

→ Inference in Linear Regression

Recall that for the Least Square Estimates $\hat{\beta}_0$ and $\hat{\beta}_1$ we have:

$$E(\hat{\beta}_0) = \beta_0, \text{var}(\hat{\beta}_0) = \left(\frac{1}{n} + \frac{\bar{x}^2}{\sum (x_i - \bar{x})^2} \right) \sigma^2 \quad \text{and} \quad E(\hat{\beta}_1) = \beta_1, \text{var}(\hat{\beta}_1) = \frac{\sigma^2}{\sum (x_i - \bar{x})^2}$$

where $\sum_{i=1}^n (x_i - \bar{x})^2$ is commonly denoted as S_{xx} ("sum of squares")

Since the errors are normally distributed (i.e. $\varepsilon_i \sim N(0, \sigma^2)$) and $\hat{\beta}_0, \hat{\beta}_1$ are both linear combinations of normally distributed RV's, then $\hat{\beta}_0, \hat{\beta}_1$ are also normally dist.

i.e. $\hat{\beta}_0 \sim N(\beta_0, \text{var}(\hat{\beta}_0))$ and $\hat{\beta}_1 \sim N(\beta_1, \text{var}(\hat{\beta}_1))$

For instance, $\hat{\beta}_1 \sim N(\beta_1, \frac{\sigma^2}{S_{xx}})$ and $se(\hat{\beta}_1) = \sqrt{\sigma^2 / S_{xx}}$ (standard deviation)

Then $\frac{\hat{\beta}_1 - \beta_1}{\sqrt{\sigma^2 / S_{xx}}} \sim N(0, 1)$ (★) (standardization)

Turns out: (1) $\hat{\sigma}^2 = \frac{SSE}{n-2} = \frac{\sum (y_i - \hat{y}_i)^2}{n-2} = \frac{\sum (y_i)^2 - \hat{\beta}_0 \sum y_i - \hat{\beta}_1 \sum x_i y_i}{n-2} \sim \frac{\sigma^2 \chi_{n-2}^2}{n-2}$

(2) $\hat{\sigma}^2$ is independent of $\hat{\beta}_1$.

(★) $\frac{\hat{\beta}_1 - \beta_1}{\sigma / \sqrt{S_{xx}}} = (\hat{\beta}_1 - \beta_1) \cdot \frac{\sqrt{S_{xx}}}{\sigma} \sim N(0, 1)$ and we want to replace σ by $\hat{\sigma}$ (unknown estimate)

$$\Rightarrow (\hat{\beta}_1 - \beta_1) \cdot \frac{\sqrt{S_{xx}}}{\sigma} \cdot \frac{\sigma}{\hat{\sigma}} \left(\frac{\frac{\hat{\beta}_1 - \beta_1}{\sigma / \sqrt{S_{xx}}} \sim N(0, 1)}{\frac{\hat{\sigma}}{\sigma}} \sim \sqrt{\frac{\chi_{n-2}^2}{n-2}}} \right) \sim t_{n-2}$$

$$\Rightarrow \frac{\hat{\beta}_1 - \beta_1}{\hat{\sigma} / \sqrt{S_{xx}}} \sim t_{n-2}$$

⇒ Therefore, we will use the statistic $T = \frac{\hat{\beta}_1 - \beta_1}{se(\hat{\beta}_1)} \sim t_{n-2}$ for Regression inference.

• 2-sided CI for β_1 :

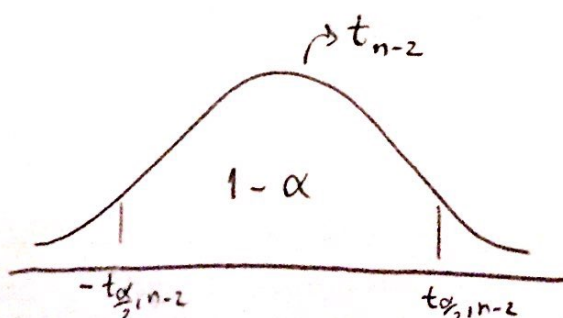
$$1 - \alpha = P(-t_{\alpha/2, n-2} \leq \frac{\hat{\beta}_1 - \beta_1}{se(\hat{\beta}_1)} \leq t_{\alpha/2, n-2})$$

$$= P(\hat{\beta}_1 - t_{\alpha/2, n-2} se(\hat{\beta}_1) \leq \beta_1 \leq \hat{\beta}_1 + t_{\alpha/2, n-2} se(\hat{\beta}_1))$$

