

9PM Thursday ~ 9PM Friday  
(10/8) (10/9)

→ HW6 Q1, Q2

Warm up: A drawer contains  $S$  black socks and  $S$  white socks ( $S > 0$ ). I pull two socks out at random without replacement and call that my first pair. Then I pull out two socks at random without replacement and call that my second pair. I proceed in this way until I have  $S$  pairs and the drawer is empty. Find the expected number of pairs in which two socks are of different colors.

Step 1  $X = \text{Number of pairs out of } S \text{ that mismatch}$

Step 2  $I_j = \begin{cases} 1 & \text{if } j\text{th pair is mismatch} \\ 0 & \text{else} \end{cases}, j=1, \dots, S$

Step 3  $p = P(I_j = 1) \stackrel{\text{Symmetry}}{=} P(I_1 = 1) = \frac{\binom{S}{1}\binom{S}{1}}{\binom{2S}{2}} \quad \begin{array}{l} \text{draw 2 black} \\ \text{draw 1 white} \end{array} = 2 * \frac{S}{2S} * \frac{S}{2S-1}$

Step 4  $X = I_1 + I_2 + \dots + I_S$  (Hypergeometric formula)

Step 5  $E(X) = S \cdot p$

Alternatively  
White - Black  
Black - white

Step 3

pop

$p = P(\text{1st pair is mismatch})$

$= P(WB) + P(BW)$

$= P(W)P(B|W) + P(B)P(W|B)$

$= \frac{S}{2S} * \frac{S}{2S-1} + \frac{S}{2S} * \frac{S}{2S-1}$

$= 2 * \frac{S}{2S} * \frac{S}{2S-1}$

$P(\text{1st trial is White sock})$

$= \frac{S}{2S}$

Step 5

$X = I_1 + I_2 + \dots + I_S$

$E(X) = E(I_1 + I_2 + \dots + I_S)$

$= E(I_1) + E(I_2) + \dots + E(I_S)$

$= p + p + \dots + p$

$= S \cdot p$

## 5.4. Unbiased Estimators

**Preliminary: Linear Function Rule** Let  $X$  be a random variable and let  $Y = aX + b$ . Then  $Y$  is a linear function of  $X$ . Then

$$E(Y) = E(aX + b) = \sum_{\text{all } x} (ax + b)P(X = x)$$

$$E(g(x)) = \sum_{\text{all } x} g(x)P(X=x) \quad \text{with } g(x) = ax+b$$

$$= a \underbrace{\sum_{\text{all } x} xP(X=x)}_{E(X)} + b \underbrace{\sum_{\text{all } x} P(X=x)}_1 = aE(X) + b$$

**Terminology** Data scientists often want to estimate a parameter of a population.

- A **parameter** is a fixed unknown number associated with the population.
- A **statistic** is a number based on the data in your sample.
- An **estimator** is a statistic used to approximate a parameter.   
 *< estimate*
- An **unbiased estimator** of a parameter is an estimator whose expected value is equal to the parameter.

Sample mean as an estimator of population mean

Ex Estimate the average annual income in California,  $\mu$ .

*parameter*  
 $\rightarrow E(X_i) = \mu$  for  $i=1 \rightarrow n$  Each  $X_i$  from the same distribution (population)

Suppose you draw a random sample of size  $n$ .  $X_1, \dots, X_n$  are sample incomes. The sample average is the statistic  $\bar{X}$  defined as the function

$$\bar{X} = g(X_1, X_2, \dots, X_n) = \frac{1}{n} \sum_{i=1}^n X_i. \quad \text{- statistic}$$

Important:  $\mu$  is constant, unknown;  $\bar{X}$  estimator, random variable

$\bar{X}$  is unbiased if  $E(\bar{X}) = \mu$ . In fact,

*repeated samples*  $\rightarrow$  *Long run avg. value of  $\bar{X}$ , on avg, it is equal to  $\mu$*

$$\text{PF } E(\bar{X}) = E\left(\frac{1}{n} \sum_{i=1}^n X_i\right) = \frac{1}{n} E\left(\sum_{i=1}^n X_i\right) = \frac{1}{n} \sum_{i=1}^n E(X_i) = \frac{1}{n} n \cdot \mu = \mu.$$

*$E(aX) = aE(X)$  for any const  $a$*   
*Adding  $E(X+Y) = E(X) + E(Y)$*

$n=3$  *pop*

Today  $X_1 = 50K, X_2 = 60K, X_3 = 70K \rightarrow \bar{X} = 60K$   
 Tomorrow  $X_1 = 65K, X_2 = 50K, X_3 = 60K \rightarrow \bar{X} = \frac{50K + 60K + 65K}{3}$   
 ...

