

Quiz Day !!  
7pm (not Berkeley time)

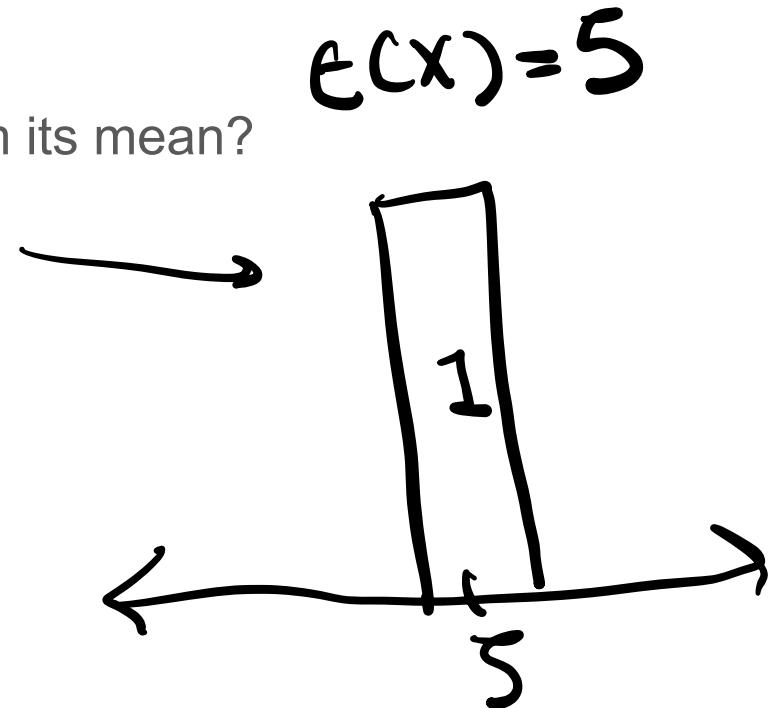
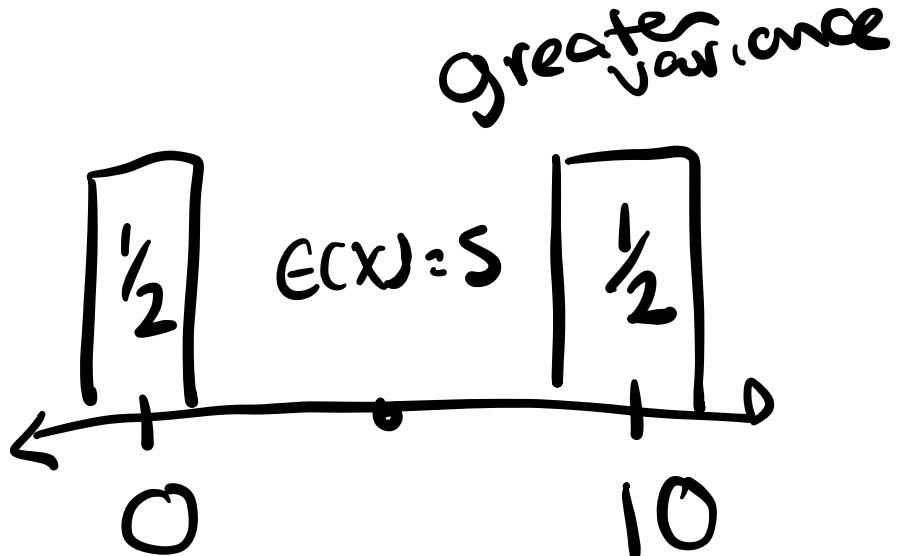
# Lecture 5C:

# Variance, Covariance & Correlation

UC Berkeley CS70  
Summer 2023  
Nikki Suzani

# What is variance?

How much does a random variable deviate from its mean?



