

* Announcement

① HW 7 due today (11:59 PM PT)

② HW 8 ~ 10/26

③ Quiz 6. Ch6

After today's lecture

→ HW 8 Q1, Q2

Q3, Q4

STAT 88: Lecture 23

Contents

Section 7.2: Sampling Without Replacement

Warm up: (Exercise 7.4.5) The number of typos on the cover page of an exam has a distribution given by

value	0	1
probability	0.8	0.2

The number of misprints in the rest of the exam has the Poisson(3) distribution, independently of the cover page. Find the expectation and SD of the total number of misprints on the exam.

$$\begin{aligned} X_1 &= \# \text{ of typos cover} \sim \text{Bernoulli}(0.2) \\ X_2 &= \# \text{ of misprints exam.} \sim \text{Pois}(3) \end{aligned} \quad \left. \vphantom{\begin{aligned} X_1 &= \# \text{ of typos cover} \sim \text{Bernoulli}(0.2) \\ X_2 &= \# \text{ of misprints exam.} \sim \text{Pois}(3) \end{aligned}} \right\} \text{ independent.}$$

$$X = X_1 + X_2$$

$$\begin{aligned} E(X) &= E(X_1 + X_2) = E(X_1) + E(X_2) \\ &= 0.2 + 3 \\ &= 3.2 \end{aligned}$$

$$\begin{aligned} \text{Var}(X) &= \text{Var}(X_1 + X_2) = \text{Var}(X_1) + \text{Var}(X_2) \\ &= 0.2(0.8) + 3 \\ &= 3.16 \end{aligned}$$

$$\text{SD}(X) = \sqrt{3.16}$$

Last time

Sums of independent random variables: If X and Y are independent random variables, then

$$\text{Var}(X + Y) = \text{Var}(X) + \text{Var}(Y).$$

If $X \sim \text{Bernoulli}(p)$,

$$\text{SD}(X) = \sqrt{p(1-p)}.$$

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

"
sum of n indep. Ber RVs

$$\text{SD}(X) = \sqrt{np(1-p)}.$$

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

$\downarrow \begin{matrix} n \rightarrow \infty \\ p \rightarrow 0 \\ np \rightarrow \mu \end{matrix}$

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

If $X \sim \text{Geom}(p)$,

$$\text{SD}(X) = \frac{\sqrt{1-p}}{p}.$$

Example: (Exercise 7.4.10) A non-negative integer valued random variable has expectation 50 and SD 10. Could the random variable have a binomial distribution?

$$E(X) = 50$$

$$SD(X) = 10.$$

Let's assume $X \sim \text{Binom}(n, p)$

$$E(X) = np$$

$$SD(X) = \sqrt{np(1-p)}$$

$$\rightarrow \begin{cases} np = 50 \\ np(1-p) = 100 \end{cases}$$

$$\rightarrow 50(1-p) = 100$$

$$\rightarrow 1-p = 2$$

$$\rightarrow p = -1 \quad \text{Not possible.}$$

$\Rightarrow X$ cannot have Binomial

7.2. Sampling without Replacement

The draws in a SRS are **not independent** of each other and this makes computing the SD of the hypergeometric distribution more complicated than for binomial distribution.

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	$1 - P(A)$	$P(A)$

$$I_A = \begin{cases} 1 & \text{if event } A \text{ occurs} \\ 0 & \text{else} \end{cases}$$

$\sim \text{Bernoulli}(P(A))$

$$I_A^2 = I_A$$

Find $E(I_A)$ and $E(I_A^2)$.

$$\begin{aligned} & \text{" } P(A) \\ & \text{" } E(I_A) \\ & \text{" } P(A) \end{aligned}$$

Now let I_B be the indicator for event B . What values does $I_A I_B$ take?

