

Since the Midterm: Inference

Given a population parameter, or a deterministic function of population parameters, we can test a statement about it or give an interval estimate for it.

↑
hypothesis test

↑
confidence interval

① Hypothesis test

- (i) proportion
- (ii) difference in proportions
- (iii) mean
- (iv) difference in means
- (v) slope in a linear model

② Confidence intervals

- (i) proportion
- (ii) difference in proportions
- (iii) mean
- (iv) difference in means
- (v) slope in a linear model
- (vi) the expected value of the response

③ Prediction intervals

④ Connection between hypothesis test and confidence interval

① Hypothesis test

② Settings

$$(i) H_0: \pi = \pi_0 \text{ vs } H_A: \pi \neq \pi_0$$

$$(ii) H_0: \pi_1 - \pi_2 = 0 \text{ vs } H_A: \pi_1 - \pi_2 \neq 0$$

$$(iii) H_0: \mu = \mu_0 \text{ vs } H_A: \mu \neq \mu_0$$

$$(iv) H_0: \mu_1 - \mu_2 = 0 \text{ vs } H_A: \mu_1 - \mu_2 \neq 0$$

$$(v) H_0: \beta_1 = 0 \text{ vs } H_A: \beta_1 \neq 0$$

③ Procedure

goal: to decide whether to reject H_0 or not based on a sample from the population, at significance level α

1. identifying test statistic and its null distribution

↑ function of the sample

desired properties

- We want it to contain useful info about the population parameters of interest

e.g. \bar{X} for $H_0: \mu = \mu_0$

\hat{p} for $H_0: \pi = \pi_0$

- We want its distribution to be known under H_0

(so that we can compute useful probabilities under H_0)

e.g. $\frac{\bar{X} - \mu_0}{S/\sqrt{n}} \sim N(0, 1)$ under H_0 (for setting (iii))

2. Computing observed test statistic

sample (random): $x_1, \dots, x_n \rightarrow \frac{\bar{x} - \mu_0}{S/\sqrt{n}} =: T_{\text{random}}$

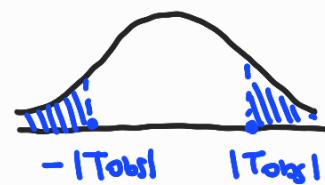
observed sample: $x_1, \dots, x_n \rightarrow \frac{\bar{x} - \mu_0}{S/\sqrt{n}} =: T_{\text{obs}}$

↑ e.g. for setting (iii)

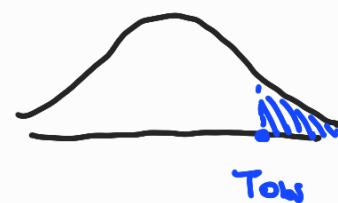
3. Knowing when to reject, given an observed sample

↳ by comparing the p-value with α
 (reject if the p-value is less than α)

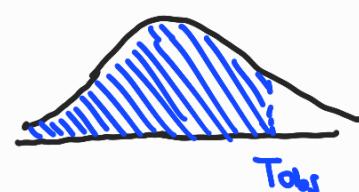
two-sided: $P_{H_0}(|T| \geq |T_{obs}|)$
 (e.g. $H_0: \mu = \mu_0$)
 $= P_{H_0}(T \geq |T_{obs}|) + P_{H_0}(T \leq -|T_{obs}|)$
 $= 2P_{H_0}(T \geq |T_{obs}|)$



one-sided,
 greater than: $P_{H_0}(T \geq T_{obs})$
 (e.g. $H_0: \mu > \mu_0$)



one-sided,
 less than: $P_{H_0}(T \leq T_{obs})$
 (e.g. $H_0: \mu < \mu_0$)



Summary of test statistics for different settings

| | | | |
|------------------------|-------------------------------|--------------------------|-------------------------------|
| (i) $H_0: \pi = \pi_0$ | (ii) $H_0: \pi_1 - \pi_2 = 0$ | (iii) $H_0: \mu = \mu_0$ | (iv) $H_0: \mu_1 - \mu_2 = 0$ |
|------------------------|-------------------------------|--------------------------|-------------------------------|

Binomial
test

$$X \stackrel{H_0}{\sim} \text{Bin}(n, \pi_0)$$

Z-test

$$\frac{\hat{\pi}_1 - \hat{\pi}_2}{\sqrt{\hat{\pi}(1-\hat{\pi})(\frac{1}{n_1} + \frac{1}{n_2})}}$$

pooling proportion
 If $\hat{\pi}_1 = \frac{x_1}{n_1}, \hat{\pi}_2 = \frac{x_2}{n_2}$,
 Then $\hat{\pi} = \frac{x_1 + x_2}{n_1 + n_2}$

t-test

$$\frac{\bar{X} - \mu_0}{S/\sqrt{n}}$$

random
 (sample standard deviation)

