

# CMPT 210: Probability and Computation

## Lecture 18

---

Sharan Vaswani

July 15, 2022

I am away at a conference next week (from 18-25 July), so there will be no office hours on Tuesday (19 July). Please go to the TA office hours on Thursday (21 July) for questions.

On Tuesday (19 July), I will do an online/recorded lecture on the application of probability to machine learning. The TA (Yasaman) will oversee the lecture in person and help answer questions. She is an expert on the topic, so please come to class.

On Friday (22 July), we will have a recorded lecture on analyzing Randomized QuickSort. I will release the lecture today on Piazza and you can watch it from home (no class). I can answer questions about it on Tuesday (26 July).

Late submission for Assignment 3 is on Tuesday, 19 July (submit to Yasaman). The solutions will be released after class, so no submissions are allowed after that.

Assignment 4 will be released on Friday, 22 July, and will be due on Tuesday, 2 August with late submission on August 5 (last class).

**Tail inequalities** bound the probability that the r.v. takes a value much different from its mean.

**Markov's Theorem:** If  $X$  is a non-negative random variable, then for all  $x > 0$ ,  
 $\Pr[X \geq x] \leq \frac{\mathbb{E}[X]}{x}$ .

**Chebyshev's Theorem:** For a r.v.  $X$  and all  $x > 0$ ,  $\Pr[|X - \mathbb{E}[X]| \geq x] \leq \frac{\text{Var}[X]}{x^2}$ .

**Q:** Suppose there is an election between two candidates  $A$  and  $B$ , and we are hired by candidate  $A$ 's election campaign to estimate the chances of  $A$  winning the election. In particular, we want to estimate  $p$ , the fraction of voters favoring  $A$  before the election. We conduct a voter poll – selecting (typically calling) people uniformly at random (with replacement so that we can choose a person twice) and try to estimate  $p$ . What is the number of people we should poll to estimate  $p$  reasonably accurately and with reasonably high probability?

Let us define  $X_i$  to be the indicator r.v. which is equal to 1 if person  $i$  that we called favors candidate  $A$ . The  $X_i$  r.v.'s are mutually independent since the people we poll are chosen randomly, and we assume that they are identically distributed meaning that  $X_i = 1$  with probability  $p$ .

Suppose we poll  $n$  people and define  $S_n := \sum_{i=1}^n X_i$  as the r.v. equal to the total number of people who prefer candidate  $A$  (amongst the people we polled).  $\frac{S_n}{n}$  is the fraction of polled voters who favor candidate  $A$  and is the *statistical estimate* of  $p$ .

**Q:** What is the distribution of  $S_n$ ? **Ans:** Since the voters are independent and each has probability  $p$  of favoring candidate  $A$ ,  $S_n \sim \text{Bin}(n, p)$ .

Hence, we want to find for what  $n$  is our estimate for  $p$  accurate up to an error  $\epsilon > 0$  and with probability  $1 - \delta$  (for  $\delta \in (0, 1)$ ). Formally, for what  $n$  is,

$$\Pr \left[ \left| \frac{S_n}{n} - p \right| < \epsilon \right] \geq 1 - \delta$$

Since  $S_n \sim \text{Bin}(n, p)$ ,  $\mathbb{E}[S_n] = np$  and hence,  $\mathbb{E} \left[ \frac{S_n}{n} \right] = p$ , meaning that our estimate is *unbiased* – in expectation, the estimate is equal to  $p$ . Hence, the above statement is equivalent to,

$$\Pr \left[ \left| \frac{S_n}{n} - \mathbb{E} \left[ \frac{S_n}{n} \right] \right| < \epsilon \right] \geq 1 - \delta$$

Hence, we can use Chebyshev's Theorem for the r.v.  $\frac{S_n}{n}$  with  $x = \epsilon$  to bound the LHS

$$\Pr \left[ \left| \frac{S_n}{n} - \mathbb{E} \left[ \frac{S_n}{n} \right] \right| < \epsilon \right] = 1 - \Pr \left[ \left| \frac{S_n}{n} - \mathbb{E} \left[ \frac{S_n}{n} \right] \right| \geq \epsilon \right] \geq 1 - \frac{\text{Var}[S_n/n]}{\epsilon^2}.$$

