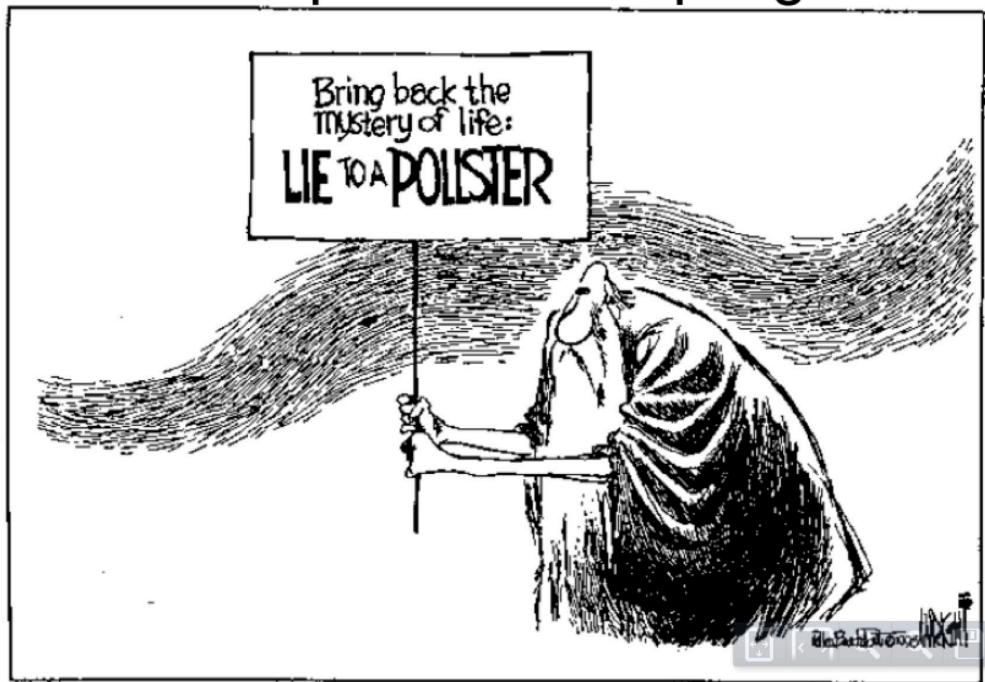


# Stat 88: Prob. & Math. Statistics in Data Science



Lecture 18: 3/29/2022

Sampling without replacement, the law of averages,  
distribution of a sample sum

7.2, 7.3, 8.1

## Review: variance, SD, inequalities when we don't know dsn

$$E(X) = \mu \quad E((X - \mu)^2) = E(X^2) - \mu^2$$

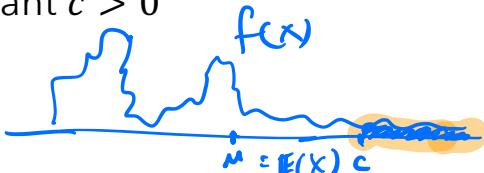
1. The variance of a rv  $X$ :  $Var(X) = \sigma^2 = E[(X - E(X))^2] = E(X^2) - (E(X))^2$
2. The SD or standard deviation is given by  $SD(X) = \sigma = \sqrt{Var(X)} \geq 0$

If  $X$  and  $Y$  are *independent*, then

$$Var(X_1 + X_2 + \dots + X_n) = Var(X_1) + Var(X_2) + \dots + Var(X_n)$$

3.  $Var(X + Y) = Var(X) + Var(Y)$  &  $SD(X + Y) = \sqrt{Var(X) + Var(Y)}$
4. **Markov's Inequality:** For a nonnegative rv  $X$ , and constant  $c > 0$

