

STATS 412: 4 P.M. SKB 2500

- ~~HT/CIs~~ for population proportion
- ~~HT/CIs~~ for diff of population means
- ~~HT/CIs~~ for mean of population diff (same as 1-pop)
- ~~HT/CIs~~ for diff between population proportions
- "General Inference Topics" ??
→
- Correlation Coeff + LSRL
- Uncertainties in data
- Check assumptions + transform data

HYPOTHESIS TESTS FOR POP. PROPORTION:

- Pop. proportion \equiv pop. mean for 0s, 1s population

$$\rightarrow X \sim N(\mu = np, \sigma^2 = np(1-p)) \quad [\text{Binomial}]$$

$$\rightarrow \hat{p} = \frac{X}{n} \sim N(\mu = p, \sigma^2 = \frac{p(1-p)}{n})$$

→ USE Z-TEST

→ Assumptions:

- RANDOM SAMPLE (no interdependence)

- Approx Normal

→ # successes, # failures ≥ 10

→ $np \geq 10$; $n(1-p) \geq 10$

$$- z = \frac{\hat{p} - p_0}{\sqrt{p_0(1-p_0)/n}} \quad \leftarrow \text{std dev of original}$$

① Let p be the population proportion of all Gen Zers who would say that "... is a good thing for our society.

$$H_0: \hat{p} \leq 0.52 \quad \leftarrow \text{population proportion}$$

$$H_1: \hat{p} > 0.52$$

② We would like a RS from Bernoulli population and distribution of \hat{p} can be approximated w/ Normal dist.

③ R.S.: Not stated. We hope that any one individual sampled had no effect on any other individual's answer. Peer pressure may be a concern.

• 1178 is 5% of 23560. Reasonable that pop. size of Gen Zers is > 23560

④ Is Normal Approx Good?: $np = 1178(0.52) = 612.6 \geq 10$; $n(1-p) = 1178(1-0.52) = 565.4 \geq 10$ } then, \hat{p} approx Normal. So, can use z-procedures

⑤ Step 3: Calculate test statistic

$$z = \frac{\hat{p} - p_0}{\sqrt{p_0(1-p_0)/n}} = \frac{\frac{730}{1178} - 0.52}{\sqrt{(0.52)(0.48)/1178}} = 6.85 \quad \leftarrow \text{very large}$$

$$P\text{-value} = P[\hat{p} \geq 0.52]$$

$$= P[z \geq 6.85]$$

$$\approx 0 \quad \leftarrow \text{reject } H_0$$

Step 4: As we rejected H_0 , there is strong ($p \approx 0$) evidence to suggest that the population proportion of all Gen-Zers who would say that increased diversity is higher than 52%, which is the corresponding popul. prop of Gen X-ers.

C.I.s for POP. PROPORTION:

- AGRESTI-COULL ADJUSTMENT

$$\rightarrow \tilde{n} = n + 4 \quad ; \quad \tilde{p} = \frac{X + 2}{n + 4} \quad (\text{add 4 observations})$$

$$- \text{CI: } \hat{p} \pm z_{\alpha/2} \sqrt{\frac{\tilde{p}(1-\tilde{p})}{n}}$$

for 95% this would be 2.576

- Assumptions:

- ONLY RANDOM SAMPLE

- no need to check normality

HYPOTHESIS TESTS FOR DIFF OF POP MEANS ($\bar{x} - \bar{y}$) - not paired

- Assumptions:

- 2 independent R.S. from normally dist. pops
 - ① R.S. (for both)
 - ② independent samples
 - ③ normality
- } "sufficiently large"

- Use t-test unless both pop. stds known (rare)

$$T = \frac{(\bar{x} - \bar{y}) - \Delta_0}{\sqrt{\frac{S_x^2}{n_x} + \frac{S_y^2}{n_y}}} \quad (\text{df complicated, given an exam})$$

- CIs:

$$(\bar{x} - \bar{y}) \pm t_{n, \alpha/2} \sqrt{\frac{S_x^2}{n_x} + \frac{S_y^2}{n_y}}$$

POPULATION MEAN OF DIFFERENCES ($\overline{x-y}$) (paired)

- Assumptions:

→ RS from Normal pop (same as single pop)

- $T = \frac{\bar{D} - \Delta_0}{S_D / \sqrt{n}}$

- $\bar{d} \pm t_{n-1, \alpha/2} \frac{S_D}{\sqrt{n}}$

DIFFERENCE BETWEEN TWO PROPORTIONS:

• Assumptions:

- Independence between samples
- R.S.
- \sim Normal (10 succ/fail in each sample)

