

# ECON 3740: INTRODUCTION TO ECONOMETRICS

INSTRUCTOR: CHAOYI CHEN  
Department of Economics and Finance, University of Guelph

## Lecture 14

Last lecture, we studied the the hypothesis test with a single parameter.  
Today, we will

- Construct confidence interval
- Test hypotheses about a single linear combination of the parameters
- Test multiple linear restrictions
  - Test exclusion restrictions
  - Unrestricted model and restricted model
  - $F$  test
  - Test overall significance
  - Test general linear restrictions

# MLR: Computing $p$ -Values for $t$ Tests

- Recall that  $\frac{\hat{\beta}_j - \beta_j}{se(\hat{\beta}_j)} \sim t_{n-k-1}$ . Hence,

$$P\left(\frac{\hat{\beta}_j - \beta_j}{se(\hat{\beta}_j)} > c_{0.05}\right) = 0.025$$

$$P\left(\frac{\hat{\beta}_j - \beta_j}{se(\hat{\beta}_j)} < -c_{0.05}\right) = 0.025$$

$$P(c_{0.05} \leq \frac{\hat{\beta}_j - \beta_j}{se(\hat{\beta}_j)} \leq c_{0.05}) = 0.95$$

where  $c_{0.05}$  is the 5% critical value of two-sided test

- Simple manipulation of above result shows

$$\begin{aligned} P\left(\underbrace{\hat{\beta}_j - c_{0.05} * se(\hat{\beta}_j)}_{\text{lower bound of the CI}} \leq \beta_j \leq \underbrace{\hat{\beta}_j + c_{0.05} * se(\hat{\beta}_j)}_{\text{upper bound of the CI}}\right) \\ = P\left(\left|\frac{\hat{\beta}_j - \beta_j}{se(\hat{\beta}_j)}\right| \leq c_{0.05}\right) = 0.95 \end{aligned}$$

- The confidence interval is  $[\hat{\beta}_j - c_{0.05} * se(\hat{\beta}_j), \hat{\beta}_j + c_{0.05} * se(\hat{\beta}_j)]$ . 0.95 is called the **confidence level**.

# MLR: Confidence Interval - an example

- The fitted regression line is

$$\widehat{\log(rd)} = \begin{matrix} -4.38 \\ (0.47) \end{matrix} + \begin{matrix} 1.084\log(\text{sales}) \\ (0.06) \end{matrix} + \begin{matrix} 0.0217\text{profmarg} \\ (0.0128) \end{matrix} \quad \text{where}$$

$rd$  = firms spending on RD

$\text{sales}$  = annual sales

$\text{profmarg}$  = profits as percentage of sales

- $df = n - k - 1 = 32 - 2 - 1 = 29$ . Hence,  $c_{0.05} = 2.045$ .
- The 95% CI for  $\beta_{\log(\text{sales})}$  is  $[1.084 - 2.045 * (0.06), 1.084 + 2.045 * (0.06)] = [0.961, 1.21]$ . The effect of  $\log(\text{sales})$  on  $\log(rd)$  is relatively precisely estimated as the interval is narrow. Moreover, the effect is significantly different from zero because zero is outside the interval.
- The 95% CI for  $\beta_{\text{profmarg}}$  is  $[0.0217 - 2.045 * 0.0128, 0.0217 + 2.045 * 0.0128] = [-0.0045, 0.0479]$ . The effect of  $\text{profmarg}$  on  $\log(rd)$  is imprecisely estimated as the interval is very wide. It is not even statistically significant because zero lies in the interval.

# MLR: Testing Hypotheses about a Single Linear Combination of the Parameters: An Example

- To investigate whether returns to education at 2-Year vs. at 4-Year colleges are equal or not, one propose a model as following

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

where

$jc$ =years years of education at 2-year colleges,  $univ$ = at 4-year colleges

- Suppose we want to test

$$H_0 : \beta_1 - \beta_2 = 0 \text{ vs } H_1 : \beta_1 - \beta_2 < 0$$

- A possible test statistic would be

$$t = \frac{\hat{\beta}_1 - \hat{\beta}_2}{se(\hat{\beta}_1 - \hat{\beta}_2)}$$

- However, here appear some validity problems relating to above  $t$  statistic. First, we know standardized  $\hat{\beta}_1$  follows  $t$  distribution, and standardized  $\hat{\beta}_2$  follows  $t$  distribution. But these cannot make sure that standardized  $\hat{\beta}_1 - \hat{\beta}_2$  will also follow  $t$  distribution. Secondly, assuming standardized  $\hat{\beta}_1 - \hat{\beta}_2 \sim t_{n-k-1}$ , it is still impossible to compute such a  $t$  statistic since one cannot compute  $se(\hat{\beta}_1 - \hat{\beta}_2)$ , which relates to the  $\widehat{cov}(\hat{\beta}_1, \hat{\beta}_2)$  and is usually not available in regression output.

