

Examples :

① $X_1, X_2, \dots, X_n \sim N(\mu, \sigma^2)$ iid

$\bar{X} \sim N(\mu, \frac{\sigma^2}{n})$ (even if n is small)

② $X_1, \dots, X_n \sim N(\mu_X, \sigma_X^2)$ iid

$Y_1, \dots, Y_m \sim N(\mu_Y, \sigma_Y^2)$ iid

$\bar{X} - \bar{Y} \sim N(\mu_X - \mu_Y, \frac{\sigma_X^2}{n} + \frac{\sigma_Y^2}{m})$ ← properties of the normal & linear combinations of normal r.v.

③ $X_1, \dots, X_n \sim \text{Bernoulli}(p_X)$ iid $\left\{ \begin{array}{l} n, m \text{ large} \\ Y_1, \dots, Y_m \sim \text{Bernoulli}(p_Y) \end{array} \right.$

$\bar{X} \approx N(p_X, \frac{p_X q_X}{n})$, $\bar{Y} \approx N(p_Y, \frac{p_Y q_Y}{m})$ $\begin{array}{l} E(X_i) = p_X \\ \text{Var}(X_i) = p_X q_X \end{array}$
 ↑
 approx by CLT

Uniform	$f(x) \geq 0$	$F(x) = \int_{-\infty}^x f(t) dt$
Exponential	$\int_{-\infty}^{\infty} f(x) dx = 1$	$F(x)$, $(0 \leq F(x) \leq 1)$
Normal		- continuous & increasing on $[-\infty, \infty)$

$P(a < X < b) = \int_a^b f(t) dt$

$E(X), \text{Var}(X), \text{SD}(X) = F(b) - F(a)$

Example 1

2 independent samples :

Sample 1 from Fresno, size 400

Sample 2 from Irvine size 500

We have average income from sample 1 is $\bar{X} = \$70,000$
sample SD = \$40,000

Sample 2: mean $\bar{Y} = \$80,000$

sample SD = \$50,000

$\text{Var}(\bar{X})$ can be estimated as $\frac{(40,000)^2}{400}$

$\text{Var}(\bar{Y})$ " " " $\frac{(50,000)^2}{500}$

$$\bar{X} - \bar{Y} = \$10,000$$

$$\text{SD}(\bar{X} - \bar{Y}) \approx \sqrt{\frac{(40,000)^2}{400} + \frac{(50,000)^2}{500}} = \$3,000$$

95% CI for the true difference in mean incomes b/w Irvine & Fresno \bar{X}

$$= (\bar{X} - \bar{Y}) \pm 2 \times \text{SD}(\bar{X} - \bar{Y})$$

$$= (-10,000 \pm 2 \times 3000) = -10,000 \pm 6000$$

$$= (-\$16,000, \$-4,000)$$

Example 2 . Perform a hypothesis test
on the true mean difference in incomes being 0.

$$H_0: \mu_X = \mu_Y$$

$$H_A \text{ or } H_1: \mu_X < \mu_Y \text{ or } \mu_Y - \mu_X > 0$$

If you use $\bar{X} - \bar{Y}$ & $\mu_X - \mu_Y < 0$ as your test



$$\bar{Y} - \bar{X} \underset{\text{by CLT}}{\sim} N(0, \text{Var}(\bar{Y} - \bar{X}))$$

p-value

$$\text{SD}(\bar{Y} - \bar{X}) \approx 3000 \text{ (estimated above)}$$

$$Z = \frac{(\bar{Y} - \bar{X}) - E(\bar{Y} - \bar{X} | H_0)}{\text{SD}(\bar{Y} - \bar{X})} = \frac{10,000 - 0}{3000} = \frac{10}{3}$$

$$p\text{-value} = P(Z > 10/3) = 1 - \Phi(10/3) \approx 0.00043 \approx 0.04\%$$

Reject the null hypothesis.

C.I. for proportions

2 indep samples from Fresno size 400 & Irvine size 500

proportion in Fresno sample that want Gov. Newsom recalled = 49%

proportion in Irvine sample that want CN recalled is 48%

$$\bar{X} = 0.49 \quad \text{Var}(\bar{X}) \approx \frac{(0.49)(0.51)}{400}$$

$$\bar{Y} = 0.48 \quad \text{Var}(\bar{Y}) \approx \frac{(0.48)(0.52)}{500}$$

$$\text{SD}(\bar{X} - \bar{Y}) = \sqrt{\text{Var}(\bar{X}) + \text{Var}(\bar{Y})} \approx 0.0335$$

$$\bar{X} - \bar{Y} = 0.49 - 0.48 = 0.01$$

95% C.I. for true difference in means (proportions)
(by CLT) $0.01 \pm 2 \times 0.0335$



$$= (-0.057, 0.077) = (-5.7\%, 7.7\%)$$


Duality b/w C.I & 2-sided hypothesis tests

tells us that we would not reject a null hyp of ^{true} mean difference being 0 at 5% significance level since 0 is in this 95% C.I.

Test for equality of proportions




H_0 : prop. of ^{Fresno} voters favoring recall = prop. of Irvine voters favoring recall.

$$p_x = p_y = p$$

H_1 : $p_x > p_y$ ($p_x - p_y > 0$ )

$$T = \bar{X} - \bar{Y}$$

H_1

- $p_x \neq p_y$ 
- $p_x > p_y$ 
- $p_x < p_y$ 

I need a value for the common value of p_x, p_y

Estimate this using total sample $\hat{p} = \frac{\text{total count}}{\text{total sample}}$ ^{estimate}
(weighted average of \hat{p}_x, \hat{p}_y)

$$\hat{p} = \frac{(0.49)(400) + (0.48)(500)}{400 + 500} \approx 0.4844$$

Under H_0 , $\bar{X} - \bar{Y} \sim N(0, \sigma^2)$

$$\sigma = \sqrt{\frac{pq}{400} + \frac{pq}{500}} \approx 0.0335$$

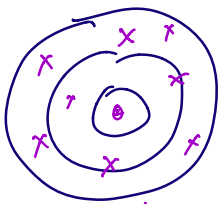
$$Z = \frac{(\bar{X} - \bar{Y}) - (0)}{SD(\bar{X} - \bar{Y})} = \frac{0.01 - 0}{0.0335}$$

$$P\text{-value} = P(Z > \frac{0.01 - 0}{0.0335}) \approx 0.382$$

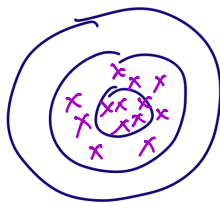
Fail to reject the Null.

Section 11.1 Bias & Variance

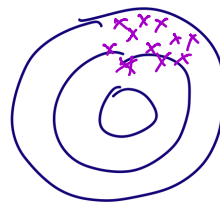
Data to estimate a population parameter
 θ : true value, T : estimate



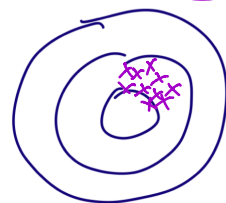
unbiased
high variance



unbiased
low variance



Biased
Low variance
precise (low variance)



low bias
low variance

unbiased
accurate

unbiased $E(T) = \theta$

biased $E(T) \neq \theta$