• POOLED PROPORTION:

$$\hat{p} = \frac{X + Y}{n_x + n_y}$$

$$Z = \frac{\hat{p}_x - \hat{p}_y}{\sqrt{\hat{p}(1-\hat{p})\left(\frac{1}{n_x} + \frac{1}{n_y}\right)}}$$

$$\hat{p}_x = \frac{X}{n_x}$$

$$\hat{p}_y = \frac{Y}{n_y}$$

$$\hat{p} = \frac{X + Y}{n_x + n_y}$$

• CONFIDENCE INTERVALS

→ No need to check sample sizes/normality

- AGRESTI - COULL

→ Add 4 samples still

→ 2 to each, 1 S / 1 f

$$\rightarrow \tilde{n}_x = n_x + 2 \quad \tilde{p}_x = \frac{X+1}{n_x+2}$$

$$\tilde{n}_y = n_y + 2 \quad \tilde{p}_y = \frac{Y+1}{n_y+2}$$

$$- (\tilde{p}_x - \tilde{p}_y) \pm z_{\alpha/2} \sqrt{\frac{\tilde{p}_x(1-\tilde{p}_x)}{n_x} + \frac{\tilde{p}_y(1-\tilde{p}_y)}{n_y}}$$

- Fixed-level testing: decide α before data

- TYPE I / TYPE II errors:

→ If H_0 is true and you reject → TYPE I

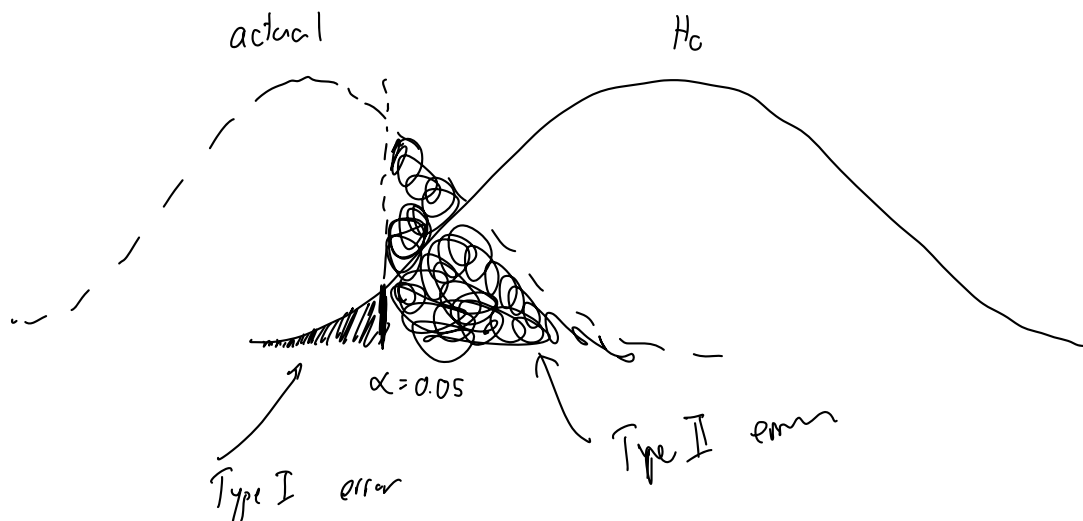
→ If H_0 is false and you accept → TYPE II

I: H_0 true

II: H_1 true

→ $P(\text{type I}) = \alpha$ (for single-sided test)

traditionally, we focus on minimizing Type I



As n increases, $P(\text{error}) \downarrow$

INTRO TO LINEAR REGRESSION:

① Form

② Direction

③ Strength

④ Outliers / Unusual

} describing scatterplot

$$\mu_y = \beta_0 + \beta_1 x$$

$$\hat{y} = \beta_0 + \beta_1 x$$

Residual: $e = \text{Actual} - \text{Pred}$
 $y_i - \hat{y}_i$ (negative \rightarrow overestimate)

- $b_1 = r \cdot \frac{s_y}{s_x}$

- $b_0 = \bar{y} - b_1 \bar{x}$

- Use Multiple R^2 from R output

- Do not extrapolate

- R^2 : ^(%) amount of variation explained by regression line

$$\frac{s_y^2 - s_{\text{resid}}^2}{s_y^2} = R^2 = 0.879$$

• TYPES OF OUTLIERS

- Non-leverage: weird y , but in x -range
- Leverage: Out of x -range
 - Influential: changes line (y weird)
 - Non-influential: does not change line (y fine for that x)

• Checking Conditions

- L: Linearity
- I: Independent Samples (no time series / time-pattern)
- N: Nearly - Normal Residuals (no outliers, look @ Q-Q plot)
- E: Equal Variability; variability around LSRL \sim same
(look @ residuals vs. fitted plot)
no fanning