

Last time:

Sec 10.3 Exponential

Sec 10.10 Normal

Today,

finish normal

Bias / Unbiased Estimator.

Sec 10.4 (Continued)

C) & test for difference between population means ($\mu_1 - \mu_2$)
proportions (p_1, p_2)

setup	size	sample proportion	population proportion
City X	$n = 400$	37%	p_1
City Y	$m = 600$	28%	p_2

For Bern(p), mean = p, Var = p(1-p), sd = $\sqrt{p(1-p)}$.

sd of sample mean = $\frac{1}{\sqrt{n}}$ sd = $\sqrt{\frac{p(1-p)}{n}}$

By CLT, $\bar{X} \sim N(p_1, \frac{p_1(1-p_1)}{n})$

$\bar{Y} \sim N(p_2, \frac{p_2(1-p_2)}{m})$ (3%)

$\bar{X} - \bar{Y} \sim N(p_1 - p_2, \frac{p_1(1-p_1)}{n} + \frac{p_2(1-p_2)}{m})$

$\bar{X} - \bar{Y} = 37\% - 28\% = 9\%$

sd($\bar{X} - \bar{Y}$) = $\sqrt{\frac{37\% \times 63\%}{400} + \frac{28\% \times 72\%}{600}} = 3\%$

95% CI for $p_1 - p_2$ is

$9\% \pm 2 \times 3\% = [3\%, 15\%]$

test,

H_0 : City X and City Y has the same proportion

H_a : City X has a higher proportion.

What is the difference between a 95% CI

and a test with significant level 0.05?

H_0 : $p_1 = p_2$ $p_1 - p_2 = 0$

H_a : $p_1 > p_2$ $p_1 - p_2 > 0$

test statistic = $\bar{X} - \bar{Y}$

distr under H_0 : $N(0, (3\%)^2)$

observed value: $\bar{X} - \bar{Y} = 9\% \approx +3$ SD

Conclusion: reject H_0

Bias

"Estimator"

R.V.

T is an estimator to θ

Computed from data

fixed

unknown parameter.

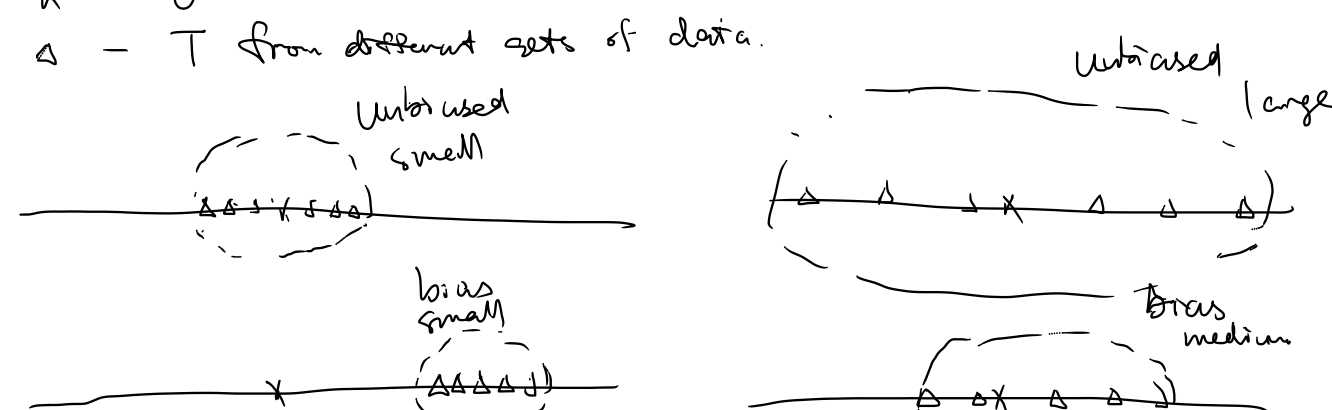
Why are we using these estimators?

