

* Announcement

① HW11 due today (11:59 PM PT)

② HW12 ~ 11/23

③ Quiz 10 : Ch 10.3 ~ Ch 11.1

After today's lecture (HW12)

→ Q1, Q2, Q3

(Sec 10.4)

Q4, Q5

(Sec 11.1 ~ 11.2)

STAT 88: Lecture 34

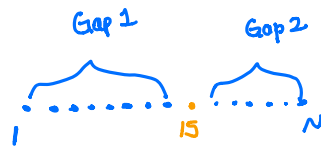
Contents

Section 11.1: Bias and Variance

Section 11.2: The German Tank Problem, Revisited

Warm up:

N consecutive positive integers, $1, 2, \dots, N$ are in a hat, with N unknown. You randomly sample one number from the hat and get 15. Estimate N .



$$\rightarrow \text{Gap 1} + \text{Gap 2} + 1 = N$$

$$\rightarrow E(\text{Gap 1}) + E(\text{Gap 2}) + 1 = N$$

$$\underbrace{\quad}_{\text{G}} \quad \underbrace{\quad}_{\text{G}}$$

(by symmetry, $E(\text{Gap 1}) = E(\text{Gap 2})$)

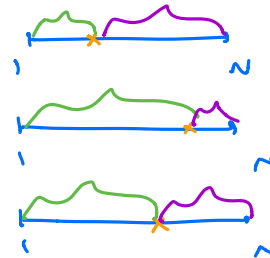
$$\rightarrow 2G + 1 = N$$

$$\uparrow$$

$$14$$

\Rightarrow Estimate G by 14

\Rightarrow Estimate N by $2 \cdot 14 + 1 = 29$



parameter
 $E(N) = N$

Last time

Suppose you have two independent samples, X_1, X_2, \dots, X_n i.i.d. with mean μ_X and SD σ_X , and Y_1, Y_2, \dots, Y_m are i.i.d. with mean μ_Y and SD σ_Y . By CLT, $\bar{X} - \bar{Y}$ is approximately distributed as

$$\mathcal{N}\left(\mu_X - \mu_Y, \frac{\sigma_X^2}{n} + \frac{\sigma_Y^2}{m}\right).$$

An approximate 95% CI for $\mu_X - \mu_Y$ is given by

$$\bar{X} - \bar{Y} \pm 2\sqrt{\frac{\sigma_X^2}{n} + \frac{\sigma_Y^2}{m}}.$$

If we wish to test $H_0 : \mu_X = \mu_Y$ vs $H_A : \mu_X > \mu_Y$, $T = \bar{X} - \bar{Y}$ is our test statistic. Under H_0 , we have

$$T \sim \mathcal{N}\left(0, \frac{\sigma_X^2}{n} + \frac{\sigma_Y^2}{m}\right).$$

If t is an observed value of our test statistic, p-value is given by

$$\text{p-val} = P(T \geq t) = P\left(Z \geq \frac{t}{\sqrt{\frac{\sigma_X^2}{n} + \frac{\sigma_Y^2}{m}}}\right) = 1 - \Phi\left(\frac{t}{\sqrt{\frac{\sigma_X^2}{n} + \frac{\sigma_Y^2}{m}}}\right).$$

Next, you have two independent populations of 0's and 1's, so X_1, X_2, \dots, X_n are i.i.d. from $\text{Bernoulli}(p_X)$, and Y_1, Y_2, \dots, Y_m are i.i.d. from $\text{Bernoulli}(p_Y)$. By CLT, $\bar{X} - \bar{Y}$ is approximately distributed as

$$\bar{X} - \bar{Y} \sim \mathcal{N}\left(p_X - p_Y, \frac{p_X(1-p_X)}{n} + \frac{p_Y(1-p_Y)}{m}\right).$$