(equal variance)

$$\frac{\bar{X}_1 - \bar{X}_2}{S_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}}$$

(pooled sd)

(Welch)

$$\frac{\bar{X}_1 - \bar{X}_2}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}}$$

complicated

(paired)

$$\frac{\bar{D}}{S/\sqrt{n}}$$

$$D_i = X_{1,i} - X_{2,i}$$

S = sample standard deviation of $\{D_1, \dots, D_n\}$

$$(V) H_0: \beta_1 = \gamma$$

t-test

$$\frac{\hat{\beta}_1 - \gamma}{\text{SE}(\hat{\beta}_1)} \stackrel{H_0}{\sim} t_{n-2}$$

② Confidence intervals

- interval estimates for population parameter of interest

IP (the parameter is in the 1-d level) **(confidence interval)** = 1 - α

random (before we plug in values computed from the observed sample)

- no notion of null hypothesis or alternative hypothesis
we're not making any assumption on the parameter of interest

Terminology

(1 - α) - confidence interval for γ :

population parameter of interest

$$(PE - CV \cdot \hat{SE}, PE + CV \cdot \hat{SE})$$

MOE **MOE**

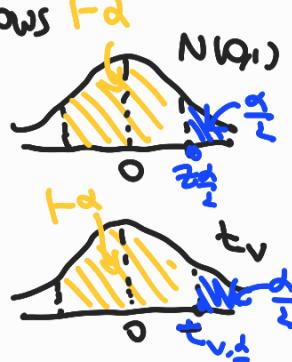
PE: point estimate (estimate for γ)

CV: critical value

↳ CV is found using confidence level ($= 1 - \alpha$)

and the distribution that the PE follows

- if $PE \sim N(0, \sigma^2)$ then $CV = z_{\frac{1-\alpha}{2}}$



\hat{SE} : estimate of the standard error of PE

↳ standard error of PE is often unknown, and thus we use an estimate of it

MOE: margin of error, $CV \cdot \hat{SE}$

Summary of CI for different settings

| | (i) π | (ii) $\pi_1 - \pi_2$ | (iii) M | (iv) $M_1 - M_2$ |
|----------------------|--|--|---|--|
| Z -CI (Wald CI) | PE: \hat{p} CV: $\frac{\sqrt{n}}{2\alpha}$ SE: $\sqrt{\frac{\hat{p}(1-\hat{p})}{n}}$ | PE: $\hat{p}_1 - \hat{p}_2$ CV: $\frac{\sqrt{n}}{2}$ SE: $\sqrt{\frac{\hat{p}_1(1-\hat{p}_1)}{n_1} + \frac{\hat{p}_2(1-\hat{p}_2)}{n_2}}$ | PE: \bar{x} CV: $t_{n-1, \frac{\alpha}{2}}$ SE: $\frac{s}{\sqrt{n}}$ | PE: \bar{x} CV: $t_{n-1, \frac{\alpha}{2}}$ SE: $\frac{s}{\sqrt{n}}$ |
| Agresti-Coull CI | Let $\hat{p} = \frac{x+2}{n+4}$ PE: $\hat{p} := \frac{x+2}{n+4}$ CV: $\frac{\sqrt{n}}{2\alpha}$ SE: $\sqrt{\frac{\hat{p}(1-\hat{p})}{n+4}}$ | Let $\hat{p}_1 = \frac{x_1+2}{n_1}$ and $\hat{p}_2 = \frac{x_2+2}{n_2}$ PE: $\hat{p}_1 - \hat{p}_2$, where $\hat{p} = \frac{\hat{p}_1 + \hat{p}_2}{2}$ CV: $\frac{\sqrt{n}}{2}$ SE: $\sqrt{\frac{\hat{p}_1(1-\hat{p}_1)}{n_1+2} + \frac{\hat{p}_2(1-\hat{p}_2)}{n_2+2}}$ | PE: \bar{x} CV: $t_{n-1, \frac{\alpha}{2}}$ SE: $\frac{s}{\sqrt{n}}$ | PE: \bar{x} CV: $t_{n-1, \frac{\alpha}{2}}$ SE: $\frac{s}{\sqrt{n}}$ |
| t -CI | | | PE: \bar{x} CV: $t_{n-1, \frac{\alpha}{2}}$ SE: $\frac{s}{\sqrt{n}}$ | PE: \bar{x} CV: $t_{n-1, \frac{\alpha}{2}}$ SE: $\frac{s}{\sqrt{n}}$ |
| | | | With Welch PE: $\bar{x}_1 - \bar{x}_2$ CV: $t_{\nu, \frac{\alpha}{2}}$ compared SE: $\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}$ | With Welch PE: $\bar{x}_1 - \bar{x}_2$ CV: $t_{n_1+n_2-2, \frac{\alpha}{2}}$ SE: $\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}$ |
| | | | (equal variance) PE: $\bar{x}_1 - \bar{x}_2$ CV: $t_{n_1+n_2-2, \frac{\alpha}{2}}$ SE: $\sqrt{\frac{s^2}{n_1+n_2}}$ | (equal variance) PE: $\bar{x}_1 - \bar{x}_2$ CV: $t_{n_1+n_2-2, \frac{\alpha}{2}}$ SE: $\sqrt{\frac{s^2}{n_1+n_2}}$ |