*Avg  $\mu$*

Which of these estimators of  $\mu$  is unbiased?

(a)  $X_{15}$ .

Draw a sample  $X$  from your pop. whose mean  $= \mu$ .

(b)  $(X_1 + X_{15})/15$ .

Example  $E(X) = \mu$ .

$\hookrightarrow X_1 \sim \text{Binomial}(n, p)$   $\rightarrow np$

$E(X_1) = np$

$X_2 \sim \text{Binomial}(n, p) \rightarrow E(X_2) = np$

(c)  $(X_1 + 2X_{100})/3$ .

(a)  $E(X_{15}) = \mu$  Unbiased

(b)  $E\left(\frac{X_1 + X_{15}}{15}\right) = \frac{1}{15} E(X_1 + X_{15}) = \frac{E(X_1) + E(X_{15})}{15} = \frac{2\mu}{15}$  biased

(c)  $E\left(\frac{X_1 + 2X_{100}}{3}\right) = \frac{1}{3} E(X_1 + 2X_{100}) = \frac{1}{3} * (\mu + 2\mu) = \mu$  Unbiased

If we have a biased estimator how can we make it unbiased?

Let's make  $\frac{X_1 + X_{15}}{3}$  unbiased.

$$\begin{aligned} E\left(\frac{X_1 + X_{15}}{3}\right) &= \frac{E(X_1) + E(X_{15})}{3} \\ &= \frac{\mu + \mu}{3} \\ &= \frac{2\mu}{3} \end{aligned}$$

$\hookrightarrow \frac{3}{2} E\left(\frac{X_1 + X_{15}}{3}\right) = \mu$

$\hookrightarrow E\left(\frac{3}{2} * \left(\frac{X_1 + X_{15}}{3}\right)\right) = \mu$

//  
 $E\left(\frac{X_1 + X_{15}}{2}\right)$

Unbiased estimator

(a special case of sample mean when the pop. consists of 0s and 1s)

## Sample proportion as an estimator of population proportion

When the population consists of zeros and ones, the population mean is the population proportion of ones.

Example

population				
0	0	1	1	1

population mean =  $p = \frac{3}{5}$

$X_1, \dots, X_n$  samples w/ replacement from pop.

$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i = \text{Sample proportion}$ ,  $E(\bar{X}) = p$

Example You roll a die 30 times and find the sample proportion of sixes. The population consists of  $\{0, 0, 0, 0, 0, 1\}$ . Repeat experiment 20,000 times and plot distribution of sample proportions.

draw w/ replacement

Exp 1:  $X_1, X_2, \dots, X_{30} \rightarrow \bar{X}^{(1)}$

Exp 2:  $X_1, X_2, \dots, X_{30} \rightarrow \bar{X}^{(2)}$

20,000 different

observed values of  $\bar{X}$

$n = 30$

$p = 0.1667$

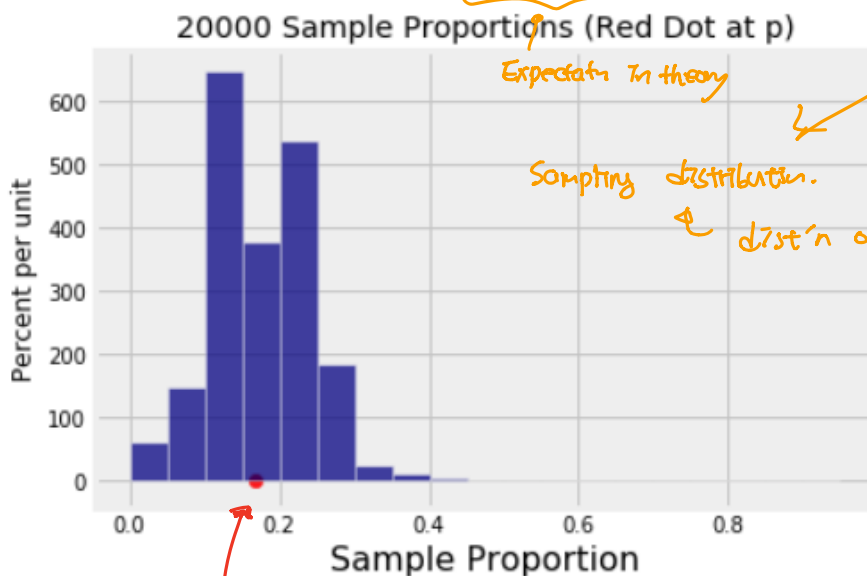
Exp 20,000:  $X_1, \dots, X_{30} \rightarrow \bar{X}^{(20000)}$

