1}{\varepsilon^2} \ln(1/\delta)$, one can get an (ε, δ) -estimate est .
- Analysis: **Chernoff Bounds!** (Hard to overstate its importance)

