

Stat 68 Lec 24

Warm Up 2:00-2:10

Let I_A be the indicator of the event A . Then the distribution of I_A is given by

value	0	1
probability	$1 - P(A)$	$P(A)$

a) Find $E(I_A) = 0 \cdot (1 - P(A)) + 1 \cdot P(A) = P(A)$ indicator for an event

b) Find $E(I_A^2) = 0^2 \cdot (1 - P(A)) + 1^2 \cdot P(A) = P(A)$

c) What is the distribution table for $I_A I_B$

d) Find $E(I_A I_B) = P(AB)$

$\begin{matrix} 0 \\ 1 \\ \hline I_{AB} \end{matrix}$

$$\begin{array}{c|c} 0 & 1 \\ \hline 1-P(AB) & P(AB) \end{array}$$

Last time

Sec 7.1 Sums of independent random variables

If X_1, X_2 are independent RVs

$$\text{Var}(X_1 + X_2) = \text{Var}(X_1) + \text{Var}(X_2)$$

$X \sim \text{Bernoulli}(\rho)$

$$SD(X) = \sqrt{\rho q}$$

$X \sim \text{Binomial}(n, \rho)$

$$SD(X) = \sqrt{n \rho q}$$

$X \sim \text{Poisson}(\mu)$

$$SD(X) = \sqrt{\mu}$$

$X \sim \text{Geometric}(\rho)$

$$SD(X) = \frac{\sqrt{q}}{\rho}$$

Today

- (1) Sec 7.2 Sampling without replacement
- (2) Sec 7.3 The Law of Averages

(3) Sec 7.2 Sampling without replacement

Whereas a binomial RV is a sum of independent indicators, a hypergeometric RV is a sum of dependent indicators. This makes it more complicated to find the formula for SD (hypergeometric).

Squares and products of indicators

Let I_A be the indicator of the event A . Then the distribution of I_A is given by

value	0	1
probability	$1 - P(A)$	$P(A)$

We saw in Warmup:

if I_B is the indicator for event B

$$\text{then } I_A I_B = I_{AB} = \begin{cases} 1 & \text{if } A \text{ and } B \text{ true} \\ 0 & \text{else} \end{cases}$$

$\overbrace{\quad\quad\quad}^{\text{P}(AB)}$

and $E(I_A I_B) = \text{P}(AB)$

SD of Hypergeometric

Let $X \sim HG(N, G, n)$

$X = \# \text{ good in a SRS from pop size } N \text{ w/}$

$G \text{ good, } P = \frac{G}{N}$

$X = I_1 + \dots + I_n \quad \text{where } I_2 = \begin{cases} 1 & \text{if 2nd draw good} \\ 0 & \text{else.} \end{cases}$

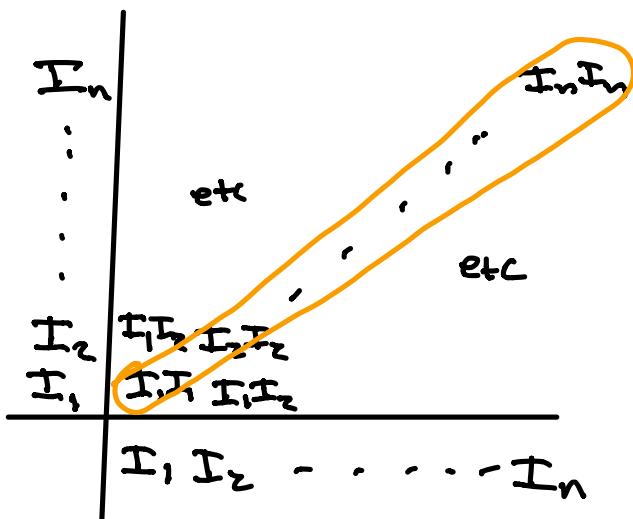
Recall

$$E(X) = n \left(\frac{G}{N} \right)$$

Find $\text{Var}(X)$

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

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



$$(I_1 + \dots + I_n)^2 = \text{sum of diag terms} + \text{sum of off diag terms},$$