## Variance (formally)

$$\text{value } E[cX] = cE[X]$$

$$E[c] = c$$

**Definition 16.1** (Variance). For a r.v.  $X$  with expectation  $E[X] = \mu$ , the variance of  $X$  is defined to be

$$\text{Var}(X) = E[(X - \mu)^2].$$

mean

Squared deviation from the mean

The square root  $\sigma(X) := \sqrt{\text{Var}(X)}$  is called the standard deviation of  $X$ .

$$\begin{aligned} E[(X - E[X])^2] &= E[X^2 - 2XE[X] + E[X]^2] \\ &= E[X^2] - \underbrace{E[2XE[X]]}_{E[2E[X]]} + \overbrace{E[E[X]^2]} \\ &= E[X^2] - 2E[X]E[X] + \overbrace{E[X]}^2 \\ &= E[X^2] - (E[X])^2 \end{aligned}$$

# Variance of a Coin Toss

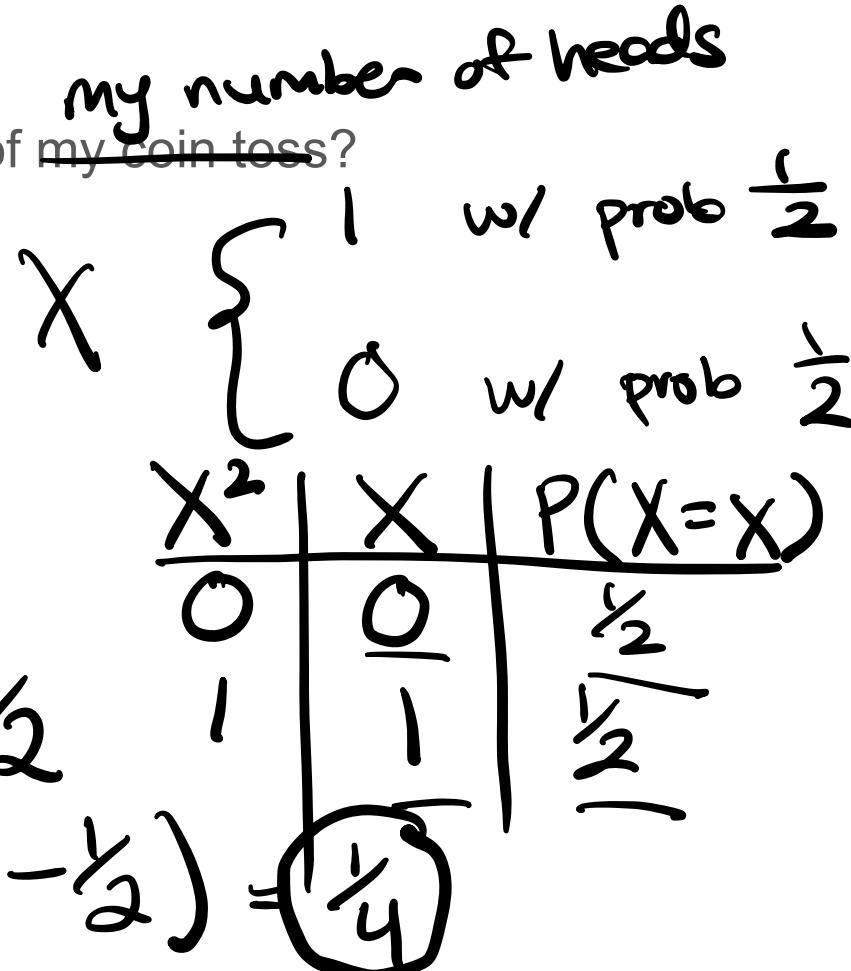
I flip a fair coin. What is the variance of my coin toss?

$X$ : # of heads  
on my flip

$$E[X^2] - E[X]^2$$

$$E[X^2] = 0 \cdot \frac{1}{2} + 1 \cdot \frac{1}{2} = \frac{1}{2}$$

$$\frac{1}{2} - \left(\frac{1}{2}\right)^2 = \frac{1}{2}(1 - \frac{1}{2}) = \frac{1}{4}$$



# Variance of a Dice Roll

I roll a single six-sided dice. What is the variance of my dice roll?