| (V) BI | | (VI) $E[Y_* X_*$] | |
|--------|-----------------------------|----------------------|--|
| PE | $\hat{\beta}_1$ | PE | $\hat{Y}_* \leftarrow \text{predicted response for } X_*$ |
| CV | $t_{n-2, \frac{\alpha}{2}}$ | CV | $t_{n-2, \frac{\alpha}{2}}$ |
| SE | $\hat{S.E.C}(\hat{\beta})$ | RE | $\hat{\sigma} = \sqrt{\frac{1}{n-2} \sum_{i=1}^n (Y_i - \hat{Y}_i)^2}$ $\hat{\sigma} = \sqrt{\frac{1}{n-2} \sum_{i=1}^n \frac{1}{n} + \frac{(X_i - \bar{X})^2}{\sum_{i=1}^n (X_i - \bar{X})^2}}$ |

* $E[Y_* | X_*]$ is a deterministic function of parameters

$$E[Y_* | X_*] = E[\beta_0 + X_* \beta_1 + \epsilon_* | X_*] \xrightarrow{\text{Ex} \sim N(0, \sigma^2)}$$

$$= \beta_0 + X_* \beta_1$$

↑ parameter ↑ parameter
 considered to be a fixed constant

③ Prediction intervals for Y_*

$$[\hat{Y}_* - t_{n-2, \frac{\alpha}{2}} \cdot \hat{\sigma} \sqrt{1 + \frac{1}{n} + \frac{(X_* - \bar{X})^2}{\sum_{i=1}^n (X_i - \bar{X})^2}}, \hat{Y}_* + t_{n-2, \frac{\alpha}{2}} \cdot \hat{\sigma} \sqrt{1 + \frac{1}{n} + \frac{(X_* - \bar{X})^2}{\sum_{i=1}^n (X_i - \bar{X})^2}}]$$

wider than the CI for $E[Y_* | X_*]$
due to the randomness in ϵ_*

$Y_* = \beta_0 + X_* \beta_1 + \epsilon_*$ is a random variable, not a deterministic function of parameters. This is why it's called a prediction interval instead of a confidence interval

④ connection between hypothesis test and confidence interval

Given an observed sample, we reject H_0 in a two-sided hypothesis test for a parameter at significance level α

$\Leftrightarrow (1-\alpha)$ level observed CI does not contain the value of the parameter specified in H_0 .

Other remarks

[Linear regression]

■ Linear Model

Given a sample $(X_1, Y_1), \dots, (X_n, Y_n)$, we assume

$$Y_i = \beta_0 + \beta_1 X_i + \varepsilon_i \quad \text{for } i \in \{1, \dots, n\},$$

where $\varepsilon_i \stackrel{\text{iid}}{\sim} N(0, \sigma^2)$.

■ Estimates of coefficients

$$\hat{\beta}_1 = r \cdot \frac{s_y}{s_x}, \quad \text{where } r \in [-1, 1] \text{ is the sample correlation coefficient}$$

$$\hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{X}$$

■ Predicted response for X_i :

$$\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 X_i$$

■ Residual for i^{th} observation

$$\hat{\varepsilon}_i = Y_i - \hat{Y}_i$$

Knowing how to read the R output is important,
e.g. what to find observed $\hat{\beta}_i$ and observed $\hat{SE}[\hat{\beta}_i]$.

[Normal vs t-distribution]



As the degrees of freedom of the t-distribution increases, the t-distribution approaches the normal distribution.

[Binomial distribution]

X is a binomial random variable if X can be thought of as the number of successes in a ^①**fixed number** of trials, where the trials are ^②**independent** and have the ^③**same success probability**.