

Stat 585 Lec 34

Warm up 2:00 - 2:10

N consecutive positive integers, $1, 2, 3, \dots, N$ are in a hat, with N unknown.

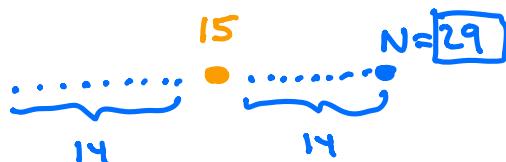
You randomly sample one number from the hat and get 15, Estimate N



lets call the blue dots to the left of the gold dot a "gap".

By symmetry we would expect an equal sized gap on the other side of the gold dot.

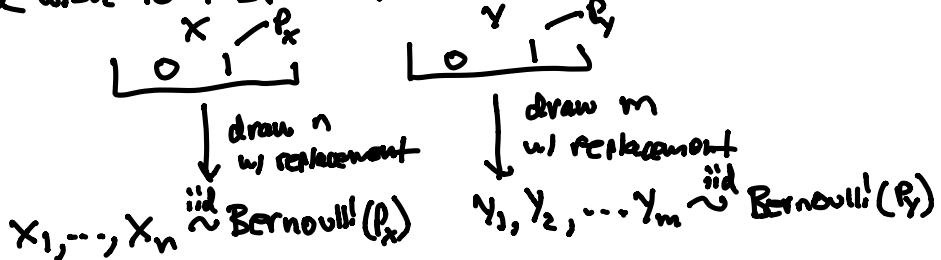
so $14 + 1 + 14 = 29$ is a reasonable estimate of N .



Last time Q4 - next Friday Chap 8, 9, 10

Sec 10.4 Test for the equality of two proportions

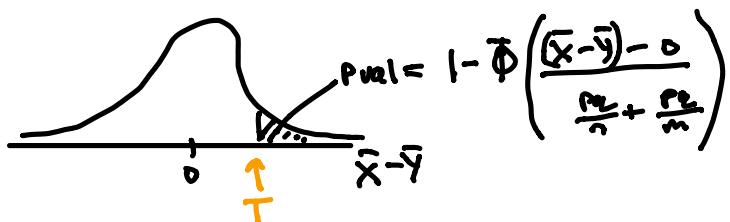
We have two independent populations of 0, 1s and wish to test that the proportion of 1s is the same.



$$H_0: p_x = p_y = p \quad \text{level } 5\% \text{ test.}$$

$$H_1: p_x > p_y$$

$T = \bar{X} - \bar{Y}$ is our T.S., and we reject the null if $P_{\text{val}} < .05$



Sec 11.1 Bias and Variance

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

$$\text{mse}_{\theta}(\tau) = E_{\theta}((\tau - \theta)^2) = \text{Bias}_{\theta}^2(\tau) + \text{Var}_{\theta}(\tau)$$

\uparrow estimator
Parameter

where $\text{Bias}_{\theta}(\tau) = E(\tau) - \theta$ is the bias,

$$\text{Var}_{\theta}(\tau) = E_{\theta}((\tau - E(\tau))^2) \text{ is the variance,}$$

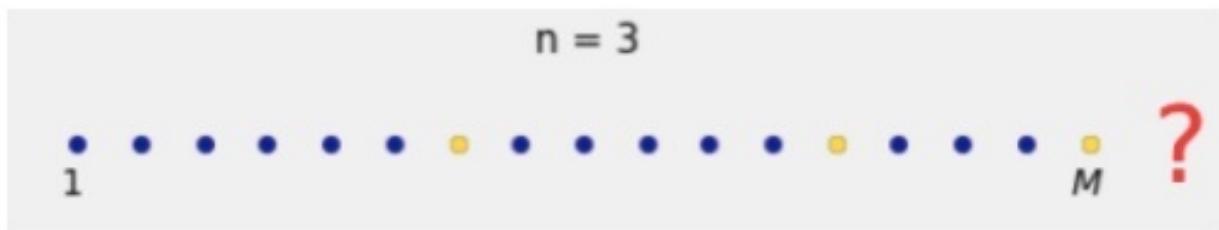
Today

- (1) Sec 11.2 The German Tank Problem,

① Sec 11.2 The German Tank Problem

The Allies during WW II 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 .

Suppose they capture 3 tanks with serial numbers 6, 10, 16.