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

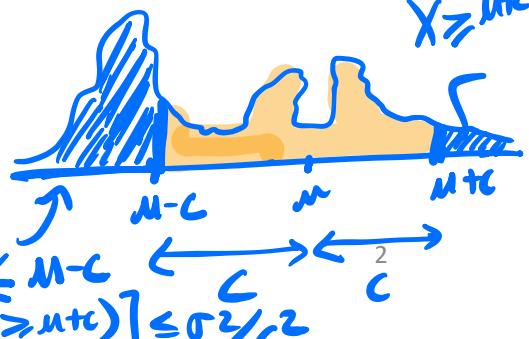


5. **Chebyshev's inequality:** For any random variable  $X$  (not necessarily non-negative), with mean  $\mu$  and standard deviation  $\sigma$ , for any positive constant  $c > 0$ , we have:

$$P(|X - \mu| \geq c) \leq \frac{\sigma^2}{c^2} = \frac{Var(X)}{c^2}$$

dist b/w

$X \& \mu$



$$P[(X \leq \mu - c) \text{ or } (X \geq \mu + c)] \leq \frac{\sigma^2}{c^2}$$

$$P(|X-\mu| \geq k\sigma) \leq \frac{1}{k^2}$$

Example

$$P(\mu - c < X < \mu + c) = 1 - P\left(\frac{|X-\mu|}{\sigma} \geq \frac{c}{\sigma}\right)$$

$$\geq 1 - \frac{\sigma^2}{c^2}$$

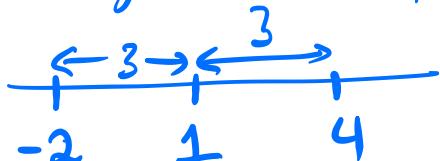
A list of non negative numbers has an average of 1 and an SD of 2. Let  $p$  be the proportion of numbers over 4. To get an upper bound for  $p$ , you should:

- a) Assume a binomial distribution.
- b) Use Markov's inequality.
- c) Use Chebyshev's inequality
- d) None of the above.

$p = \text{proportion of } \#s \text{ over } 4$   
 If  $X$  is a # picked at random from list.  $E(X) = 1, SD(X) = 2$

(b)  $P(X > 4) \leq \frac{1 - E(X)}{4}$

(c) Chebyshev's inequality



$$P(|X-1| \geq 3) =$$

$$P(X \leq -2) + P(X \geq 4) \geq P(X \geq 4)$$

$$P(X \geq 4) \leq [P(X \leq -2) + P(X \geq 4)] \leq \frac{\sigma^2}{c^2} = \frac{2^2}{3^2} = \frac{4}{9}$$

Chebyshov's Ineq  $P(|X-\mu| \geq c) \leq \frac{\sigma^2}{c^2}$

$$P(X > 4)$$



$$P(|X-1| \geq 3)$$

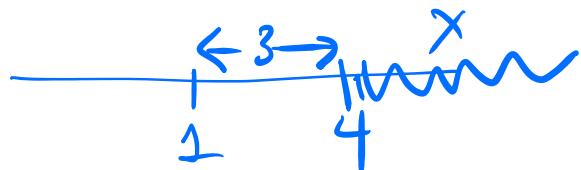
$$= P(X \leq -2) + P(X \geq 4)$$

$$\geq P(X \geq 4) \geq P(X > 4)$$

$$P(X > 4) \leq P(X \geq 4) \leq \frac{\sigma^2}{c^2} = \frac{2^2}{3^2} = \frac{4}{9}$$

want

$$P(|X-1| \geq 3)$$



Answer: (b) Markov's inequality gives a better bound.

## Recap: sums of iid random variables

$$X_1 + X_2 + X_3 + \dots + X_{15}$$

- Let  $X_1, X_2, \dots, X_n$  be independent and identically distributed random variables with mean  $\mu$  and variance  $\sigma^2$ . Define  $S_n = X_1 + X_2 + \dots + X_n$
- $E(S_n) = \sum E(X_k) = n\mu$  &  $Var(S_n) = Var(X_1) + \dots + Var(X_n) = n\sigma^2$
- $SD(S_n) = \sqrt{n}\sigma$  (Square root law: sd grows by a factor of  $\sqrt{n}$ , not  $n$ )
- $X \sim Bin(n, p)$ :  $E(X) = np, Var(X) = np(1 - p), SD(X) = \sqrt{np(1 - p)}$   
 $\sqrt{npq}, q = 1 - p$
- $X \sim Poisson(\mu)$ :  $E(X) = Var(X) = \mu, SD = \sqrt{\mu}$
- $X \sim \text{Geometric}(p)$  distribution:  $E(X) = \frac{1}{p}, Var(X) = \frac{1-p}{p^2}$   
$$\left. \begin{array}{l} \text{Fact. we} \\ \text{need to} \\ \text{prove.} \end{array} \right\}$$