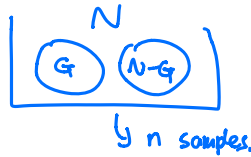
$$I_B \sim \text{Bernoulli}(P(B))$$

$$\begin{aligned} & \text{A and B occurs} \\ & \text{" } \\ & \text{A occurs and B occurs} \\ & \text{" } \\ & = \begin{cases} 1 & \text{if } I_A=1 \text{ and } I_B=1 \\ 0 & \text{else.} \end{cases} \end{aligned}$$

Find $E(I_A I_B)$.

$$\begin{aligned} & \text{" } \\ & E(I_{A \cap B}) \\ & \text{" } \\ & P(A \cap B) \end{aligned}$$

$$= I_{A \cap B}$$



SD of Hypergeometric

Let $X \sim \text{HG}(N, G, n)$. So X measures the number of good elements in a simple random sample of size n drawn from a population of N elements of which G are good. Then we can write

$$X = I_1 + \cdots + I_n,$$

where

$$I_j = \begin{cases} 1 & \text{if } j\text{th draw is good} \\ 0 & \text{else} \end{cases} \quad p = \frac{G}{N} \quad (\text{Symmetry})$$

Recall $E(X) = n \frac{G}{N}$. We want to find $\text{Var}(X) = E(X^2) - (EX)^2$.

We start with the following equation:

$$X^2 = (I_1 + \cdots + I_n)^2 = \sum_{j=1}^n I_j^2 + \sum_{j \neq k} I_j I_k.$$

n terms n^2 - n = n(n-1) terms

Taking expectation on both side,

$$\begin{aligned} E(X^2) &= E\left(\sum_{j=1}^n I_j^2\right) + E\left(\sum_{j \neq k} I_j I_k\right) \\ &\stackrel{\text{Additivity}}{=} \sum_{j=1}^n E(I_j^2) + \sum_{j \neq k} E(I_j I_k) \\ &\stackrel{\text{Symmetry}}{=} nE(I_1^2) + n(n-1)E(I_1 I_2). \end{aligned}$$

$\frac{G}{N} \cdot \frac{G-1}{N-1}$
 \parallel
 $P(I_1=1) P(I_2=1 | I_1=1)$
 \parallel
 $P(I_1=1 \text{ and } I_2=1)$
 \parallel
 $P(I_1 I_2=1)$
 \parallel

We know from our calculation above that $E(I_1^2) = \frac{G}{N}$ and $E(I_1 I_2) = \frac{G}{N} \cdot \frac{G-1}{N-1}$. So,

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

Since $\text{Var}(X) = E(X^2) - (EX)^2$, it follows that

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

E(X^2) (EX)^2

After some boring algebra (see the last page of the note), we can simplify it to

$$\text{Var}(X) = n \frac{G}{N} \cdot \frac{N-G}{N} \cdot \frac{N-n}{N-1}, \text{ and } \text{SD}(X) = \sqrt{n \frac{G}{N} \cdot \frac{N-G}{N} \cdot \frac{N-n}{N-1}} = \sqrt{np(1-p)} \sqrt{\frac{N-n}{N-1}}$$

SD of Binomial(n, p)

Example: Draw 5 cards from a deck. Let X be the number of hearts in your hand. Find $E(X)$ and $\text{SD}(X)$.

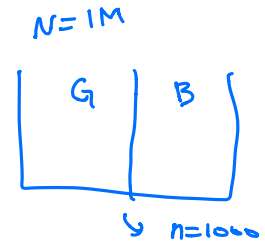
$$X \sim \text{HG}(52, 13, 5)$$

$$E(X) = 5 \cdot \frac{1}{4}, \quad \text{Var}(X) = 5 \cdot \frac{1}{4} \cdot \frac{3}{4} \cdot \frac{52-13}{52-1}$$

