

taking log because they are thinking that the fit is non-linear

4 Testing Hypotheses about a Single Linear Combination of the Parameter

Consider

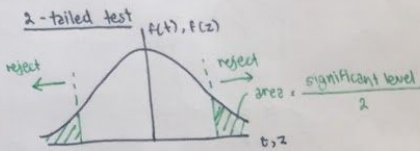
$$\log(\text{wage}) = \beta_0 + \beta_1 jc + \beta_2 \text{univ} + \beta_3 \text{exper} + u$$

where jc = number of years attending a two-year college
 $univ$ = number of years at a four-year college
 $exper$ = months in the workforce.

We want to test whether $\beta_1 = \beta_2$. \rightarrow if the return from 1 more year of education at a Junior college is the same as that of the university

$$H_0: \beta_1 = \beta_2 \rightarrow H_0: \beta_1 - \beta_2 = 0$$

$$H_a: \beta_1 \neq \beta_2 \rightarrow H_a: \beta_1 - \beta_2 \neq 0$$



$$t = \frac{(\hat{\beta}_1 - \hat{\beta}_2) - 0}{\text{s.e.}(\hat{\beta}_1 - \hat{\beta}_2)}$$

we compute this t-statistic and compare with the critical value

where $\text{s.e.}(\hat{\beta}_1 - \hat{\beta}_2) = \sqrt{\text{Var}(\hat{\beta}_1 - \hat{\beta}_2)}$

"not very straight forward to calculate"

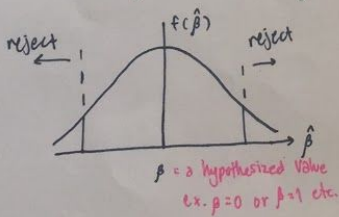
$$= \sqrt{\text{Var}(\hat{\beta}_1) + \text{Var}(\hat{\beta}_2) - 2 \text{Cov}(\hat{\beta}_1, \hat{\beta}_2)}$$

we use a variable transformation trick \rightarrow see notes

March 12 Inference \rightarrow hypothesis testing about " β " the true parameter

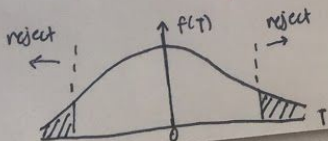
$$\text{wage} = \beta_0 + \beta_1 \text{educ} + \beta_2 \text{experience} + \dots + u$$

We want to test hypothesis about the true impact (β) of each X variables (educ, experience) on the dependent variable (y) BUT. we don't know what the true β are. So we use $\hat{\beta}$ (estimator) and $\text{se}(\hat{\beta})$ to test the hypothesis



① test if β = some number
 eg. $\beta_3 = 0 \rightarrow x_3$ has no impact only
 $\beta_3 = 1 \rightarrow 1$ unit \uparrow in x_3 correspond to 1 unit \uparrow in y

$$\rightarrow t\text{-test } \frac{\hat{\beta}_j - \beta_j}{\text{se}(\hat{\beta}_j)} \sim t_{d.f.}$$



significant level = total area in the rejection region

assume d.f. = 100
 $\text{area} = 2 \times (0.5 - 0.4803) = 2 \times 0.0197$ (FLT)
 $\text{p-value} \rightarrow 0.0394$

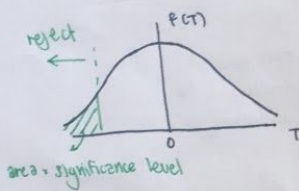
• suppose we calculate a t-statistic
 $\frac{\hat{\beta}_j - \beta_j}{\text{se}(\hat{\beta}_j)} = 2.06$
 • suppose, we are testing
 $H_0: \beta_j = 0$ $H_a: \beta_j \neq 0$
 \rightarrow 2-tailed test
 $t = 2.06 = \frac{(\hat{\beta}_j - \beta_j)}{\text{se}(\hat{\beta}_j)}$ \rightarrow p-value = total shaded area
 p-value = significant level which we will reject the H_0 or prob that will reject H_0
If the p-value < significance level \rightarrow reject H_0
 rule of thumb

another possible hypothesis test (one-tailed alternative)