• 1-sided CI for β_1 :

$$\beta_1 \in (-\infty, \hat{\beta}_1 + t_{\alpha, n-2} se(\hat{\beta}_1)) \text{ (Upper)}$$

$$\beta_1 \in (\hat{\beta}_1 - t_{\alpha, n-2} se(\hat{\beta}_1), \infty) \text{ (Lower)}$$

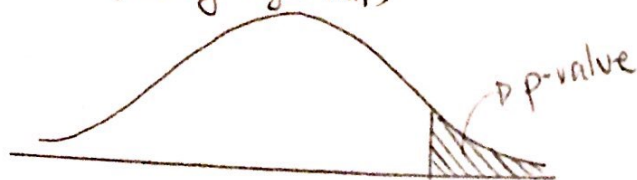


Understanding P-value

A P-value is just a "trigger" to decide whether you reject the null hypothesis or not.

Before continuing, understand the difference between a one-tailed and two-tailed hypothesis test.

One-tailed t-test
(Usually right tail)



$$H_0: \mu = 30$$

$$H_1: \mu > 30 \text{ (right tail)}$$

$$\text{or } H_1: \mu < 30 \text{ (left tail)}$$

For a Confidence level of:

Two-tailed t-test



$$H_0: \mu = 30$$

$$H_1: \mu \neq 30$$

(ie $\mu > 30$ or $\mu < 30$)

t-distribution

$(t_{\alpha, 1})$

One-tailed t-test

3.078

6.314

31.821

notice the #?

$(t_{\alpha/2, 1})$

Two-tailed t-test

± 6.314

± 12.706

± 63.657

df's

1

1

1

most t-tables will bring
upper tail prob. of the dist.
(ie. t s.t. $P(X > t) = 1 - \alpha$)

So, you construct your hypothesis testing first. Then, you calculate, based on your data, the probability of obtaining a result more "extreme" than those observed in your data assuming H_0 is true. This probability ($P(|t| > |t_{obs}|)$) is called the p-value. When it's very small, it means that the prob. of getting a result that confirms the null hypothesis is also very small, so you do not have enough evidence in the data to accept the null hypothesis (roughly speaking, you "reject" the null hypothesis based on your data).

In our previous example, p-value was 0.175.

Confidence Level

0.99

0.95

0.9

Accept/Reject H_0 ?

Accept
Accept
Accept
Reject

Notice that for a α of 0.99, I want strong evidence that H_0 is not true, whereas for a α of 0.8, I'm willing to reject H_0 if I have "weak" evidence that H_0 is not true.

Multiple Regression Example:

After running a multiple regression with 2 covariates (X_1 and X_2) in R, the following table was obtained ($n=20$).

Predictor	Coef.	SE coef.	T
constant	3.324	3.111	1.07
X_1	3.7681	0.6142	<input type="text"/>
X_2	5.0796	0.6655	<input type="text"/>

① Compute t-value corresponding to X_1 and X_2 .

Answer: $t_{X_1} = \frac{3.7681}{0.6142} = 6.13$ $t_{X_2} = \frac{5.0796}{0.6655} = 7.63$

② Test whether the predictor variables significantly affect Y at level 0.05.

$n=20$ so $t_{n-p-1} = t_{20-2-1} = t_{17}$

$t_{17, \frac{1-0.05}{2}} = t_{17, 0.975} = 2.109816$ (in R $qt(0.975, 17)$)

since both t_{X_1} and t_{X_2} are $>$ $t_{critical}$, both t_{X_1} and t_{X_2} are significant at 0.05 level

③ Using the model, estimate the value of Y at $X_1 = 5$ and $X_2 = 16$.

$$\hat{Y} = 3.324 + 3.7681(\underset{\substack{\downarrow \\ 5}}{X_1}) + 5.0796(\underset{\substack{\downarrow \\ 16}}{X_2}) = 103.44$$

... Coming next class (10/06) $\left. \begin{array}{l} \bullet \text{ CI for estimated average} \\ \bullet \text{ Prediction interval} \end{array} \right\}$