

* Hypothesis testing (t-test)

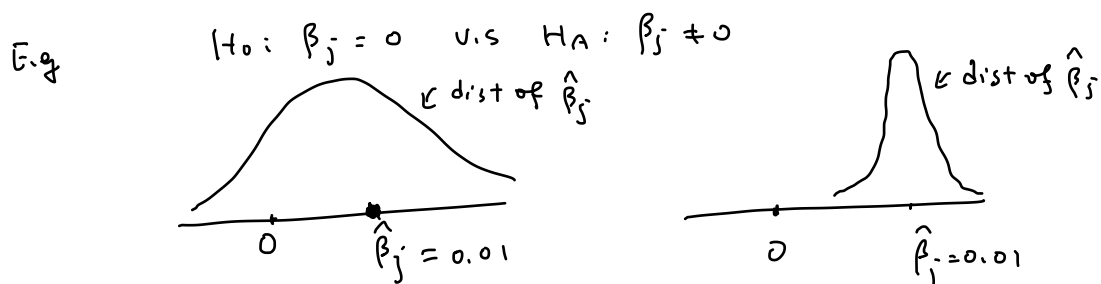
problem: for linear regression model with p predictors.

①

$$Y_i = \beta_0 + \beta_1 X_{i,1} + \beta_2 X_{i,2} + \dots + \beta_p X_{i,p} + e_i$$

Want to test $\underbrace{H_0: \beta_j = \beta_j^*}_{\text{null hypothesis}} \text{ v.s. } \underbrace{H_A: \beta_j \neq \beta_j^*}_{\text{Alternative Hypothesis}}, j \in \{1, \dots, p\}$

where β_j^* is prespecified, such as $\beta_j^* = 0$ for many applications



§1. * Review problem (from intro. to Stats.)

(*) $Y_1, \dots, Y_n \stackrel{\text{independent}}{\sim} N(\beta_0, \sigma^2)$

Want to test $H_0: \beta_0 = 0 \text{ v.s. } H_A: \beta_0 \neq 0$

Note that (*) is equivalent to the regression model with only intercept term in the mean function.

$i=1, \dots, n \quad Y_i = \beta_0 + e_i$

if $e_i \sim N(0, \sigma^2) \Rightarrow Y_i \sim N(\beta_0, \sigma^2)$

✓ Z-test (if σ^2 is known)

Z test statistic = $\frac{\bar{Y}}{\sigma/\sqrt{n}}$

\bar{Y} ← estimator of β_0

σ/\sqrt{n} ← s.e. of this estimator.

where $\bar{Y} = \frac{1}{n} \sum_{i=1}^n Y_i$

σ/\sqrt{n} is the standard error of the estimator \bar{Y} (s.e.)

\uparrow

$s.e.(\bar{Y}) = \sqrt{\text{Var}(\bar{Y})} = \sqrt{\frac{\sigma^2}{n}} = \frac{\sigma}{\sqrt{n}}$

practical problem

Generally
Recall that in Hypothesis testing, we want to control

$$\text{Type I error} = P(\text{Reject } H_0 \mid H_0 \text{ is true})$$

Smaller than significant level α
(\leq)

(typically $\alpha = 0.05$ in many applications)