$H_0: \beta_1 = \beta_2 \rightarrow H_0: \beta_1 - \beta_2 = 0$

$H_a: \beta_1 < \beta_2 \rightarrow H_a: \beta_1 - \beta_2 < 0$

It is assumed that β_1 would not be more than β_2 (returns to a 2-year college would never be more than returns to university education)



$$t = \frac{(\hat{\beta}_1 - \hat{\beta}_2) - 0}{s.e.(\hat{\beta}_1 - \hat{\beta}_2)}$$

"then go to extra note"

in class exercise

consider the multiple regression model, assume MLR 1-6 are satisfied

$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3 + u$

You would like to test the $H_0: \beta_1 - 3\beta_2 = 1$

1st) write the t-statistic for testing H_0

$$t = \frac{(\hat{\beta}_1 - 3\hat{\beta}_2) - 1}{s.e.(\hat{\beta}_1 - 3\hat{\beta}_2)}$$

2nd) Define $\theta_1 = \hat{\beta}_1 - 3\hat{\beta}_2 \rightarrow H_0: \theta_1 = 1, H_a: \theta_1 \neq 1$

$t = \frac{\hat{\theta}_1 - 1}{s.e.(\hat{\theta}_1)}$ → we need our regression to have θ_1 in it. So, STATA or OLS estimation will automatically give $\hat{\theta}_1$ & $s.e.(\hat{\theta}_1)$

Now $\hat{\beta}_1 = \hat{\theta}_1 + 3\hat{\beta}_2$

or $\beta_1 = \theta_1 + 3\beta_2$

sub in the main regression and get

$Y = \beta_0 + (\theta_1 + 3\beta_2)X_1 + \beta_2 X_2 + \beta_3 X_3 + u$

$= \beta_0 + \theta_1 X_1 + 3\beta_2 X_2 + \beta_2 X_2 + \beta_3 X_3 + u$

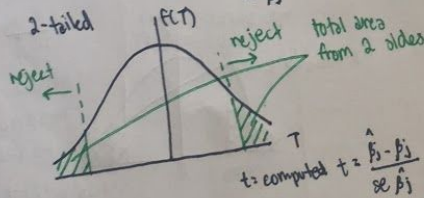
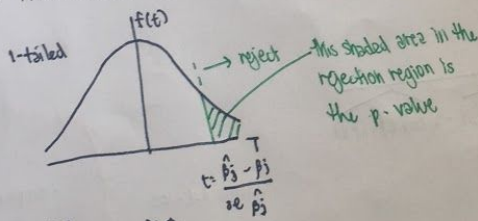
$= \beta_0 + \theta_1 X_1 + \beta_2 (X_2 + 3X_1) + \beta_3 X_3 + u$

* now, the explanatory variables are going to be $X_1, X_2 + 3X_1,$ and X_3

• we can calculate $t = \frac{\hat{\theta}_1 - 1}{s.e. \hat{\theta}_1}$

5 Computing p-Values for t-Tests

- What is the significance level given the computed t-statistics?

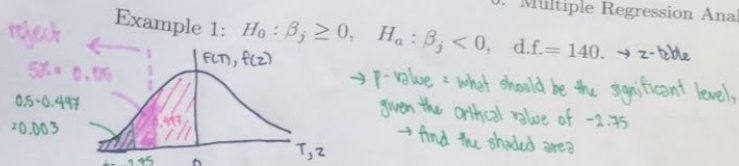


• p-value : $P(|T| > |t|)$

T = t-distributed random variable with d.f. = $n - k - 1$

t = computed t-statistic

→ P-value = probability that a random T value will be greater (in the $| |$ term) than our t in the H_0 test