Estimate N in your breakout room.

$$\text{average gap size} = \frac{\text{total # blue dots}}{\# gaps} = \frac{m-3}{3}$$

$$= \frac{16-3}{3} = \frac{13}{3}$$

By symmetry there should be a gap of size $\frac{13}{3}$ to the right of M $\Rightarrow N \approx m + \frac{m-3}{3} = \frac{4m-1}{3} = \frac{4}{3} \cdot 16 - 1 = 20.22$

We will see the estimator $\frac{4}{3}m-1$ is unbiased.

So we will now assume, as the Allies did, that the serial numbers of the observed tanks are random variables X_1, X_2, \dots, X_n drawn uniformly at random without replacement from $\{1, 2, 3, \dots, N\}$. That is, we have a simple random sample of size n from the population $\{1, 2, 3, \dots, N\}$, and we have to estimate N .

We will compare several estimators.

||^T

By symmetry

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

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

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

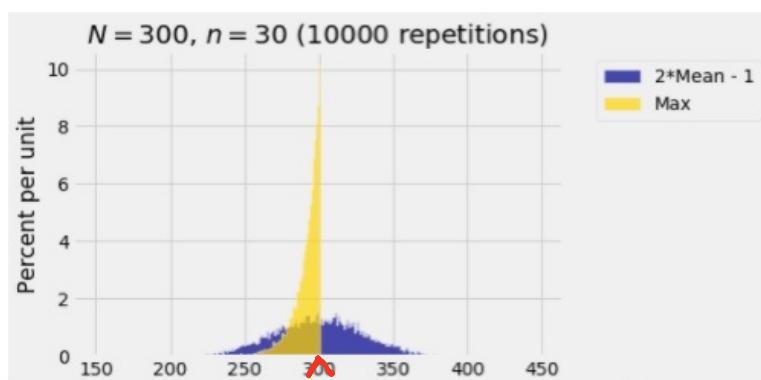
$$N = 2E(\bar{x}) - 1$$

let $T_1 = 2\bar{x} - 1$

check $E(T_1) = 2E(\bar{x}) - 1 = N \checkmark$

||^T

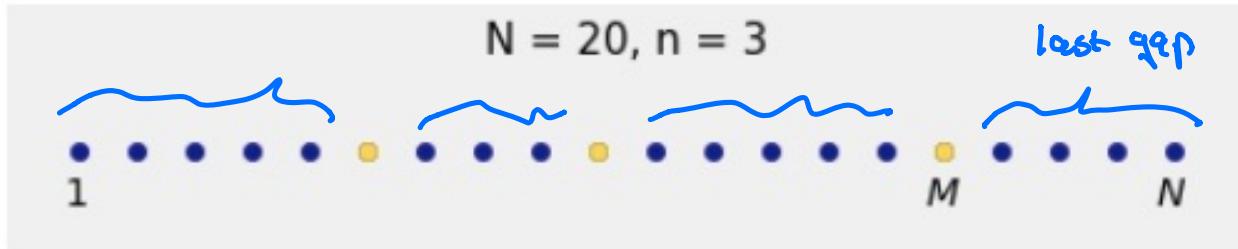
let $T_2 = \max(x_1, \dots, x_n)$



T_2 is
clearly
biased.

Compute bias of T_2 :

$$B_N(T_2) = E(M) - N$$



by symmetry, the expected length of all the gaps is the same.

$N - E(M)$ is the expected length of the last gap,

soll $N - E(M) \Rightarrow ?$ $\frac{20-3}{4} \leftarrow \begin{matrix} 3 \text{ gold} \\ 4 \leftarrow 4 \text{ gaps} \end{matrix}$
 $= E(M) - N$

so $B_N(T_2) = -\frac{(20-3)}{4}$

more generally,

$$B_N(T_2) = ? \quad \frac{-(N-n)}{n+1}$$

write $E(M)$ as a function of N breakout room

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

or $E(M) = \frac{(N-n)}{n+1} \cdot n + n = \frac{n}{n+1} (N-1) + n$
 $= \boxed{\frac{n}{n+1} (N+1)}$

T_3 (the "augmented maximum")

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

$$E(M) = \frac{n}{n+1} (N+1)$$

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

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

so $T_3 = \frac{n+1}{n} M - 1$ is an unbiased estimator

$$E(T_3) = \frac{n+1}{n} E(M) - 1 = N$$

In the case where $M=16$
 $n=3$

$$T_3 = \frac{4}{3} \cdot 16 - 1 = 20.33$$

How does $\text{Var}(T_3)$ and $\text{Var}(T_2)$ compare?