## 7.3: Sampling without replacement

- When we have a simple random sample (SRS), the draws are without replacement (like drawing cards from a deck).
- The random variables are no longer independent
- So, how do we compute the variance of the sum of draws of a SRS?
- To begin with, let's look at the squares and products of indicators
- If  $I_A$  and  $I_B$  are indicator functions, what can we say about  $I_A^2$  and  $I_A I_B$ ?

Draw tickets from a box without replacement

$I_A$  is an indicator for event A     $E(I_A) = P(A)$   
 $I_B$  "    "      "      B     $I_A = \begin{cases} 1 & \text{if } A \text{ true} \\ 0 & \text{if } A \text{ not true.} \end{cases}$   
 $E(I_B) = P(B)$      $I_A^2 = \begin{cases} 1^2, \text{A true} \\ 0^2, \text{if A false.} \end{cases}$

$$\begin{aligned} E(I_A^2) &= 1 \cdot P(A) + 0 \cdot P(A^c) \\ &= P(A) \end{aligned}$$
$$I_A \cdot I_B = \begin{cases} 1, & \text{if BOTH A \& B are true} \\ 0, & \text{o/w} \end{cases}$$

$E(I_A I_B) = P(A \cap B)$

$$\mathbb{E}(I_j I_k) = P(I_j=1 \& I_k=1) = P(I_j=1) P(I_k=1 | I_j=1)$$

Variance of a hypergeometric random variable

*total popn.* # of "good" tickets, n=sample size

- Let  $X \sim HG(N, G, n)$ ; then can write  $X = I_1 + I_2 + \dots + I_n$ , where  $I_k$  is the indicator of the event that the kth draw is good.

$$\mathbb{E}(X) = \mathbb{E}(I_1) + \mathbb{E}(I_2) + \dots = \frac{G}{N} + \frac{G}{N} + \dots + \frac{G}{N} = \frac{nG}{N} = G + G + \dots + \frac{G}{N}$$

- We can compute the expectation of  $X$  using symmetry:  $E(X) = \frac{nG}{N}$
- But what about variance?
- Since the indicators are not independent, we can't just add the variances
- Let's just use the formula:  $Var(X) = E(X^2) - (E(X))^2 = E(X^2) - \left(\frac{nG}{N}\right)^2$

$$X^2 = (I_1 + I_2 + \dots + I_n)^2 = \sum_{k=1}^n I_k^2 + \sum_j \sum_{k,k \neq j} I_j I_k$$

$$E(X^2) = nE(I_k^2) + n(n-1)E(I_j I_k)$$

$$= n \frac{G}{N} + n(n-1)P(I_j = 1)P(I_k = 1 | I_j = 1)$$

$$E(X^2) = n \frac{G}{N} + n(n-1) \frac{G}{N} \cdot \frac{G-1}{N-1}$$

$$X = I_1 + I_2 + I_3 + \dots + I_n, \quad E(X) = n \cdot \frac{G}{N}$$

$$\text{Var}(X) = E(X^2) - (E(X))^2 = E(X^2) - \left(\frac{nG}{N}\right)^2$$

$$X^2 = (I_1 + I_2 + \dots + I_n)^2$$

For example, if  $n=2$

$$(I_1 + I_2)^2 = I_1^2 + I_2^2 + I_1 I_2 + I_2 I_1$$

$$X^2 = \underbrace{I_1^2 + I_2^2 + \dots + I_n^2}_{j \neq k} + \sum_{j=1}^n \sum_{k=1}^n I_j I_k$$