while having a good statistical power

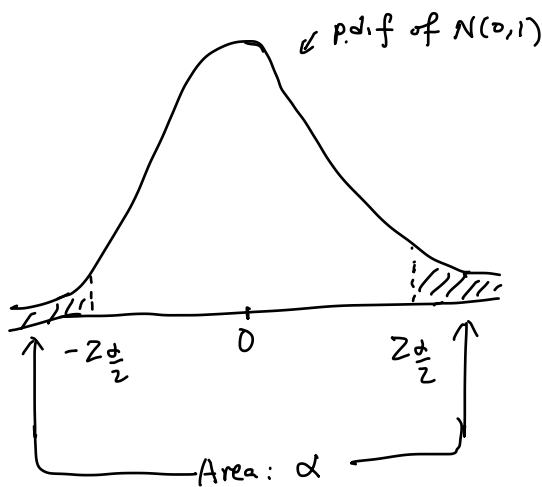
"
1 - Type II error

$$\text{with Type II error} = P(\text{Accept } H_0 \mid H_A \text{ is true})$$

For Z-test statistic, we have

When H_0 is true,

$$Z\text{-statistic} = \frac{\bar{Y} - \mu_0}{\sigma/\sqrt{n}} \stackrel{H_0}{\sim} N(0,1)$$



Test Procedure:

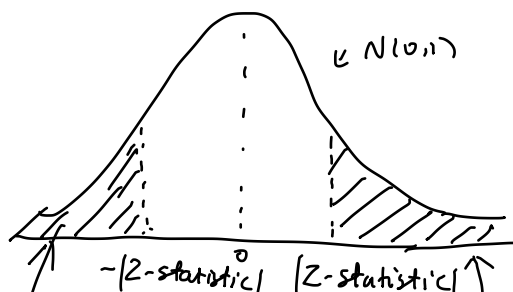
- ① compute Z -statistic from Data
- ② find a critical value $Z_{\frac{\alpha}{2}}$ such that $P(\text{Rej } H_0 \text{ if } |Z| > Z_{\frac{\alpha}{2}} \mid H_0 \text{ is true}) = \alpha$

If $\alpha = 5\%$, $Z_{\frac{\alpha}{2}} \approx 1.96$

- ③ Rej H_0 if $|Z\text{-statistic}| > Z_{\frac{\alpha}{2}}$
Accept H_0 otherwise.

or equivalently, we can use the p-value to perform the above test:

Replace ②&③ by ②' calculate p-value = $P(\underbrace{N(0,1)}_{\substack{\text{dist of test statistic} \\ \text{under } H_0}} > \underbrace{|Z\text{-statistic}|}_{\text{test statistic}} \text{ or } < -\underbrace{|Z\text{-statistic}|}_{\text{test statistic}})$



(computed from Data)
 Total Area: 100%
 p-value $\in [0, 1]$

③' Ref H_0 if p-value $< \alpha$
 Accept H_0 otherwise.

✓ t-test. (when σ^2 is unknown) ← estimator of β_0

$$t\text{-statistic} = \frac{\bar{y}}{\hat{\sigma}/\sqrt{n}} \quad \leftarrow \text{estimated s.e.}$$

Where $\hat{\sigma} = \sqrt{\frac{\sum_{i=1}^n (y_i - \bar{y})^2}{n-1}}$ ← sample variance of $\{y_1, \dots, y_n\}$.

Note that $\hat{\sigma}/\sqrt{n}$ is an estimated standard error of \bar{y} .

Under H_0 (H_0 is true),

$$t\text{-statistic} = \frac{\bar{y}}{\hat{\sigma}/\sqrt{n}} \underset{H_0}{\sim} t\text{-distribution}$$

With degrees of freedom (d.f.) $n-1$.

When $n \rightarrow \infty$ t-dist with $n-1$ d.f. $\rightarrow N(0, 1)$.

t-test \approx Z-test.

§2. Regression Setting.

mean function: $E[Y_i | x_i] = \beta_0 + \beta_1 x_{i,1} + \dots + \beta_j x_{i,j} + \dots + \beta_p x_{i,p}$

(problem ①) : $H_0: \beta_j = \beta_j^*$ v.s $H_A: \beta_j \neq \beta_j^*$.