$X$ : Value of  
my dice roll

$$\begin{aligned} \text{Var}(X) &= E(X^2) - E(X)^2 \\ &= \left( \frac{1+4+9+16+25+36}{6} \right) - \left( \frac{1+2+3+4+5+6}{6} \right)^2 \\ &= \frac{35}{12} \end{aligned}$$

$X^2$	$X$	$P(X=x)$
1	1	$\frac{1}{6}$
4	2	$\frac{1}{6}$
9	3	$\frac{1}{6}$
16	4	$\frac{1}{6}$
25	5	$\frac{1}{6}$
36	6	$\frac{1}{6}$

$$\mathbb{E}(X) = \sum_{k=0}^{\infty} k P(X=k)$$

$$\mathbb{E}(X^2) = \sum_{k=0}^{\infty} k^2 P(X=k)$$

## Variance of a Uniform Random Variable

$X \sim \text{Uniform}(n)$

Let  $X$  be a random variable that takes on the values 1, ..., n with the same probability (1/n). What is Var(X)?

$$\begin{aligned}\text{Var}(X) &= \mathbb{E}(X^2) - \mathbb{E}(X)^2 \\ &= \sum_{k=0}^{\infty} k^2 P(X=k) - \left( \sum_{k=0}^{\infty} k P(X=k) \right)^2 \\ &= \frac{1}{n} \sum_{k=0}^{\infty} k^2 - \left( \frac{1}{n} \sum_{k=0}^{\infty} k \right)^2\end{aligned}$$

## Variance of a Uniform Random Variable (cont.)

$$\frac{1}{n} \cdot \frac{n(n+1)(2n+1)}{6} - \left( \frac{1}{n} \cdot \frac{n(n+1)}{2} \right)^2$$
$$= \frac{n^2 - 1}{12}$$

## Example: Variance of a Dice Roll (again)

Roll a fair six-sided dice. What is the variance?

$\{1, \dots, 6\}$   
each w/ prob  $\frac{1}{6}$

$X \sim \text{Uniform}(6)$

$$\text{Var}(X) = \frac{n^2 - 1}{12} = \frac{36 - 1}{12} = \frac{35}{12}$$

$$SD = \sqrt{Var}$$

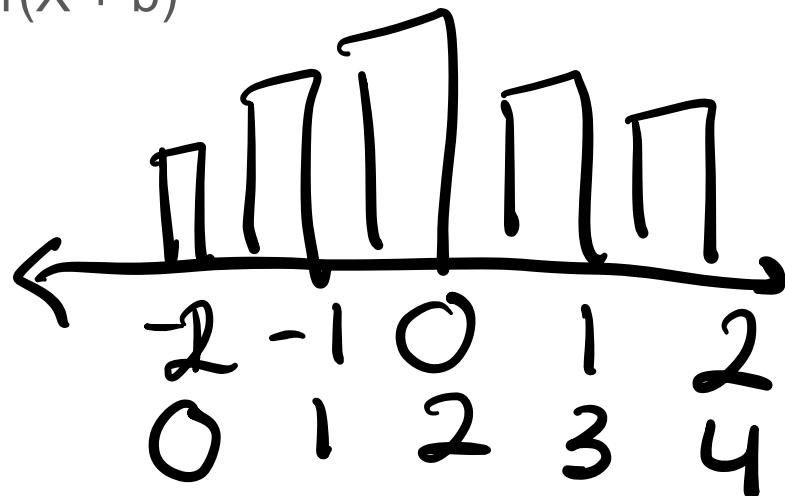
$$E[cX] = cE[X]$$

Properties of Variance

$$\begin{aligned} \text{Var}(cX) &= E[(cX)^2] - (E[cX])^2 \\ \downarrow \\ SD(cX) &= E[c^2 X^2] - (cE[X])^2 \\ &= c^2 E[X^2] - c^2 E[X]^2 \\ &= c^2 (E[X^2] - E[X]^2) \\ &= c^2 (\text{Var}(X)) \rightarrow cSD(X) \end{aligned}$$