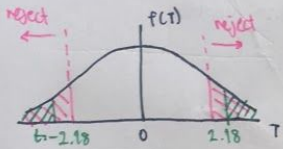
suppose the calculated $t_{\beta_j} = -2.75$

$$t_{\beta_j} = \frac{(\hat{\beta}_j - \beta_j)}{se(\hat{\beta}_j)}$$

- From the z-table, the value -2.75 corresponds to area = 0.003
- Thus, p-value = 0.003

- Would we reject H_0 if we use the significance level = 5%? **Yes**
 * Rule we reject H_0 if p-value < sig level

Example 2: $H_0: \beta_j = a_j, H_a: \beta_j \neq a_j, d.f. = 18 \leftarrow \text{use } t\text{-table}$



suppose the calculated $t_{\beta_j} = -2.18$

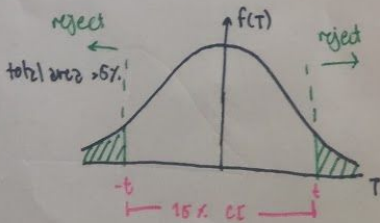
- From the t-table, the value -2.18 corresponds to area = 0.02 to 0.05
- Thus, p-value = is between 0.02-0.05

- Would we reject H_0 if we use the significance level = 5%?
 Yes, reject H_0 because the area is less than 0.05 or p-value < 0.05

6 Confidence Intervals (CI)

- Confidence Intervals for the POPULATION PARAMETER (β_j)

- A 95% CI of β_j is given by \rightarrow the range of values that would capture the true β_j at a 5% chance

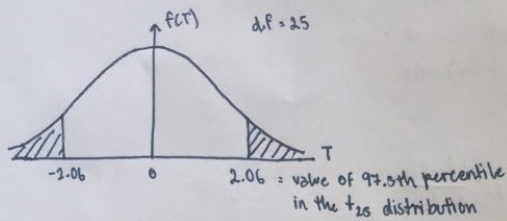


$$CI = \hat{\beta}_j \pm c \times se(\hat{\beta}_j)$$

c is the \square percentile in the t -distribution with $n-k-1$ d.f.

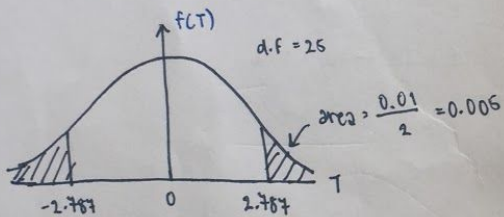
\uparrow 97.5 (in this case)

Example 1: 95% CI



The 95% CI for $\hat{\beta}_j = [\hat{\beta}_j - 2.06 \cdot se(\hat{\beta}_j), \hat{\beta}_j + 2.06 \cdot se(\hat{\beta}_j)]$

Example 2: 99% CI



The value of 99.5 percentile
or area = $1 - 0.005 = 0.995$

The 99% CI for $\hat{\beta}_j = [\hat{\beta}_j - 2.787 \cdot se(\hat{\beta}_j), \hat{\beta}_j + 2.787 \cdot se(\hat{\beta}_j)]$

7 Testing Multiple Linear Restrictions: The F-test

Suppose the model is specified by

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + u$$

$$H_0 : \beta_1 = 0 \text{ and } \beta_2 = 0 \rightarrow \text{want to test if } x_1 \text{ and } x_2$$

$$H_2, H_1 : H_0 \text{ is not true} \quad \text{BOTH have no impact on } y$$

We can use the F-test to test this type of "multiple hypotheses".