$$= \sum_{j=1}^n I_j^2 + \sum_{j=1}^n \sum_{\substack{k=1 \\ j \neq k}}^n I_j I_k$$

so

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

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

$$= \underbrace{nE(I_1^2) + n(n-1)E(I_1 I_2)}$$

$$E(I_1^2) = \frac{6}{N}, \quad E(I_1 I_2) = \frac{6}{N} \cdot \frac{6-1}{N-1}$$

← $P(I_1=1)$ ← $P(I_1=1 \text{ and } I_2=1)$
 $P(I_1=1)$ $P(I_2=1 | I_1=1)$

so,

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

Since,

$$Var(X) = E(X^2) - (E(X))^2$$

$$Var(X) = \underbrace{n \frac{G}{N} + n(n-1) \frac{G}{N} \cdot \frac{G-1}{N-1}}_{E(X^2)} - \underbrace{\left(n \frac{G}{N}\right)^2}_{E(X)^2}$$

 boring algebra
(*) (see appendix to lecture)

$$Var(X) = n \frac{G}{N} \cdot \frac{N-G}{N} \cdot \frac{N-n}{N-1}$$

$$= npq \cdot \left[\frac{N-n}{N-1} \right]$$

$$\Rightarrow SD(X) = \sqrt{npq} \cdot \sqrt{\frac{N-n}{N-1}}$$

↑
SD of Binomial(n, p)

called Finite Population Correction (FPC)

Ex Draw 5 cards from a deck.

$X = \#$ of hearts in your hand,

$$\text{Find } E(X) = 5 \frac{13}{52} = 1.25$$

$$SD(X) = \sqrt{\frac{5 \frac{13}{52} \cdot \frac{39}{52}}{52}} \cdot \sqrt{\frac{52-5}{52-1}} = .93$$

Size of the FPC

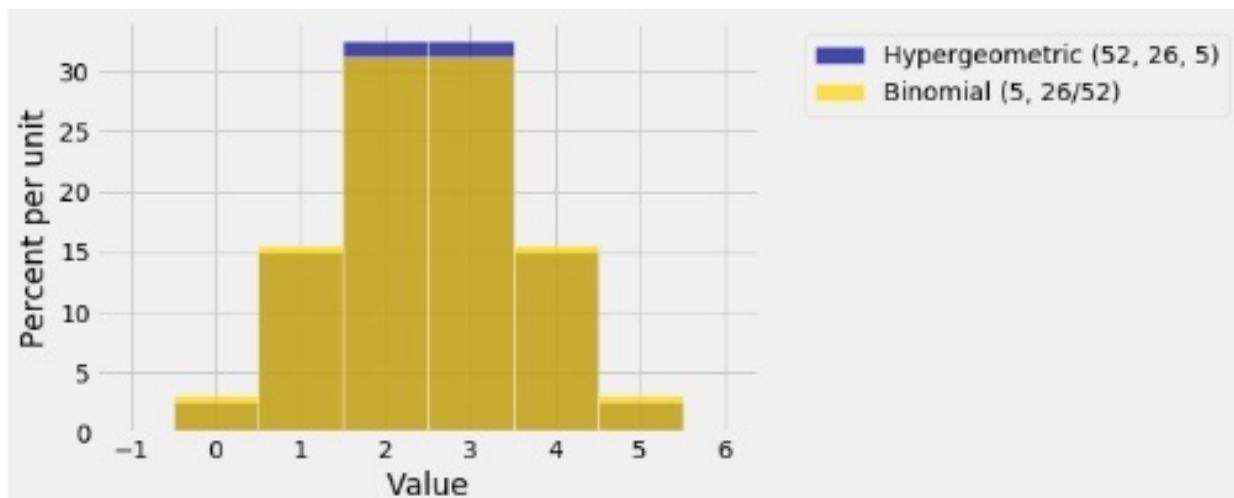
$$FPC = \sqrt{\frac{N-n}{N-1}} \leq 1$$

if $n \ll N$ (less than 1%)
 then like drawing with replacement
 and $FPC=1$

$$SD(HG) = SD(\text{Binomial}) \cdot FPC$$

$$\text{says } SD(HG) \leq SD(\text{Binomial})$$

Ex



Polls are made without replacement because the count of # democrats is more accurate (i.e. smaller SD than polling with replacement).