## Properties of Variance

$$\text{Var}(X + b)$$



$$\text{Var}(X + 2)$$

$$E((X+b)^2)$$

$$- E(X+b)^2$$

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

$$= \text{Var}(X)$$

# Properties of Variance

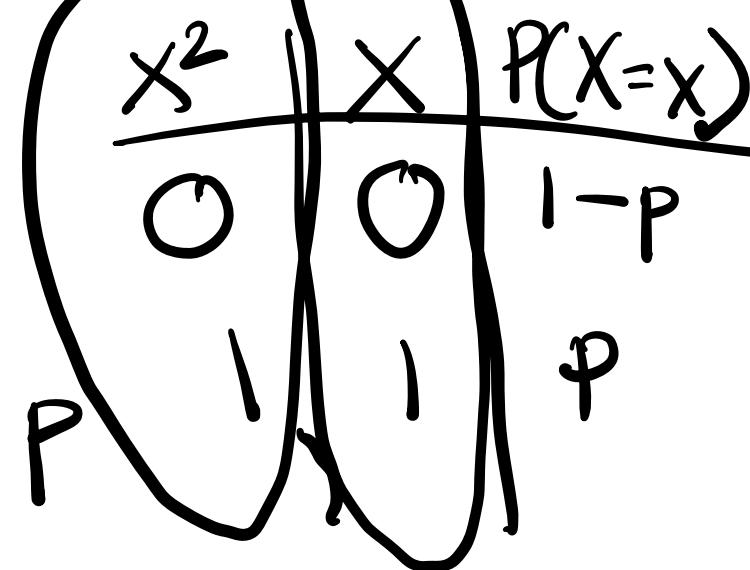
$$\text{Var}(X+b)$$

$$0 \cdot 1 - p + 1 \cdot p = p$$

Variance of the Bernoulli Distribution

Let  $X \sim \text{Bernoulli}(p)$ . What is  $\text{Var}(X)$ ?

$$X : \begin{cases} 1 & \text{w/ probability } p \\ 0 & \text{w/ probability } 1-p \end{cases}$$



$$\begin{aligned} \text{Var}(X) &= E(X^2) - E(X)^2 = \boxed{p(1-p)} \\ &= E(X) - E(X)^2 \\ &\equiv E(X)(1 - E(X)) \end{aligned}$$

# Making Serves

Nikki plays tennis (badly). Her serve ends up in the “legal zone” 20% of the time.  
What is the variance of the ~~number of times~~ <sup>one serve</sup> her serve is legal?

serve one time

$$X: \begin{cases} 1 & \text{w/ prob 0.2} \\ 0 & \text{w/ prob 0.8} \end{cases} \quad X \sim \text{Bernoulli}(0.2)$$

$$\begin{aligned} \text{Var}(X) &= p(1-p) = 0.2(1-0.2) \\ &= 0.2 \cdot 0.8 = 0.16 \end{aligned}$$

# Variance of Independent Random Variables

If X and Y are independent random variables,  $\text{Var}(X + Y) = \text{Var}(X) + \text{Var}(Y)$

*possibly true when not independent*

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

all independent trials

## Variance of the Binomial Distribution

Let  $X \sim \text{Binomial}(n, p)$ . What is  $\text{Var}(X)$ ?

$X_i \sim \text{Bernoulli}(p)$

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

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

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

$$= p(1-p) + p(1-p) + \dots + p(1-p)$$
$$= [np(1-p)]$$

Example: Variance of 10 coin tosses

$$P(1-P) \\ (1-P)P$$

You toss a coin ten times that flips heads with probability  $\frac{2}{3}$ . What is the variance of the number of heads?

$$X: \begin{matrix} \# \text{ of} \\ \text{heads} \end{matrix} \quad X = X_1 + X_2 + \dots + X_{10}$$

$$X_i \sim \begin{cases} 1 & \text{if got a heads} \\ 0 & \text{otherwise} \end{cases}$$

$$X_i \sim \text{Bernoulli}\left(\frac{2}{3}\right)$$

$$X \sim \text{Binomial}(10, \frac{2}{3}) = 10 \cdot \frac{2}{3} \left(1 - \frac{2}{3}\right) \\ = \frac{20}{9}$$

# Example: More tennis

Carlos Alcaraz wins matches against Novak Djokovic 1/3 of the time. If they play five times next year, what is the variance of the number of matches Alcaraz wins?