Average of observed sample proportions = 0.1664

$0.1667 = p = E(\bar{X}) \approx \frac{1}{20000} \sum_{i=1}^{20000} \bar{X}^{(i)} = 0.1664$

$0.1667 = p = E(\bar{X})$

Expectation in numerical exp.



Expectation in theory

Sampling distribution.

distribution of your estimator ( $= \bar{X}$ )

histogram in numerical exp.

$p$  Since  $E(\text{sample proportion}) = p$

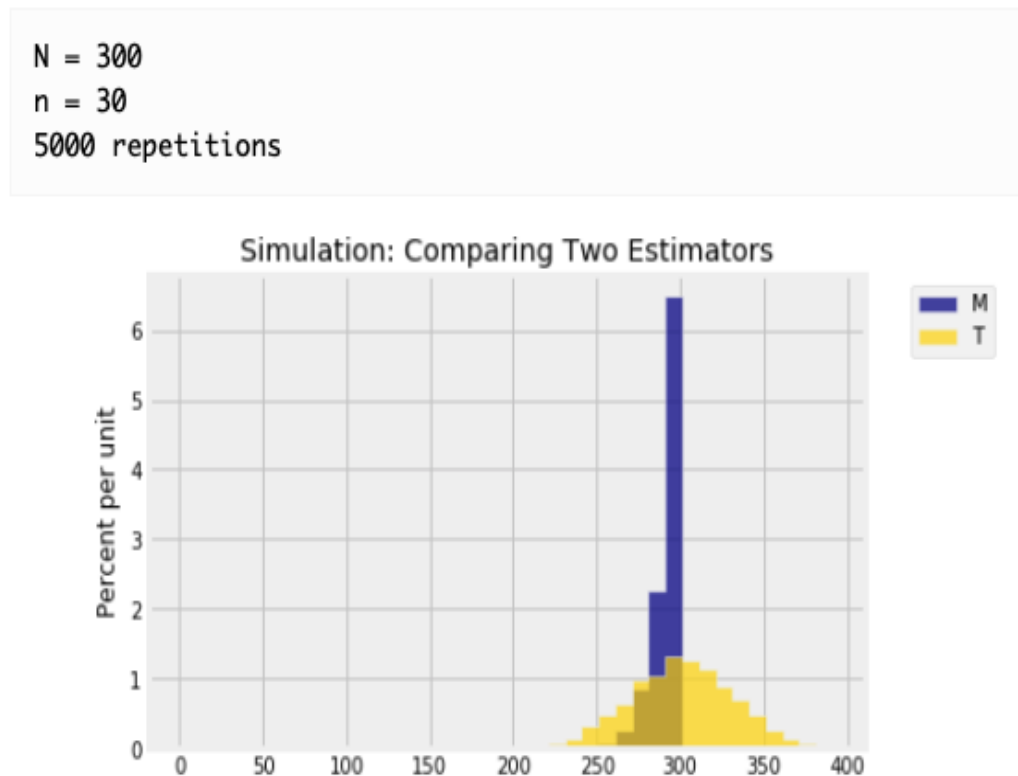
## Estimating the largest possible value

Let  $X_1, X_2, \dots, X_n \stackrel{i.i.d.}{\sim} \text{Uniform}\{1, 2, \dots, N\}$  for some fixed but unknown  $N$ . To estimate  $N$ , you may think  $M = \max\{X_1, \dots, X_n\}$  and this is an estimator but we want an unbiased estimator.

The population mean is  $\mu = (N + 1)/2$  and  $E(\bar{X}) = (N + 1)/2$  since it is unbiased. What is an estimator such that

$$E(\text{estimator}) = N?$$

Lets look at sampling distribution of (1)  $T = 2\bar{X} - 1$  and (2)  $M = \max(X_1, \dots, X_n)$ .



The histograms show that both estimators have pros and cons.

$M$  - Pros: small spread of values; Cons: biased.

$T$  - Pros: unbiased; Cons: big spread of values.

Unbiasedness is a good property, but so is low variability. Bias-variance tradeoff