An approximate 95% CI for $p_X - p_Y$ is given by

$$\bar{X} - \bar{Y} \pm 2\sqrt{\frac{p_X(1-p_X)}{n} + \frac{p_Y(1-p_Y)}{m}}.$$

Estimate $p_X \leftarrow \bar{X}$
 $p_Y \leftarrow \bar{Y}$

If we wish to test $H_0 : p_X = p_Y = p$ vs $H_A : p_X > p_Y$, $T = \bar{X} - \bar{Y}$ is our test statistic. Under H_0 , we have

$$T \sim \mathcal{N}\left(0, \frac{p(1-p)}{n} + \frac{p(1-p)}{m}\right).$$

Estimate p by $\frac{n}{n+m}\bar{X} + \frac{m}{n+m}\bar{Y}$

If t is an observed value of our test statistic, p-value is given by

$$\text{p-val} = P(T \geq t) = P\left(Z \geq \frac{t}{\sqrt{\frac{p(1-p)}{n} + \frac{p(1-p)}{m}}}\right) = 1 - \Phi\left(\frac{t}{\sqrt{\frac{p(1-p)}{n} + \frac{p(1-p)}{m}}}\right).$$

11.1. Bias and Variance

Bias-Variance Decomposition

Some estimators for a parameter θ are better than others. A good estimator is one with a small mean square error.

where

$$\text{MSE}_\theta(T) = E_\theta((T - \theta)^2) = B_\theta^2(T) + \text{Var}_\theta(T),$$

\swarrow estimator
 \uparrow parameter
 \swarrow Bias
 \searrow Variance

$$B_\theta(T) = E_\theta(T) - \theta \text{ and } \text{Var}_\theta(T) = E_\theta((T - E_\theta(T))^2).$$

11.2. The German Tank Problem

The Allies during WWII needed to estimate how many Tanks N the Germans had produced. The idea was to model the observed serial numbers as random draws from $1, 2, \dots, N$ and then estimate N .

So we will now assume, as the Allies did, that the serial numbers of the observed tanks are random variables X_1, \dots, X_n drawn uniformly at random without replacement from $\{1, 2, \dots, N\}$. That is, we have a simple random sample of size n from the population $\{1, 2, \dots, N\}$ and we have to estimate N . $X_1, \dots, X_n \sim \text{Unif}\{1, 2, \dots, N\}$

Now we will compare several estimators.

T_1 : By symmetry, for each i ,

$$E(X_i) = \frac{N+1}{2}.$$

\swarrow Expectation of $\text{Unif}\{1, \dots, N\}$
 $= \text{pop. mean}$

Since \bar{X} is an unbiased estimator for the pop. mean,

$$\frac{X_1 + \dots + X_n}{n} \sim \bar{X}$$

$$E(\bar{X}) = \frac{N+1}{2}.$$

$\rightarrow 2E(\bar{X}) - 1 = N$
 $\rightarrow E(2\bar{X} - 1) = N$

This is a linear function of N so we can find an unbiased estimator for N :

$$E(\underbrace{2\bar{X} - 1}_{=T_1}) = N.$$

If we define

$$T_1 = 2\bar{X} - 1, \quad \rightarrow \quad E(T_1) = N$$

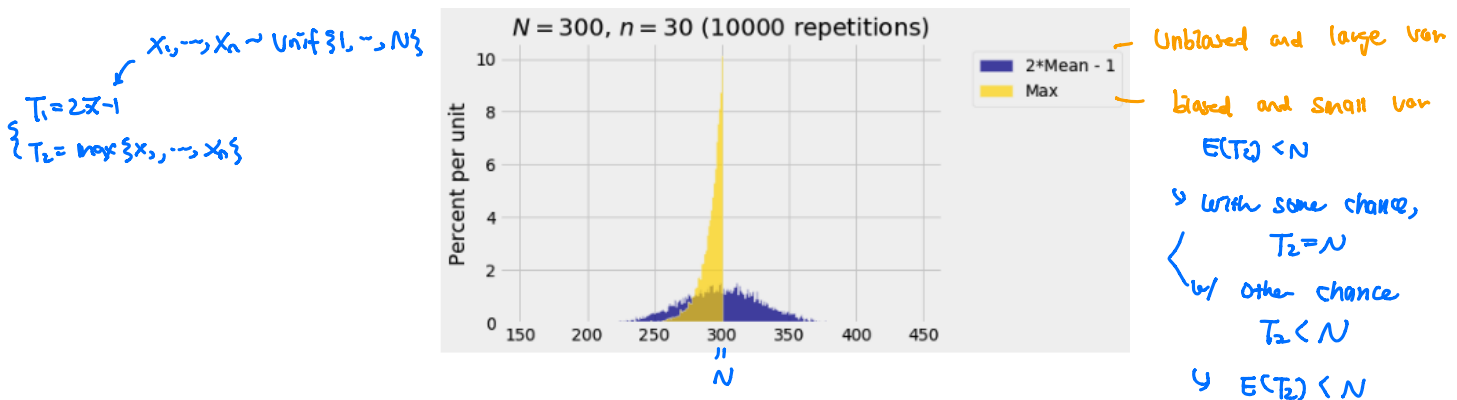
it is an unbiased estimator of N .