Year	Event	Surface	RND	Winner	Result
2023	Wimbledon Great Britain	Outdoor Grass	F	Carlos Alcaraz	16 76 <sup>6</sup> 61 36 64
2023	Roland Garros France	Outdoor Clay	SF	Novak Djokovic	63 57 61 61
2022	ATP Masters 1000 Madrid Spain	Outdoor Clay	SF	Carlos Alcaraz	67 <sup>5</sup> 75 76 <sup>5</sup>

$$X = X_1 + X_2 + \dots + X_5$$

$$X_i \sim \text{Bernoulli}\left(\frac{1}{3}\right)$$

$$X_i \sim \text{Binomial}(5, \frac{1}{3})$$

$$\text{Var}(X) = 5 \cdot \frac{1}{3} \cdot \left(1 - \frac{1}{3}\right)$$

$$= \frac{1}{q}$$

## What is covariance?

$\text{Cov}(X, Y)$  gives us information about **how related** X and Y are to each other.

Covariance is positive  
when X is big, Y is big

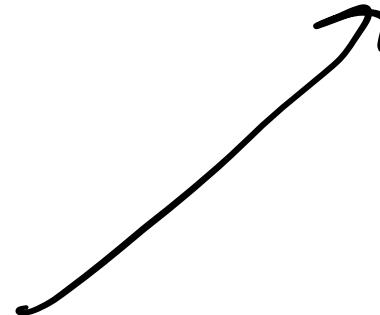
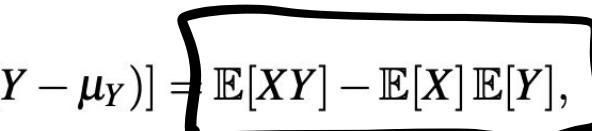
Covariance is negative  
when X is small, Y is big

# Covariance (Formally)

**Definition 16.2** (Covariance). *The covariance of random variables  $X$  and  $Y$ , denoted  $\text{Cov}(X, Y)$ , is defined as*

$$\text{Cov}(X, Y) = \mathbb{E}[(X - \mu_X)(Y - \mu_Y)] = \mathbb{E}[XY] - \mathbb{E}[X]\mathbb{E}[Y],$$

where  $\mu_X = \mathbb{E}[X]$  and  $\mu_Y = \mathbb{E}[Y]$ .



“Mean of the product minus the product of the means”

$\underset{\uparrow}{ab\text{Cov}(X,Y)}$

Covariance is Bilinear

$$\begin{aligned} \text{Cov}(aX + bY, aX + bY) &= \text{Cov}(aX, aX) + \text{Cov}(aX, bY) \\ &\quad + \text{Cov}(bY, aX) + \text{Cov}(bY, bY) \\ &\xrightarrow{\text{abCov}(Y,X)} \\ a\text{Cov}(X, aX) &= a^2\text{Cov}(X, X) \\ &\quad \swarrow b^2\text{Cov}(Y, Y) \end{aligned}$$

If  $X$  &  $Y$  are independent, they are uncorrelated

$$\rightarrow \text{Cov}(X, Y) = 0$$

$E[XY]$  =  $E[X]E[Y]$  if  $X$  and  $Y$  are independent

$$\begin{aligned} E[XY] &= \sum_x \sum_y xy \underbrace{P(X=x \wedge Y=y)}_{\nearrow} \\ &= \sum_x \sum_y xy \underbrace{P(X=x)P(Y=y)}_{\nearrow} \\ &= \underbrace{\sum_x x P(X=x)}_{\leftarrow} \sum_y y P(Y=y) \end{aligned}$$

NOTE: If two events are uncorrelated, this doesn't mean they're independent.

$$= E(X) E(Y)$$

# Covariance Relation to Variance

$$\text{Cov}(X, X) = \text{Var}(X)$$



$$E[XX] - E[X]E[X]$$

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

Covariance is Symmetric

$$\text{Cov}(X, Y) = \text{Cov}(Y, X)$$

$$E[XY] - E[X]E[Y]$$

$$E[YX] - E[Y]E[X]$$

equivalent

What does this mean about  $\text{Var}(X+Y)$ ?