$$n=3$$

$$(I_1 + I_2 + I_3)(I_1 + I_2 + I_3)$$

$$I_1^2 + I_1 I_2 + I_1 I_3 + I_2 I_1 + I_2^2 + I_2 I_3 + I_3 I_1 + I_3 I_2 + I_3^2$$

$$E(I_1^2) = E(I_i) = \frac{G}{N}$$

$$= n \cdot E(I_1^2) + \sum_{j=1}^n \sum_{\substack{k=1 \\ j \neq k}}^n E(I_j I_k)$$

$$= n \cdot \underbrace{E(I_1^2)}_{\text{all the same}} + n(n-1) \underbrace{E(I_j I_k)}$$

$$= n \cdot \frac{G}{N} + n(n-1) \underbrace{\frac{G}{N} \cdot \frac{G-1}{N-1}}_{P(I_j=1)P(I_k=1|I_j=1)}$$

$$E(X^2) = n \cdot \frac{G}{N} + n(n-1) \cdot \frac{G}{N} \cdot \frac{G-1}{N-1}$$

$$\text{Var}(X) = E(X^2) - (E(X))^2$$

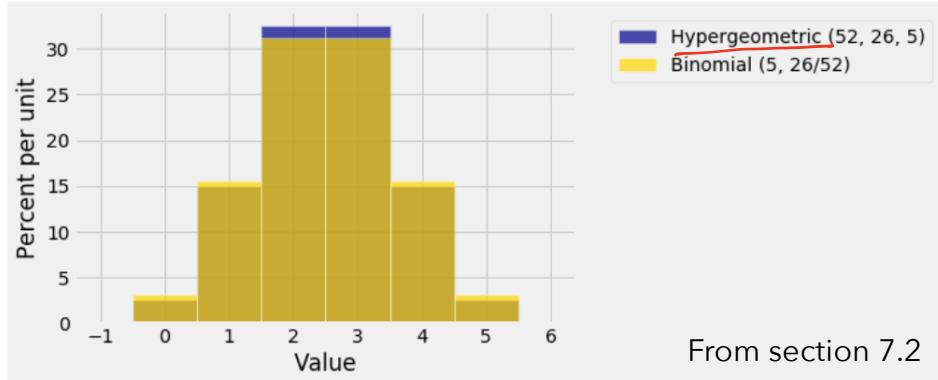
$$= \frac{nG}{N} + n(n-1) \frac{G}{N} \cdot \frac{G-1}{N-1} - \left(\frac{nG}{N}\right)^2$$

$$\mathbb{E}(I_j I_k) = P(I_j=1 \& I_k=1) = P(I_j=1) P(I_k=1 | I_j=1)$$

Variance of a hypergeometric random variable

$$\begin{aligned}
 \text{Var}(X) &= \frac{nG}{N} \left[ 1 + \frac{(n-1)(G-1)}{\frac{(N-1)}{N}} - \frac{nG}{N} \right] \\
 &= \frac{nG}{N} \left[ \frac{\cancel{N^2-N} + \cancel{(nN-n)(G-1)}}{\cancel{N(N-1)} + \cancel{(n-1)(G-1)}N} - \frac{nGN+nG}{nG(N-1)} \right] \\
 &= \frac{nG}{N} \left[ \cancel{N^2-N} + nNG - NG - nN + \cancel{N-nG} + \cancel{nG} \right] \\
 &= \frac{nG}{N} \left[ \frac{\cancel{N(N-G)} - \cancel{n(N-G)}}{\cancel{N^2-NG} - nN + nG} \right] = \frac{(N-n)}{(N-G)} \\
 &= \frac{nG}{N} \cdot \frac{(N-n)(N-G)}{N \cdot (N-1)} = n \left( \frac{G}{N} \right) \left( \frac{N-G}{N} \right) \left( \frac{N-n}{N-1} \right)
 \end{aligned}$$