Poll

$$T_3 = \frac{n+1}{n} T_2 - 1$$

$$\text{Var}(T_3) = \left(\frac{n+1}{n}\right)^2 \text{Var}(T_2)$$

$$\text{So } \text{Var}(T_3) > \text{Var}(T_2)$$

but almost equal for large n

How does $B_n(T_3)$ and $B_n(T_2)$ compare?

$$\frac{n}{n+1} \cdot \frac{(n-n)}{n}$$

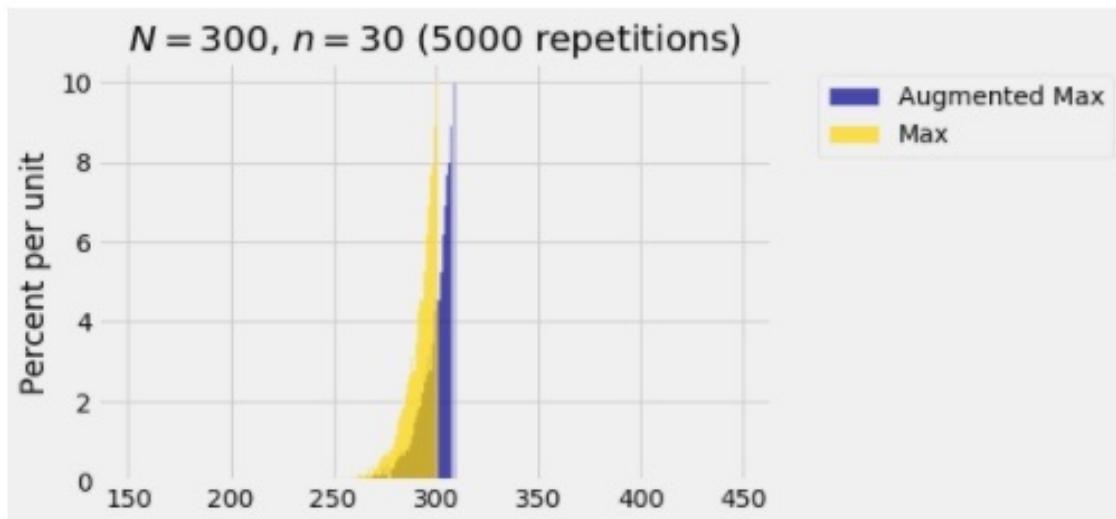
large n , T_3 is a better estimator

Average of Augmented Maxes: 300.1858733333335

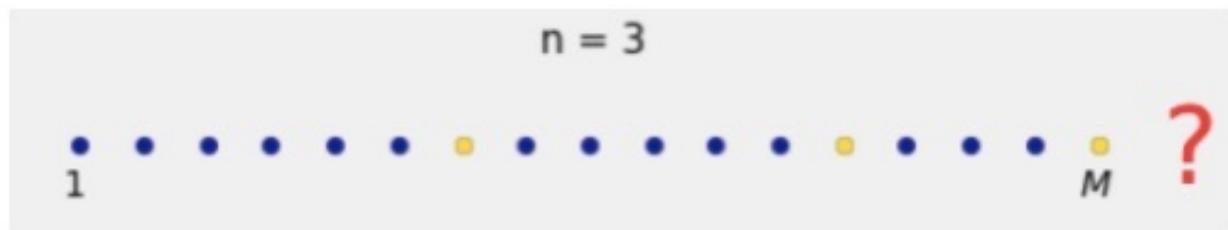
SD of Augmented Maxes: 8.947086216787126

Average of Maxes: 291.4702

SD of Maxes: 8.65847053237464



Another way to think of the "augmented max"



estimated gap length ? $\frac{M-n}{n}$

$$N \approx ? M + \frac{M-n}{n} = \frac{n+1}{n} M - 1$$

which is what we derived
earlier,

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 $\text{Var}(T_2) \approx \text{Var}(T_3)$).

The estimator with the smallest $\text{MSE} = \text{bias}^2 + \text{var}$ is best.