$$\text{Cov}(X+Y, X+Y) = \text{Cov}(X, X) + \underline{\text{Cov}(X, Y)}$$

$$\underbrace{\text{Cov}(X, Y)}_{\text{Cov}(Y, X)} + \underline{\text{Cov}(Y, X) + \text{Cov}(Y, Y)}$$

$$= \text{Cov}(X, X) + 2\text{Cov}(X, Y) + \text{Cov}(Y, Y)$$

~~$= \text{Var}(X) + 2\text{Cov}(X, Y) + \text{Var}(Y)$~~

if  $X$  &  $Y$   
are independent  
 $= \text{Var}(X) + \text{Var}(Y)$

# Covariance of Two indicator variables

Let  $X$  and  $Y$  be two not-necessarily-independent Bernoulli random variables. What is  $E[XY]$ ?

$$E(XY) = 0$$

~~+ 0.25~~

~~+ 0.25~~

~~+ 0.25~~

~~+ 0.25~~

$$+ 1 \cdot P(X=1 \wedge Y=1)$$

$XY$	$X$	$Y$	$P(X=x \wedge Y=y)$
0	0	0	$P(X=0 \wedge Y=0)$
0	0	1	$P(X=0 \wedge Y=1)$
0	1	0	$P(X=1 \wedge Y=0)$
1	1	1	$P(X=1 \wedge Y=1)$

$$E(XY) = P(X=1 \wedge Y=1)$$

# Covariance Method with Identically Distributed Variables

Find  $\text{Var}(X_1 + X_2 + \dots + X_n)$  where all the  $X_i$  are identically distributed variables.

$$\text{Cov}(X_1 + X_2 + \dots + X_n, X_1 + X_2 + \dots + X_n)$$

$$= n \text{Cov}(X_i, X_i) + n(n-1) \text{Cov}(X_i, X_j)$$

$$= n(E[X_i^2] - E[X_i]^2) + n(n-1)(E[X_i X_j] - E[X_i] E[X_j])$$

$$= n(P(X_i=1) - (P(X_i=1))^2) + n(n-1) \left( P(X_i=1, X_j=1) - \frac{E[X_i]}{E[X_j]} \right)$$

# Example: Envelopes Problem

You have  $n$  envelopes and  $n$  letters. What is the variance of the number of letters that end up in the correct envelope?

# Example: Envelopes Problem (cont.)

# Envelopes Problem (another approach)

Let  $X = X_1 + X_2 + \dots + X_n$

$$E(X^2) = \sum_i E(X_i^2) + \sum_{i \neq j} E(X_i X_j).$$

# Variance of the Hypergeometric Distribution

$$\frac{10}{20}$$

Suppose you have a box with 10 red balls and 10 green balls. You draw 10 times randomly **without replacement**. What is the variance of the number of red balls you pick?

$X$ : # of red balls

$$X = X_1 + X_2 + \dots + X_{10}$$

$$\text{Var}(X_1 + X_2 + \dots + X_{10})$$

$$X_i \begin{cases} 1 & \text{if you draw a red ball on draw } i \\ 0 & \text{otherwise} \end{cases}$$

$$\text{Cov}(X_1 + \dots + X_{10}, X_1 + \dots + X_{10})$$

$$X_i \sim \text{Bernoulli}\left(\frac{1}{2}\right)$$

## Variance of the Hypergeometric Distribution (cont.)

$$= \underbrace{10 \text{Cov}(X_i, X_i)}_{10 \text{Var}(X_i)} + 10 \cdot 9 \cdot \text{Cov}(X_i, X_j)$$

$$10 \text{Var}(X_i) + 10 \cdot 9 \cdot \text{Cov}(X_i, X_j) \xrightarrow{\text{prob.}}$$

$$10 \cdot \frac{1}{2} \left(1 - \frac{1}{2}\right) + 10 \cdot 9 \cdot \underbrace{\mathbb{E}[X_i X_j]}_{\mathbb{E}[X_i] \mathbb{E}[X_j]}$$

$$10 \cdot \frac{1}{2} \left(\frac{1}{2}\right) + 10 \cdot 9 \cdot \left(\frac{10 \times 9}{20 \times 19} - \frac{1}{4}\right)$$