Hence, the problem now is to find  $n$  such that,

$$1 - \frac{\text{Var}[S_n/n]}{\epsilon^2} \geq 1 - \delta \implies \frac{\text{Var}[S_n/n]}{\epsilon^2} < \delta$$

Let us calculate the  $\text{Var}[S_n/n]$ .

$$\begin{aligned}\text{Var}[S_n/n] &= \frac{1}{n^2} \text{Var}[S_n] && \text{(Using the property of variance)} \\ &= \frac{1}{n^2} n p (1 - p) = \frac{p(1-p)}{n} && \text{(Using the variance of the Binomial distribution)}\end{aligned}$$

Hence, we want to find  $n$  s.t.

$$\frac{p(1-p)}{n\epsilon^2} < \delta \implies n \geq \frac{p(1-p)}{\epsilon^2 \delta}$$

But we do not know  $p$ ! If  $n \geq \max_p \frac{p(1-p)}{\epsilon^2 \delta}$ , then for any  $p$ ,  $n \geq \frac{p(1-p)}{\epsilon^2 \delta}$ . So the problem is to compute  $\max_p \frac{p(1-p)}{\epsilon^2 \delta}$ . This is a concave function and is maximized at  $p = 1/2$ .

Hence, if  $n \geq \frac{1}{4\epsilon^2 \delta}$ , then  $\Pr[|\frac{S_n}{n} - p| < \epsilon] \geq 1 - \delta$  meaning that we have estimated  $p$  upto an error  $\epsilon$  and this bound is true with high probability equal to  $1 - \delta$ .

For example, if  $\epsilon = 0.01$  and  $\delta = 0.01$  meaning that we want the bound to hold 99% of the time, then, we require  $n \geq 250000$ .

## Pairwise Independent Sampling

Let  $G_1, G_2, \dots, G_n$  be pairwise independent random variables with the same mean  $\mu$  and standard deviation  $\sigma$ . Define  $S_n := \sum_{i=1}^n G_i$ , then,

$$\Pr \left[ \left| \frac{S_n}{n} - \mu \right| \geq \epsilon \right] \leq \frac{1}{n} \left( \frac{\sigma}{\epsilon} \right)^2.$$

Let us compute  $\mathbb{E}[S_n/n]$  and  $\text{Var}[S_n/n]$ .

$$\mathbb{E}[S_n] = \mathbb{E} \left[ \sum_{i=1}^n G_i \right] = \sum_{i=1}^n \mathbb{E}[G_i] = n\mu \implies \mathbb{E}[S_n/n] = \frac{1}{n} \mathbb{E}[S_n] = \mu$$

(Using linearity of expectation)

$$\text{Var}[S_n] = \text{Var} \left[ \sum_{i=1}^n G_i \right] = \sum_{i=1}^n \text{Var}[G_i] = n\sigma^2$$

(Using linearity of variance for pairwise independent r.v's)

$$\implies \text{Var}[S_n/n] = \frac{1}{n^2} \text{Var}[S_n] = \frac{\sigma^2}{n}$$

# Pairwise Independent Sampling

Using Chebyshev's Theorem,

$$\Pr \left[ \left| \frac{S_n}{n} - \mathbb{E} \left[ \frac{S_n}{n} \right] \right| \geq \epsilon \right] = \Pr \left[ \left| \frac{S_n}{n} - \mu \right| \geq \epsilon \right] \leq \frac{\text{Var}[S_n/n]}{\epsilon^2} = \frac{\sigma^2}{n\epsilon^2}$$

Hence, for arbitrary pairwise independent r.v.'s, if  $n$  increases, the probability of deviation from the mean  $\mu$  decreases.

**Weak Law of Large Numbers:** Let  $G_1, G_2, \dots, G_n$  be pairwise independent variables with the same mean  $\mu$  and (finite) standard deviation  $\sigma$ . Define  $T_n := \frac{\sum_{i=1}^n G_i}{n}$ , then for every  $\epsilon > 0$ ,

$$\lim_{n \rightarrow \infty} \Pr[|T_n - \mu| \leq \epsilon] = 1.$$

Follows from the theorem on pairwise independent sampling since

$$\lim_{n \rightarrow \infty} \Pr[|T_n - \mu| \leq \epsilon] = \lim_{n \rightarrow \infty} \left[ 1 - \frac{\sigma^2}{n\epsilon^2} \right] = 1.$$



Questions?