Tinyurl.com/march16-pt1

New Mexico has a population of 1M and California a population of 40M. The two states have the same proportion of Democrats.

SD
FPC $\propto 1$

A random sample of size 0.01% of the population is taken. The SD for the number of democrats in the sample is:

- a roughly the same in both states
- b larger in California
- c larger in New Mexico

$$.01\% \text{ of } 40M = .0001 (40M) = 4000$$

$$.01\% \text{ of } 1M = .0001 (1M) = 100$$

$$SD(\text{Dem CA}) = \sqrt{4000pq} \leftarrow \text{larger}$$

$$SD(\text{Dem NM}) = \sqrt{100pq}$$

If $n \ll N$ think of N as ∞
we see from this that SD depends on p ,
on n but not on N .

Follow up question: suppose in each state a random sample of 500 is taken. The SD of the number of democrats in the poll is:

- a) roughly the same in both states
- b) larger in CA
- c) larger in NM

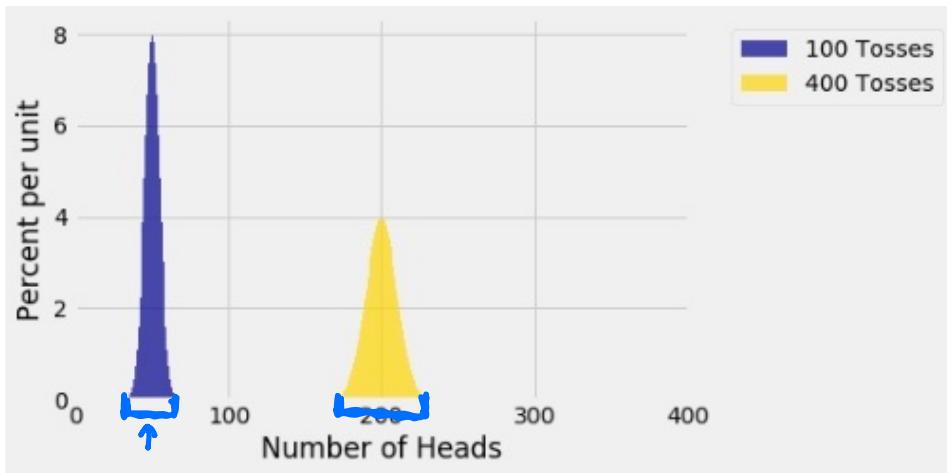
② The Law of Averages

The sample sum

You toss a fair coin n times
 Let $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} \text{Bernoulli}\left(\frac{1}{2}\right)$
 $S_n = X_1 + X_2 + \dots + X_n$ is the sample sum

What would you put your money on :

$$\begin{aligned} \text{a) } S_{100} &= 50 \quad \text{or} \quad \left(\frac{1}{2}\right)^{50} \cdot \left(\frac{1}{2}\right)^{50} \\ \text{b) } S_{400} &= 200 ? \quad \text{larger} \\ P(S_{100}=50) &= \binom{100}{50} \left(\frac{1}{2}\right)^{100} = 0.08 \\ P(S_{400}=200) &= \binom{400}{200} \left(\frac{1}{2}\right)^{400} = 1.04 \end{aligned}$$



Both histograms have area 1,
 Yellow histogram is more spread out and
 hence $P(S_{400}=200) < P(S_{100}=50)$.