- Our full model is called the **"unrestricted"** model (ur). Suppose it can be expressed as:

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + u \text{ is true} \rightarrow \text{reject } H_0$$

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \dots + \beta_k x_k + u$$

- The model which takes out x (which we think its associated $\beta = 0$) is called the **restricted model** (r). ← small model

$$y = \beta_0 + \beta_1 x_1 + u \text{ is true} \rightarrow \text{don't reject } H_0$$

o suppose there are "q" number of β that we would like to perform a joint-test of = 0

eg. in this model $q=2$

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_{k-q} x_{k-q} + u$$

$$H_0 : \beta_{k-q+1} = \beta_{k-q+2} = \dots = \beta_k = 0$$

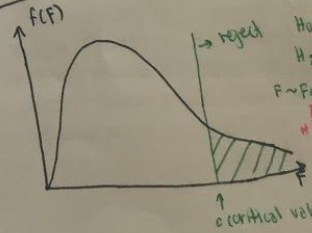
(the last q $\beta_s = 0$)

$$H_2 : H_0 \text{ is not true}$$

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_{k-q} x_{k-q} + \beta_{k-q+1} x_{k-q+1} + \beta_{k-q+2} x_{k-q+2} + \beta_k x_k + u$$

$$F \equiv \frac{(SSR_r - SSR_{ur}) / q}{SSR_{ur} / (n - k - 1)}$$

↑ this is always (+) because $SSR_r < SSR_{ur}$ everytime you add 1 more x , the model will be better explained
 ↑ d.f in ur model



$H_0 : \beta_2 = \beta_3 = \dots = 0$
 $H_2 : H_0 \text{ not true}$
 $F \sim F_{q, n-k-1}$ = d.f of ur
 if joint hypothesis being tested
 we reject H_0 if jointly no effect if $F > c$

• so, if everytime you add 1 more x variable, the $SSR \downarrow$ and $R^2 \uparrow$, why don't we just keep the additional x in the model??
 → because everytime we add 1 more x , the $\text{var}(\hat{\beta}_s)$ will increase making the prediction of β_0 less precise. so, we only keep the additional x s if it/they can improve the model enough → can $\downarrow SSR$ ($\uparrow R^2$) enough or can significantly $\downarrow SSR$ and $\uparrow R^2$

3. Some useful facts

① $R^2_{ur} > R^2_r$ because any additional x would increase R^2 (improve fit)
 $\Rightarrow SSR_{ur} < SSR_r$

② By including more x , the model is certainly better explained.
 However, we would like to reject H_0 if the inclusion of extra variables does not improve the model enough

4. Other ways to calculate the F-statistics:

\rightarrow From $R^2 = 1 - \frac{SSR}{SST}$

we have $F = \frac{(R^2_{ur} - R^2_r) / q}{(1 - R^2_{ur}) / (n - k - 1)}$
 (Annotations: q ← # of β that are set to "0", $n - k - 1$ ← intercept, n ← no. obs, k ← # of slopes)

\rightarrow If we want to test the overall significance of the model

$H_0: \beta_1 = \beta_2 = \beta_3 = \dots = \beta_k = 0$

H_2 : otherwise

$F = \frac{R^2 / k}{(1 - R^2) / (n - k - 1)}$
 (Annotations: R^2 of the model = ur , the " r " model has no x at all)

Example: Suppose we are interested in understanding the determinant of a baseball player's salary.

- y ← salary = season salary
- r ← [
 - years = years in major leagues
 - gamesyr = games per year in the league
 - bavg = career batting average
 - hrunsyr = homeruns per year
 - rbisyr = runs batted in per year
- ur ←

If we want to test whether performance has any impact on salary

$H_0: \beta_{bavg} = \beta_{hrunsyr} = \beta_{rbisyr} = 0$

H_2 : otherwise is true

- the unrestricted model (ur) is defined by

6. Multiple Regression Analysis (Inference) 73

ur model
 regress log_salary years gamesyr bavg hrunsyr rbisyr

Source	SS	df	MS			
Model	308.989208	5	61.7978416			
Residual	183.186327	347	.527914487			
Total	492.175535	352	1.39822595			

log_salary	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]
1 years	.0688626	.0121145	5.68	0.000	.0450355 .0926898
2 gamesyr	.0125521	.0026468	4.74	0.000	.0073464 .0177578
3 bavg	.0009786	.0011035	0.89	0.376	-.0011918 .003149
4 hrunsyr	.0144295	.016057	0.90	0.369	-.0171518 .0460107
5 rbisyr	.0107657	.007175	1.50	0.134	-.0033462 .0248776
_cons	11.19242	.2888229	38.75	0.000	10.62435 11.76048