What makes a good estimator?

Error = $T - \theta$

$x = \theta$

$\Delta = T$ from different sets of data.



1) bias / unbiased

2) variance

θ - fixed unknown parameter.

T - estimator

fixed
R.V.

Bias: $B_0(T) = E(T - \theta)$

$= E(T) - \theta$

unbiased: $B_0(T) = 0$, or $E(T) = \theta$

Example:

$X_1, X_2, X_3, \dots, X_n$ i.i.d with mean μ

X_1 is an unbiased estimator of μ

since $E(X_1) = \mu$

$\bar{X} = \frac{1}{n}(X_1 + \dots + X_n)$ is also an unbiased estimator of μ

$E(\bar{X}) = \frac{1}{n} \sum_{i=1}^n E(X_i) = \frac{1}{n} \cdot n \cdot \mu = \mu$

In particular for Bern(p) samples

sample proportion is an unbiased estimator of population proportion

$MS E_{\mu}(X_1) = Var(X_1)$

$MS E_{\mu}(\bar{X}) = Var(\bar{X})$

$\frac{1}{n} Var(X_1)$

Example 2

uniform $\{1, 2, \dots, N\}$

X_1, X_2, \dots, X_n

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

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

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

$\Rightarrow 2\bar{X} - 1$ is an unbiased estimator of N

WW2 German tanks

N X_1, X_2, \dots, X_n SRS from $\{1, 2, \dots, N\}$

$T_1 = 2\bar{X} - 1$

$n = 3$

$X_1 = 1, X_2 = 2, X_3 = 999$

$\bar{X} = \frac{1}{3}(1 + 2 + 999) = 334$

$2\bar{X} - 1 = 667$

not reasonable since $999 > 667$

$T_2 = M = \max(X_i) \leq N$

$E(T_2) < N \Rightarrow$ biased!

See 11.1

Mean Squared Error

$MS E_{\theta}(T) = E_{\theta}(T - \theta)^2$ just means this is a function of θ .

lower $MS E \Rightarrow$ good.

Decomposition of Error

Bias: $B_0(T) = E_{\theta}(T) - \theta$ - constant.

Deviation: $D_0(T) = T - E_{\theta}(T)$ - random

Error: $T - \theta = D_0(T) + B_0(T)$.

$E_{\theta}(T)$

$MS E_{\theta}(T) = E_{\theta}(T - \theta)^2$

$= E_{\theta}(D + B)^2$

$= E_{\theta}(D)^2 + 2E_{\theta}(D) \cdot B + B^2$

$= Var_{\theta}(T) + B_0^2(T)$

$E_{\theta}(D) = E_{\theta}(T - E_{\theta}(T))$

$= E_{\theta}(T) - E_{\theta}(T)$

$= 0$

$MS E$ is the sum of variance of the estimator and the squared bias.

Notes: in terms of M.S.E. the only aspects that we take into account are $Var_{\theta}(T)$ and $B_0(T)$.

Example Germany tank problem

$T_1 = 2\bar{X} - 1$

$T_2 = M$

G_1, G_2, G_3, G_4 has the same distr.

$E(G_1) = \dots = E(G_4) = (N - 1)$

$G_1 + G_2 + G_3 + G_4 = 15 - 3$

$E(G_4) = \frac{1}{4}(15 - 3) = 3$

Generally, $E(G_{n+1}) = \frac{1}{n+1}(N - n)$

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

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

$M - N = -G_{n+1}$

$B_N(M) = E(M - N) = -E(G_{n+1})$

$= -3$

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

$E(T_3) = N$

Compare $MS E_{\theta}(T_2)$ & $MS E_{\theta}(T_3)$

$MS E_{\theta}(T_2) = Var(T_2) + B_0^2(T_2)$

$MS E_{\theta}(T_3) = Var(T_3)$

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

\approx but very close if n is large

T_3 has slightly larger variance but is unbiased