# MLR: Testing Hypotheses about a Single Linear Combination of the Parameters: An Example Continue

- **An Alternative Method:** Define  $\theta_1 = \beta_1 - \beta_2$ . Therefore, the hypothesis test can be rewritten as

$$H_0 : \theta_1 = 0 \text{ vs } H_1 : \theta_1 < 0$$

- Now,  $\beta_1 = \theta_1 + \beta_2$ . Inserting it into the original regression, we have

$$\begin{aligned} \log(\text{wage}) &= \beta_0 + (\theta_1 + \beta_2)jc + \beta_2univ + \beta_3exper + \mu \\ &= \beta_0 + \theta_1jc + \beta_2(jc + univ) + \beta_3exper + \mu \end{aligned}$$

where  $jc + univ$  is a new regressor, representing total years of college.

# MLR: Testing Hypotheses about a Single Linear Combination of the Parameters: An Example Continue

- After estimation, the fitted regression line is

$$\log(\widehat{wage}) = \begin{matrix} 1.472 & -0.0102jc & +0.0769(jc+univ) & +0.0049exper \\ (0.021) & (0.0069) & (0.0023) & (0.0002) \end{matrix}$$

where  $n = 6763$  and  $R^2 = 0.222$

- Hence, the  $t$  statistic is

$$t = \frac{-0.0102}{0.0069} = -1.48$$

- Recall that the critical values with significance level 5% and 10% are -1.645 and -1.282 respectively for the left tail one sided test.
- Thus, the null is rejected at 10% level but not at 5% level.
- You can also compute the  $p$  value, which gives you  $P(T < -1.48) = 0.07 \in (0.05, 0.10)$ , and the 95% CI for  $\theta_1$  is  $-0.102 \pm 1.96 * (0.0069) \rightarrow (-0.0237, 0.0003)$ , which covers zero.
- Note that this method works always for **single** linear hypotheses.

# MLR: Testing Multiple Linear Restrictions: Testing Exclusion Restrictions

- The estimated restricted model is

$$\log(\text{salary}) = \beta_0 + \beta_1 \text{years} + \beta_2 \text{gamesyr} + \beta_3 \text{bavg} + \beta_4 \text{hrunsyr} + \beta_5 \text{rbisyr} +$$

where

*salary* = the 1993 total salary

*years* = years in the league

*gamesyr* = average games played per year

*bavg* = career batting average

*hrunsyr* = home runs per year

*rbisyr* = runs batted in per year

- Now, we consider following hypothesis,

$$H_0 : \beta_3 = \beta_4 = \beta_5 = 0 \text{ vs } H_1 : H_0 \text{ is not true}$$

where  $H_0$  is not true means "at least one of  $\beta_3, \beta_4, \beta_5$  is not zero.

- This is to test whether performance measures have **no** effects or can be **excluded** from regression.



# MLR: Testing Multiple Linear Restrictions: Estimation of the Unrestricted Model

- The estimated **unrestricted** model is

$$\begin{aligned} \log(\widehat{\text{salary}}) = & \quad 11.19 & \quad -0.0689\text{years} & \quad +0.012\text{gamesyr} \\ & (0.29) & (0.0121) & (0.0026) \\ & +0.00098\text{bavg} & +0.0144\text{hrunsyr} & +0.0108\text{rbisyr} \\ & (0.0011) & (0.0161) & (0.0072) \end{aligned}$$

where  $n = 353$ ,  $SSR_{ur} = 181.186$ , and  $R^2 = 0.6278$

- Note that none of these three variables is statistically significant when tested individually. (Why? compute the corresponding  $t$  statistic). However, the individual insignificance may **not** imply together they are insignificant.
- Idea**: How would the model fit (measured in  $SSR$ ) be if these variables were dropped from the regression?

# MLR: Testing Multiple Linear Restrictions: Estimation of the restricted Model

- By dropping those variables from the model, The estimated **restricted** model is  $\widehat{\log(\text{salary})} = 11.22 + 0.0713\text{years} + 0.202\text{gamesyr}$   
(0.11) (0.0125) (0.0013)  
where  $n = 353$ ,  $SSR_r = 198.311$ , and  $R^2 = 0.5971$
- The sum of squared residuals (SSR) necessarily increases in the restricted model. But, is this increase statistically significant?
- To figure out the problem, we consider to use a rigorous test statistic

$$F = \frac{(SSR_r - SSR_{ur})/q}{SSR_{ur}/(n - k - 1)} \sim F_{q, n-k-1}$$

where  $q = df_r - df_{ur}$  is the number of restrictions, and  $n - k - 1 = df_{ur}$ .

- The relative increase of the sum of squared residuals when going from  $H_1$  (unrestricted model) to  $H_0$  (restricted model) follows a  $F$  distribution.