Number of obs = 353
 F(5, 347) = 117.06
 Prob > F = 0.0000
 R-squared = 0.6278
 Adj R-squared = 0.6224
 Root MSE = .72658

the restricted model (r) is defined by

regress log_salary years gamesyr

Source	SS	df	MS			
Model	293.864058	2	146.932029			
Residual	198.311477	350	.566604221			
Total	492.175535	352	1.39822595			

log_salary	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]
years	.071318	.012505	5.70	0.000	.0467236 .0959124
gamesyr	.0201745	.0013429	15.02	0.000	.0175334 .0228156
_cons	11.2238	.108312	103.62	0.000	11.01078 11.43683

when considering each of the performance X one by one, none of them has a significant impact at 5%.

Number of obs = 353
 F(2, 350) = 259.32
 Prob > F = 0.0000
 R-squared = 0.5971
 Adj R-squared = 0.5948
 Root MSE = .75273

Now, our H_0 and H_a becomes

$$F \equiv \frac{(SSR_r - SSR_{ur}) / q}{SSR_{ur} / (n - k - 1)}$$

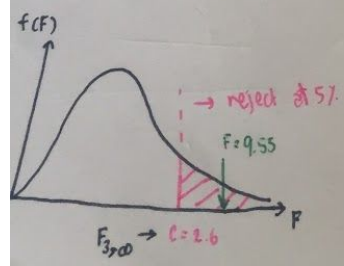
$$\equiv \frac{(198.311 - 183.186) / 3}{(183.186) / (353 - 5 - 1)} \approx 9.55$$

HW

$$F = \frac{(R^2 / q)}{(1 - R^2) / (n - k - 1)}$$

= ??

But when performing on F-tests, performance have joint impact



use 5% level of sig.
 since $F = 9.55 > 2.6$, we reject H_0 at 5% level and conclude that performances have joint effects on salary

8 How the Hypothesis Testing is done in Practice

1. Check the values of t -statistic reported by the statistical software (i.e. STATA, SPSS, SAS)

⇒ These t -statistics are to test $H_0: \beta_i = 0$

⇒ If the d.f. > 30 , then when $t > 1.96$, we can reject H_0 with 5% sig. level

⇒ When $t > 1.96$, we can say that β_i is statistically significant at 5% level. (value of $\beta_i \neq 0$)

⇒ When $t < 1.96$ we can say that β_i is not statistically significant at 5% level.

⇒ If $t < 1.96$ we can drop x_i from the model

⇒ After we drop x_i , we estimate the new regression function and obtain a new set of $\hat{\beta}$.

2. We can also perform other hypothesis testings of interest.

e.g. $H_0: \beta_i = \beta_j$

or $H_0: \beta_i = 5$ etc.

or perform an F-test for testing multiple linear restrictions

3. Usually, in economics, the estimation results are reported using this form

Dependent Variable: log(salary)			
Independent Variables	(1)	(2)	(3)
log(sales)	.224 (.027)	.158 (.040)	.188 (.040)
log(mktval)	—	.112 (.050)	.100 (.049)
profmarg	—	-.0023 (.0022)	-.0022 (.0021)
ceoten	—	—	.0171 (.0055)
comten	—	—	-.0092 (.0033)
intercept	4.94 (0.20)	4.62 (0.25)	4.57 (0.25)
Observations	177	177	177
R-squared	.281	.304	.353

↑
like a simple regression
with 1x