$$\mathbb{E}[X_i]$$
  
$$\mathbb{E}[X_j]$$

$$E[X_i] = \frac{1}{2}$$

$$E[X_j] = \frac{1}{2}$$

$$E[X_i X_j] = P(X_i = 1 \cap X_j = 1)$$

$$E[X_i X_j] = P(X_i = 1 \cap X_j = 1) = \frac{P(X_i = 1 \cap X_j = 1)}{R}$$

--- --- ---  $\frac{R}{i\text{th}}$

$R$   $R$

1st draw 2nd draw

→ on some draws we got red & on draw  $j$  we got red

$$\frac{10}{20} \times \frac{9}{19}$$

$j+1$  draw

$$10 \cdot \frac{1}{4} + 10 \cdot 9 \cdot \left(-\frac{1}{76}\right)$$

$$\frac{10}{4} - \frac{90}{76}$$

$$\frac{190}{76} - \frac{90}{76} = \frac{100}{76}$$

## Variance of the Hypergeometric (cont.)

What would the variance be if it was **with replacement**?

$$X = X_1 + \dots + X_{10}$$

$$X_i \sim \text{Bernoulli}\left(\frac{1}{2}\right)$$

independent

Binomial(10,  $\frac{1}{2}$ )

$$\text{Var}(X) = 10 \cdot \frac{1}{2} \cdot \left(1 - \frac{1}{2}\right)$$

$$= \frac{10}{2}$$

# (Formally) Variance of Hypergeometric

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



# Example of Hypergeometric: Tennis Balls

~~black~~

There are **three** ~~blue~~ tennis balls, **four** red tennis balls, and **five** green tennis balls in a box. If you draw three times randomly without replacement, what is the variance of the number of red tennis balls that you draw?

$$n: 3$$

$$G: 4$$

$$N: 3+4+5$$

$$\text{Hypergeometric} \left( \frac{12}{3}, \frac{4}{\binom{4}{12}}, \frac{3}{\binom{8}{12}} \right) \left( \frac{\binom{12-3}{3}}{\binom{11}{12-3}} \right)$$

# Correlation (Formally)

It's hard to compare  $\text{Cov}(X, Y)$  and  $\text{Cov}(A, B)$  if A & B have different units than X & Y – numbers can vary pretty heavily based on units. **Correlation** helps us standardize units, in order to cross-compare covariances.

**Definition 16.3** (Correlation). *Suppose X and Y are random variables with  $\sigma(X) > 0$  and  $\sigma(Y) > 0$ . Then, the correlation of X and Y is defined as*

$$\text{Corr}(X, Y) = \frac{\text{Cov}(X, Y)}{\sigma(X)\sigma(Y)}.$$

$$\sigma(X) = \sqrt{\text{Var}(X)}$$
$$\sigma(Y) = \sqrt{\text{Var}(Y)}$$

$$-1 \leq \text{Corr}(X, Y) \leq 1$$

# (If time, derivation) Variance of the Poisson Distribution

Let  $X \sim \text{Poisson}(\lambda)$ . What is  $\text{Var}(X)$ ?

$$\text{Var}(X) = \lambda$$

## (If time, derivation) Variance of the Geometric Distribution

Let  $X \sim \text{Geometric}(p)$ . What is  $\text{Var}(X)$ ?

$$\begin{aligned}\text{Var}(X) &= E(X^2) - E(X)^2 \\ E(X^2) &= \underbrace{p}_{1} + \underbrace{4p(1-p)}_{2^2=4} + \underbrace{q_p(1-p)^2}_{3^2=q} + \dots \\ -(1-p)E(X^2) &= q(1-p) + 4p(1-p)^2 + q_p(1-p)^3 + \dots \\ pE(X^2) &= 2E(X) - 1\end{aligned}$$

$$P(E(X^2)) = 2\left(\frac{1}{P}\right) - 1$$

$$P(E(X^2)) = \frac{2}{P} - 1$$

$$E(X^2) = \frac{\left(\frac{2-P}{P}\right)}{P} = \frac{2-P}{P^2}$$

$$\begin{aligned} \text{Var}(X) &= E(X^2) - E(X)^2 \\ &= \frac{2-P}{P^2} - \left(\frac{1}{P}\right)^2 = \frac{2-P-1}{P^2} = \frac{1-P}{P^2} \end{aligned}$$

# Recap

- Discussed the topics of variance, covariance, and correlation
  - Properties of Variance
  - Properties of Covariance
  - Properties of Correlation
- Discussed the variance of popular distributions and their derivations