## The finite population correction & the accuracy of SRS



$$n \cdot p \cdot q \cdot \left[ \frac{N-n}{N-1} \right]$$

FINITE  
POPN  
CORRECTION

$$fpc = \sqrt{\frac{N-n}{N-1}}$$

Note that  $fpc \leq 1$

So  $SD(HG) \leq SD(Bin)$

SD is less when we don't have independence

sum of SRS = sum of draws WITHOUT replacement

In general we have that the :

$$\text{SD of sum of an SRS} = \text{SD of sum WITH repl.} \times fpc$$

$$\text{SD of HG} = \text{SD of binomial} \times \sqrt{\frac{N-n}{N-1}}$$

$$\text{Var(HG)} = \text{Var(Bin)} \times \left( \frac{N-n}{N-1} \right)$$

3/28/22

Exercise : Let  $N = 10^6$   
 $n = 1000$  what is  $\sqrt{\frac{N-n}{N-1}} = 0.9995$ ,  $n > 1$

## Accuracy of samples

"Try this at home": Plug in different values of  $N$ ,  $n$  & see what is fpc.

Simple random samples of the same size of 625 people are taken in Berkeley (population: 121,485) and Los Angeles (population: 4 million). True or false, and explain your choice: The results from the Los Angeles poll will be substantially more accurate than those for Berkeley.

False

Try fpc with  $N=121485$   
 $n=625$

$$\begin{array}{l} N = 4 \times 10^6 \\ n = 625 \end{array}$$

Accuracy depends on SD & fpc

If  $fpc \approx 1$ , then only depends on SD.

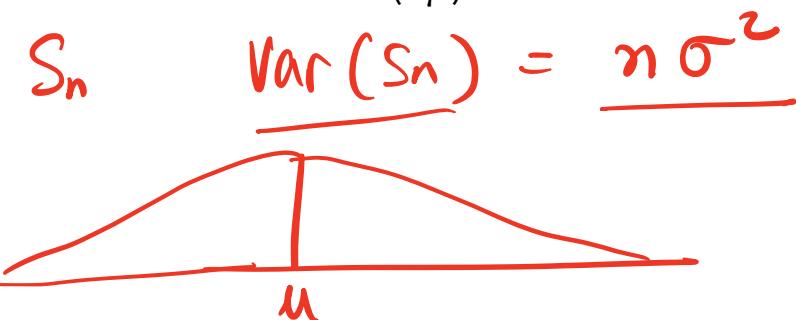
## Example (adapted from *Statistics*, by Freedman, Pisani, and Purves)

A survey organization wants to take an SRS in order to estimate the percentage of people who watched the 2022 Oscars. To keep costs down, they want to take as small a sample as possible, but their client will only tolerate a random error of 1 percentage point or so in the estimate. Should they use a sample size of 100, 2500, or 10000? The population is very large and the fpc is about 1.