T_2 : Another natural estimator is

$$X_1, \dots, X_n \sim \text{Unif}\{1, \dots, N\}$$

$$T_2 = \max\{X_1, X_2, \dots, X_n\},$$

the maximum of the observed numbers.



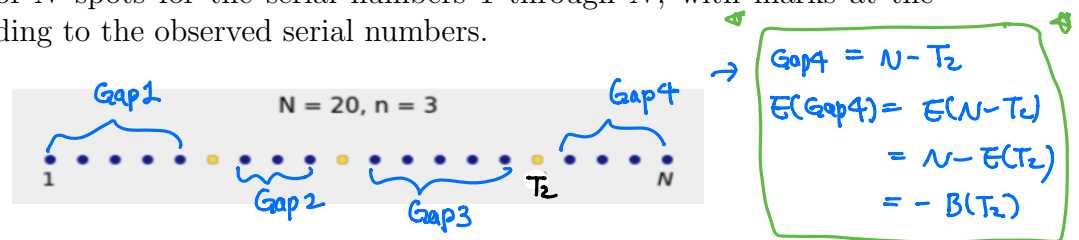
T_2 is clearly biased. We compute the bias of T_2 .

The Bias of the Sample Maximum The bias of T_2 is

$$B(T_2) = E(T_2) - N.$$

Bias of T_2

Imagine a row of N spots for the serial numbers 1 through N , with marks at the spots corresponding to the observed serial numbers.



The $n = 3$ sampled spots create $n + 1 = 4$ blue "gaps" between sampled values: one before the leftmost gold spot, two between successive gold spots, and one after the rightmost gold spot that is at position T_2 .

By symmetry, the lengths of all four gaps have the same distribution. Therefore all four gaps have the same expected length.

- The gaps are made up of $N - n = 17$ blue spots;

$$\begin{aligned} \text{Gap1} + \text{Gap2} + \text{Gap3} + \text{Gap4} &= 17 \\ \rightarrow E(\text{Gap1}) + E(\text{Gap2}) + \dots + E(\text{Gap4}) &= 17 \\ \rightarrow \text{By symmetry, } E(\text{Gap1}) = \dots = E(\text{Gap4}) &= G, \end{aligned}$$

$$\text{so } 4 \cdot G = 17$$

$$\rightarrow G = \frac{17}{4}$$

Generalize

4

$$\frac{N-n}{n+1}$$

- Since each of the four gaps has the same expected length, the expected length of a single gap is $\frac{17}{4}$.

More generally,

$$\text{expected length of gap} = \frac{N - n}{n + 1}.$$

Note that $N - E(T_2)$ is the expected length of the last gap. So,

$$\overset{\text{"E(Gap4)"}}{N - E(T_2)} = \frac{N - n}{n + 1},$$

and therefore

$$B(T_2) = E(T_2) - N = \frac{-(N - n)}{n + 1}. \quad < 0 \quad \Rightarrow \quad E(T_2) < N$$

T_3 : (the “augmented maximum”)

What is $E(T_2)$?
$$E(T_2) = N + B(T_2) = N + \frac{-(N - n)}{n + 1} = \frac{n}{n + 1}(N + 1)$$

Since $E(T_2)$ is a linear function of N , we can make a new unbiased estimator by solving for N .

$$E(T_2) = \frac{n}{n + 1}(N + 1) \Rightarrow \frac{n + 1}{n} E(T_2) - 1 = N$$

$$\Rightarrow E\left(\underbrace{\frac{n + 1}{n} \cdot T_2 - 1}_{= T_3}\right) = N.$$

If we define

$$T_3 = \frac{n + 1}{n} \cdot T_2 - 1, \quad E(T_3) = N$$

it is an unbiased estimator of N .