# MLR: Testing Multiple Linear Restrictions: F test

- Therefore, the  $F$  statistic of our example is

$$F = \frac{(198.311 - 181.186)/3}{181.186/(353 - 5 - 1)} \approx 9.55$$

- Check with the  $F$  table, we find if  $q = 3$ ,  $df = 347$ , the critical value with 1% significance level is 3.78.
- $9.55 > 3.78$ . Therefore, we reject the null.
- Alternatively, you can compare  $p$  value.  $P(F_{3,347} > 9.55) = 0.0000$ , which implies the null hypothesis is overwhelmingly rejected (even at very small significance levels).
- **Remarks:**
  - If  $H_0$  is rejected, we say that the three variables are "jointly significant".
  - They were not significant when tested individually.
  - The possible reason is multicollinearity between them.

# MLR: Testing Multiple Linear Restrictions: The R-Squared Form of the F Statistic

- Recall that

$$R^2 = -\frac{SSR}{TSS} \implies SSR = TSS(1 - R^2)$$

- Hence,

$$\begin{aligned} F &= \frac{(SSR_r - SSR_{ur})/q}{SSR_{ur}/(n - k - 1)} = \frac{[TSS(1 - R_r^2) - TSS(1 - R_{ur}^2)]/q}{TSS(1 - R_{ur}^2)/(n - k - 1)} \\ &= \frac{(R_{ur}^2 - R_r^2)/q}{(1 - R_{ur}^2)/(n - k - 1)} \end{aligned}$$

- With previous example, since  $R_{ur}^2 = 0.6278$ ,  $R_r^2 = 0.5971$ , therefore,

$$F = \frac{(0.6278 - 0.5971)/3}{(1 - 0.6278)/347} \approx 9.54$$

which is very close to the result based on SSR (difference due to rounding error).

# MLR: Testing Multiple Linear Restrictions: The F Statistic for Overall Significance of a Regression

- Consider a typical population regression model

$$y = \beta_0 + \beta_1 x_1 + \dots + \beta_k x_k + \mu$$

- Now, suppose we would like to conduct a special hypothesis test

$$H_0 : \beta_1 = \dots = \beta_k = 0 \text{ vs } H_1 : H_0 \text{ is not true}$$

- As a result, the restricted model is

$$y = \beta_0 + \mu$$

which is a regression on constant. Clearly,  $\hat{\beta}_0 = \bar{y}$  and  $R_r^2 = 0$  from the knowledge in SLR.

- Thus, the  $F$  statistic of this special hypothesis test is

$$F = \frac{(R_{ur}^2 - R_r^2) / q}{(1 - R_{ur}^2) / (n - k - 1)} = \frac{R_{ur}^2 / k}{(1 - R_{ur}^2) / (n - k - 1)} \sim F_{k, n-k-1}$$

- The test of overall significance is reported in most regression packages (also in 'lm' package in R). The null hypothesis is usually overwhelmingly rejected.

# MLR: Testing Multiple Linear Restrictions: Testing General Linear Restrictions

- Suppose you and your group member would like to test whether house price assessments are rational, where the population regression model is

$$\log(\text{price}) = \beta_0 + \beta_1 \log(\text{assess}) + \beta_2 \log(\text{lotsize}) \\ + \beta_3 \log(\text{sqrft}) + \beta_4 \text{bdrms} + \mu$$

where

*price* = house price

*assess* = the assessed housing value (before the house was sold)

*lotsize* = size of the lot, in feet

*sqrft* = square footage

*bdrms* = number of bedrooms

- Now, suppose we focus on the following hypothesis test

$$H_0 : \beta_1 = 1, \beta_2 = \beta_3 = \beta_4 = 0$$

# MLR: Testing Multiple Linear Restrictions: Testing General Linear Restrictions

- Under the null,  $\beta_1 = 1$ , which means that if house price assessments are rational, a 1% change in the assessment should be associated with a 1% change in price. (log-log model). Also,  $\beta_2 = \beta_3 = \beta_4 = 0$ , which means that in addition, other known factors should not influence the price once the assessed value has been controlled for.
- Since it is a test involving multiple linear restrictions, we consider to use  $F$  test.
- The restricted model is

$$y = \beta_0 + x_1 + \mu$$

- However, if you just regress  $y$  on  $x_1$ , we cannot ensure the restriction  $\beta_1 = 1$  hold. Hence, we consider to change the restricted model as  $y - x_1 - \beta_0 + \mu$ , which means The restricted model is actually a regression of  $y - x_1$  on a constant, and the resulting  $\hat{\beta}_0$  is the sample mean of  $y - x_1$ .

# MLR: Testing Multiple Linear Restrictions: Testing General Linear Restrictions Continue

- Suppose that after estimation with both unrestricted model restricted model, you have the following results

$$SSR_r = 1.88, \quad SSR_{ur} = 1.822, \quad n = 88$$

- Hence, the  $F$  statistic is

$$F = \frac{(SSR_r - SSR_{ur})/q}{SSR_{ur}/(n - k - 1)} = \frac{(1.88 - 1.822)/4}{1.822/(88 - 4 - 1)} \approx 0.661$$

- Checking the  $F$  table, you find the critical value with 5% significance level of a  $F_{4,83}$  distribution is 2.5, which is greater than 0.661. Therefore, we cannot reject the null