### Exercise

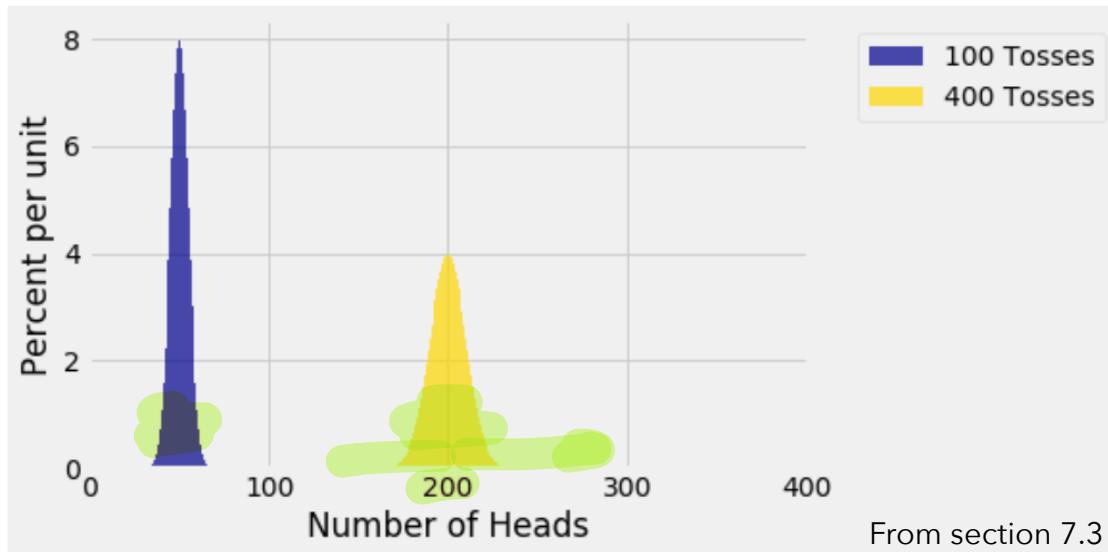
## Law of Averages

- Essentially a statement that you are already familiar with: If you toss a fair coin many times, roughly half the tosses will land heads.
- We are going to consider sample sums and sample means of iid random variables  $X_1, X_2, \dots, X_n$  where the mean of each  $X_k$  is  $\mu$  and the variance of each  $X_k$  is  $\sigma^2$ .  
$$X_k \sim \mu, \sigma^2$$
- Recall the **sample sum**  $S_n = X_1 + X_2 + \dots + X_n$ , with  $E(S_n) = n\mu$ ,  $Var(S_n) = n\sigma^2$ ,  $SD(X_n) = \sqrt{n}\sigma$
- We see here, as we take more and more draws, the variability of the sum keeps increasing, which means the values get more and more dispersed around the mean ( $n\mu$ ).



## Coin tosses

- Consider a fair coin, toss it 100 times & 400 times, count the number of H. Expect in first case, roughly 50 H, and in second, roughly 200 H.
- So do you think chance of 50 H in 100 tosses and 200 H in 400 tosses should be the same?



3/28/22

Total area = 1 in both cases  
so spread about the mean in 400 tosses is greater.

Example: Coin toss

$$S_{100} \sim \text{Bin}(100, \frac{1}{2})$$
$$S_{400} \sim \text{Bin}(400, \frac{1}{2})$$

- $SD(S_{100}) = \sqrt{npq} = \sqrt{100 \cdot \frac{1}{2} \cdot \frac{1}{2}} = 10 \cdot \frac{1}{2} = 5$
- $SD(S_{400}) = \sqrt{npq} = \sqrt{400 \cdot \frac{1}{2} \cdot \frac{1}{2}} = 20 \cdot \frac{1}{2} = 10$

- $P(200 \text{ H in } 400 \text{ tosses})$

$$= P(S_{400} = 200) = \binom{400}{200} \left(\frac{1}{2}\right)^{400}$$
$$\approx \underline{\underline{0.04}}$$

- $P(50 \text{ H in } 100 \text{ tosses})$

$$P(S_{100} = 50) = \binom{100}{50} \left(\frac{1}{2}\right)^{100} \approx 0.08$$

## Law of Averages for a fair coin

SD increase  
error increases

- Notice that as the number of tosses of a fair coin increases, the observed error (number of heads - half the number of tosses) increases. This is governed by the standard error.
- The percentage of heads observed comes very close to 50%
- Law of averages: The long run proportion of heads is very close to 50%.

Sample mean  $\longrightarrow$  prob of H as n gets large.

## Sample sum, sample average, and the square root law

$$S_n = X_1 + X_2 + \dots + X_n$$

$$A_n = \text{sample mean} = \frac{S_n}{n}$$

- Let  $A_n = \frac{S_n}{n}$ , so  $A_n$  is the average of the sample (or sample mean).
- If the  $X_k$  are indicators, then  $A_n$  is a proportion (proportion of successes)

$$\text{Note that } E(A_n) = \mu \text{ and } SD(A_n) = ?? \quad SD\left(\frac{S_n}{n}\right) = \sqrt{n} \cdot \frac{\sigma}{n} = \frac{\sigma}{\sqrt{n}}$$

as  $n$  gets large  $SD(S_n)$  gets large,  $SD(A_n)$  gets small.

- The square root law:** the accuracy of an estimator is measured by its SD, the **smaller** the SD, the **more accurate** the estimator, but if you multiply the sample size by a factor, the accuracy only goes up by the **square root** of the factor.
- In our earlier example, we \_\_\_\_\_ the accuracy by quadrupling the size.