- How does $\text{Var}(T_3)$ and $\text{Var}(T_2)$ compare? $\text{Var}(T_3) = \text{Var}\left(\frac{n + 1}{n} \cdot T_2\right) = \left(\frac{n + 1}{n}\right)^2 \text{Var}(T_2)$
- How does $B(T_3)$ and $B(T_2)$ compare? \approx If n is large, $\text{Var}(T_3) \approx \text{Var}(T_2)$

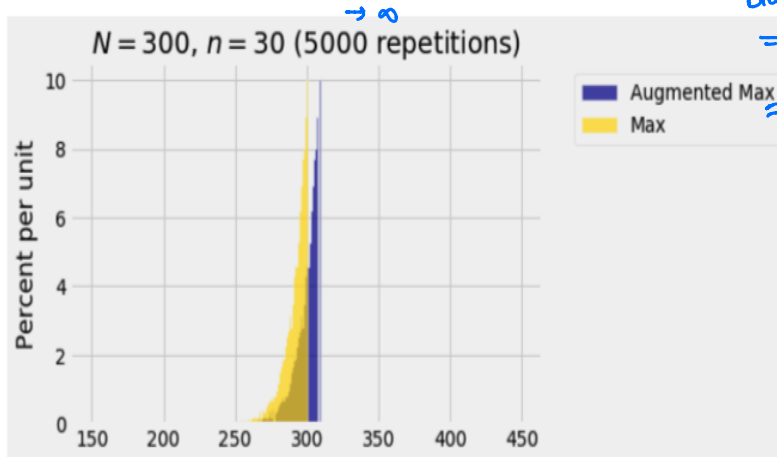
So for large n , T_3 is better than T_2 .

$$\left\{ \begin{array}{l} \text{MSE}(T_3) = \text{bias}^2 + \text{Var} \\ \quad = 0 + \text{Var}(T_3) \\ \text{MSE}(T_2) = \left(\frac{N - n}{n + 1}\right)^2 + \text{Var}(T_2) \end{array} \right. \Rightarrow T_3 \text{ has lower MSE than } T_2$$

$\overset{= E(T_3)}{\text{Average of Augmented Maxes: } 300.18587333333335}$
 $\text{SD of of Augmented Maxes: } 8.947086216787126 \overset{= SD(T_3)}{=}$
 $\text{Average of Maxes: } 291.4702 \overset{= E(T_2)}{=}$
 $\text{SD of Maxes: } 8.65847053237464 \overset{= SD(T_2)}{=}$

$\rightarrow \text{"T}_3\text{"}$
 $\rightarrow \text{Bias}^2 + \text{var} \overset{= N}{=}$
 $= (300.186 - 300)^2 + 8.95^2$
 ≈ 79.956

$\rightarrow \text{"T}_2\text{"}$
 $\text{Bias}^2 + \text{var}$
 $= (291.47 - 300)^2 + 8.66^2$
 ≈ 147.58

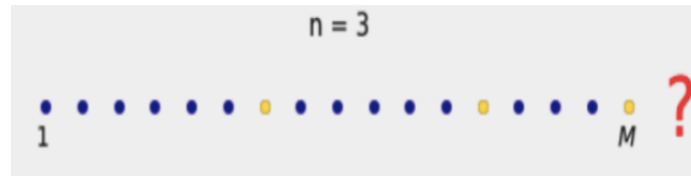


$\max\{x_1, \dots, x_n\}$

$\frac{n+1}{n} T_2 - 1$
 $\frac{2x-1}{2}$

In summary, we can have many different estimators for a parameter. In this lecture,
 T_1 was unbiased but had a large variance, T_2 was biased but had a smaller variance.
 T_3 was unbiased and had a bigger variance but for large n $\text{Var}(T_3) \approx \text{Var}(T_2)$. The
 estimator with the smallest $\text{MSE} = \text{Bias}^2 + \text{Var}$ is the best.

Another way to think of the “augmented maximum” If we could see the gap to the right of T_2 , we would see N . But we can't. So we can try to do the next best thing, which is to augment T_2 by the estimated size of that gap.



What is the estimated gap length?

Hence we can try to improve upon T_2 using the estimator

$$T_3 =$$