H_0 : mean function $E[Y_i | x_i] = \beta_0 + \beta_1 x_{i,1} + \dots + \beta_j^* x_{i,j} + \dots + \beta_p x_{i,p}$

t-test Test Statistic:

$$t\text{-statistic} = \frac{\hat{\beta}_j - \beta_j^*}{\text{s.e.}(\hat{\beta}_j)}$$

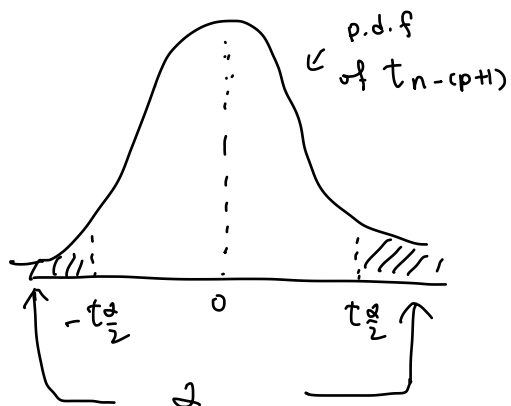
$\epsilon_{n \times 1} \sim N(0, \sigma^2 I)$
 \Updownarrow

When H_0 is true and statistical errors ($\epsilon_1, \dots, \epsilon_n$) are normally distributed.

the above t-statistic $\underset{H_0}{\sim}$ t-distribution with d.f. $n - (p+1)$

✓ Test Procedure: ① compute t-statistic from Data

② Find the critical value $t_{\frac{\alpha}{2}}$, such that



$$P(\text{Rej } H_0 \text{ if } |t\text{-stat}| > t_{\frac{\alpha}{2}} \mid H_0 \text{ is true})$$

$$\leq \alpha \quad (\text{e.g. } \alpha = 5\%)$$

$$t_{\frac{\alpha}{2}} \approx 2$$

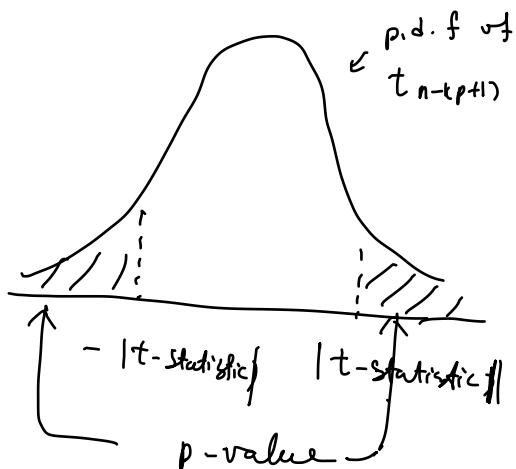
③ Rej H_0 if $|t\text{-statistic}| > t_{\frac{\alpha}{2}}$

Accept H_0 otherwise.

$$\uparrow \quad -t_{\frac{\alpha}{2}} \leftarrow \frac{\hat{\beta}_j - \beta_j^*}{\text{s.e.}(\hat{\beta}_j)} \leftarrow t_{\frac{\alpha}{2}}$$

Equivalently to ② & ③, we can use the following ②', ③'

②' we compute the p-value of the t-statistic



$$p \left(\begin{array}{l} \text{t-distribution with } n-(p+1) \text{ d.f.} \\ > |t\text{-statistic}| \text{ or } < -|t\text{-statistic}| \end{array} \right)$$

③' Rej H_0 if p-value $< \alpha$

Accept otherwise.

✓ ② & ③ are also equivalent to:

Accept H_0 if β_j^* (prespecified value) falls in the interval

$$\beta_j^* \in \left[\hat{\beta}_j - \underbrace{t_{\frac{\alpha}{2}} \cdot \text{s.e.}(\hat{\beta}_j)}_{\text{critical value in ②}}, \hat{\beta}_j + t_{\frac{\alpha}{2}} \cdot \text{s.e.}(\hat{\beta}_j) \right],$$

and Rej H_0 otherwise.

The above interval is the $(1-\alpha)$ -level Confidence Interval of β_j .
($j=0, 1, \dots, p$).

A $(1-\alpha)$ level Confidence Interval is defined as a range of values such that with $(1-\alpha)$ probability, the range will contain the true unknown value of β_j .