We can see this mathematically,

$$\begin{aligned}\text{Var}(S_n) &= \text{Var}(X_1 + \dots + X_n) \\ &= \text{Var}(X_1) + \dots + \text{Var}(X_n) \\ &= n \text{Var}(X_1) \quad \text{Var(Bernoulli } (\frac{1}{2}) \text{)} = \frac{1}{2} \cdot \frac{1}{2} \\ &= \boxed{n\sigma^2} = n \frac{1}{4} \quad = \boxed{\frac{1}{4}} \\ \Rightarrow \text{SD}(S_n) &= \sqrt{n} \sigma = \boxed{\sqrt{n} \frac{1}{2}}\end{aligned}$$

The sample average

You toss a fair coin n times

Let $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} \text{Bernoulli } (\frac{1}{2})$

$S_n = X_1 + X_2 + \dots + X_n$ is the sample sum

let $A_n = \frac{S_n}{n}$ be the sample average

Where do you put your money?

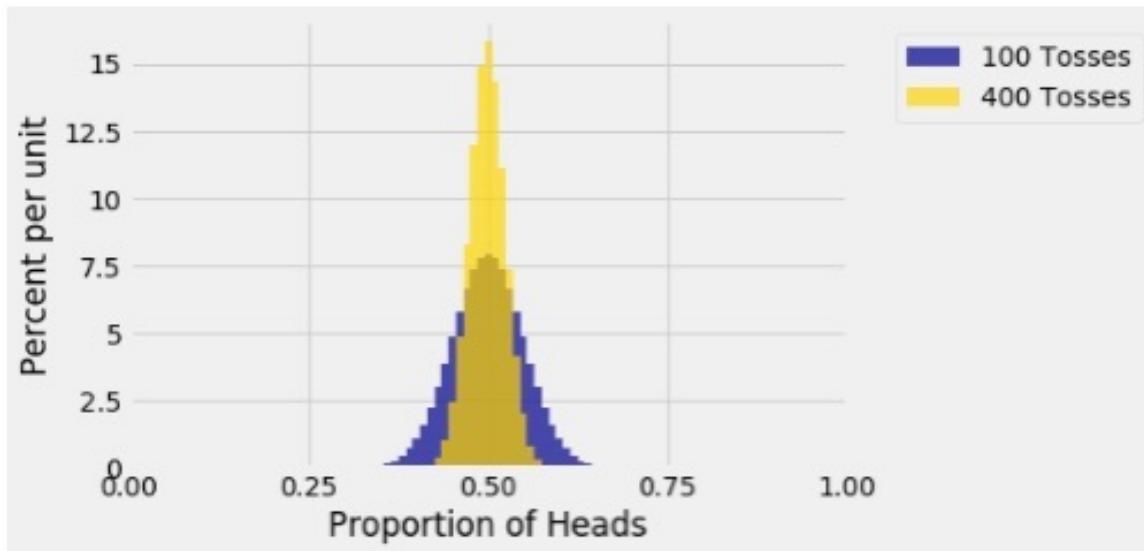
a) $A_{100} = 50$ or

b) $A_{400} = 200 ?$

$$\begin{aligned}\text{SD}(A_n) &= \text{SD}\left(\frac{S_n}{n}\right) = \frac{1}{n} \text{SD}(S_n) \quad \text{"} \sqrt{n}\sigma \\ &= \frac{1}{n} \sigma = \boxed{\frac{\sigma}{\sqrt{n}}}\end{aligned}$$

$$\text{SD}(A_{100}) = \frac{1}{10} = \frac{\sigma}{\sqrt{10}}$$

$$\text{SD}(A_{400}) = \frac{1}{20} = \frac{\sigma}{\sqrt{20}} \leftarrow \text{less error}$$



In general, the larger the sample size n , the more likely it is that the sample average A_n will be close to the population average μ .

This is the "law of averages".

* Appendix

Boring algebra details

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

Pull out $\frac{nG}{N}$

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

common denom.

$$= \frac{nG}{N} \frac{(N-1)N + N(N-1)(G-1) - nG(N-1)}{N(N-1)}$$

Factor out

$$= \frac{nG}{N} \frac{N^2 - N + Ng - N^2 - nN + N - nGN + nG}{N(N-1)}$$

reverse FOIL

$$= \frac{nG}{N} \frac{(N-G)(N-n)}{N(N-1)}$$

$$= \boxed{n \frac{G}{N} \cdot \frac{N-G}{N} \cdot \frac{N-n}{N-1}}$$