The Size of the FPC

Finite population correction or FPC is given by

$$0 < fpc = \sqrt{\frac{N-n}{N-1}} < 1$$



We saw that

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

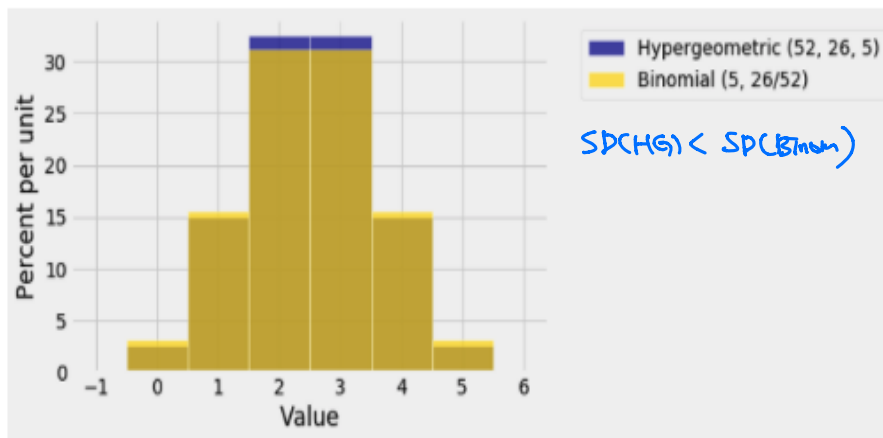
so $SD(HG) < SD(\text{Binomial})$.

($n > 1$: $fpc < 1$)

(n much smaller than N)

If $n \ll N$, $fpc \approx 1 < 1$ ($fpc \approx 0.99$)

and it's like drawing w/ replacement



$SD(HG) < SD(\text{Binom})$

Sampling with and without replacement are essentially the same when the sample size is small relative to the population size.

Read “The Accuracy of Simple Random Samples” in Ch 7.2 of the textbook.

$$p = \frac{G}{N} \text{ same between two states.}$$

Example: New Mexico has a population of 1M and California a population of 40M. The two states have the same proportion of Democrats. A random sample of size 0.01% of the population is taken. The SD for the number of democrats in the sample

is:

$$fpc \approx 1$$

$$\uparrow$$

$$n = \frac{0.01}{100} \cdot N$$

1. roughly the same in both states

$$= \frac{1}{10000} \cdot N$$

2. larger in California

$$n = \frac{0.05}{100} \cdot N$$

3. larger in New Mexico

$$= \frac{1}{2000} \cdot N$$

$X_C = \# \text{ democrats in the sample CA} \sim \text{HG}(40M, G_C, 4000)$

$X_{NM} = \# \text{ democrats in the sample New Mexico} \sim \text{HG}(1M, G_{NM}, 100)$

$$SD(X_C) \approx \sqrt{n \cdot p \cdot (1-p)}$$

$$= \sqrt{4000 \cdot p \cdot (1-p)}$$

$$SD(X_{NM}) \approx \sqrt{n \cdot p \cdot (1-p)}$$

$$= \sqrt{100 \cdot p \cdot (1-p)}$$

Example: 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:

1. roughly the same in both states
2. larger in California
3. larger in New Mexico

$$SD \approx \sqrt{n p(1-p)}$$

~ ~ ~ same .
same

Algebra details:

$$\begin{aligned}
 \text{Var}(X) &= n \frac{G}{N} + n(n-1) \frac{G}{N} \cdot \frac{G-1}{N-1} - \left(n \frac{G}{N} \right)^2 \\
 &\stackrel{\text{Pull out } n \frac{G}{N}}{=} n \frac{G}{N} \left(1 + (n-1) \cdot \frac{G-1}{N-1} - n \frac{G}{N} \right) \\
 &\stackrel{\text{Common denominator}}{=} n \frac{G}{N} \cdot \frac{N(N-1) + N(n-1)(G-1) - nG(N-1)}{N(N-1)} \\
 &= n \frac{G}{N} \cdot \frac{N^2 - \cancel{N} + n\cancel{N}G - nN - \cancel{NG} + \cancel{N} - n\cancel{N}G + nG}{N(N-1)} \\
 &= n \frac{G}{N} \cdot \frac{N^2 - nN - NG + nG}{N(N-1)} \\
 &= n \frac{G}{N} \cdot \frac{(N-G)(N-n)}{N(N-1)} \\
 &= n \frac{G}{N} \cdot \frac{N-G}{N} \cdot \frac{N-n}{N-1}.
 \end{aligned}$$