

# Lecture 6: Further Inference in the Multiple Regression Model

Shuo Liu

UCLA Summer School Econ 103

July 18, 2017

1 Joint Hypothesis Testing

2 Model Specification

# Joint Hypothesis Testing

Joint Hypothesis Testing tests a null hypothesis with **multiple conjectures**, expressed with more than one “equal sign”;

- **Example:** should a group of explanatory variables  $\{x_3, x_4, x_5\}$  be included in a particular model?

$$y = \beta_1 + \beta_2 x_2 + \beta_3 x_3 + \beta_4 x_4 + \beta_5 x_5 + e \quad (1)$$

- **Test Form:**

$$H_0 : \beta_3 = 0, \beta_4 = 0, \beta_5 = 0 \quad (2)$$

$$H_1 : “\beta_3 = 0, \beta_4 = 0, \beta_5 = 0” \text{ do not hold simultaneously} \quad (3)$$

- A joint test for whether all the three conjectures hold simultaneously

# Joint Hypothesis Testing

SALES example:

- Consider the model:

$$SALES = \beta_1 + \beta_2 PRICE + \beta_3 ADVERT + \beta_4 ADVERT^2 + e \quad (4)$$

- Test whether or not advertising has an effect on sales: If advertising does not have effect

$$H_0 : \beta_3 = 0, \beta_4 = 0 \quad (5)$$

- If advertising has effect on sales

$$H_1 : “\beta_3 \neq 0, \beta_4 = 0” \text{ or } “\beta_3 = 0, \beta_4 \neq 0” \text{ or } “\beta_3 \neq 0, \beta_4 \neq 0” \quad (6)$$

- Relative to the null hypothesis  $H_0$ , the original model (4) is called **unrestricted model**, where the restrictions in  $H_0$  have not been imposed on the original model (4);
- **Restricted model**: impose restrictions in  $H_0$  to the original model (or assume parameter restrictions in  $H_0$  are true)

$$SALES = \beta_1 + \beta_2 PRICE + e \quad (7)$$

# Joint Hypothesis Testing

How to determine which model is better (a little bit like choosing from two specific model forms)?

- We use  $SSE = \sum_{i=1}^n \hat{e}_i^2 = \sum_{i=1}^n (y_i - \hat{y}_i)^2$ , where  $\hat{y}_i$  is determined by the model form you choose (either restricted form or unrestricted form);
- We compare sum of squares of error (residuals) from unrestricted model  $SSE_U$  with that from restricted model  $SSE_R$ . Intuitively, we will choose unrestricted model form ( $H_1$ ) only when  $SSE_U$  is “too smaller” than  $SSE_R$ , because:

$$SST = SSR + SSE \quad (8)$$

and when throwing in more explanatory variables into the model,  $SSR$  will always increase, then we always have  $SSE_U \leq SSE_R$ .

- We use F test to test  $H_0 : \beta_3 = 0, \beta_4 = 0$ .

# Joint Hypothesis Testing

- Remember F test is always one-tail (right-tail) test!
- **First step:** calculate F-statistic

$$F = \frac{(SSE_R - SSE_U)/J}{SSE_U/(n - K)} \quad (9)$$

where

1.  $J$  = number of restrictions;
  2.  $n$  is sample size and  $K$  is number of parameters in the original model (unrestricted model).
- **Second step:** determine the distribution of F-statistic above under  $H_0$  is true

$$F \sim F_{(J, n-K)} \quad (10)$$

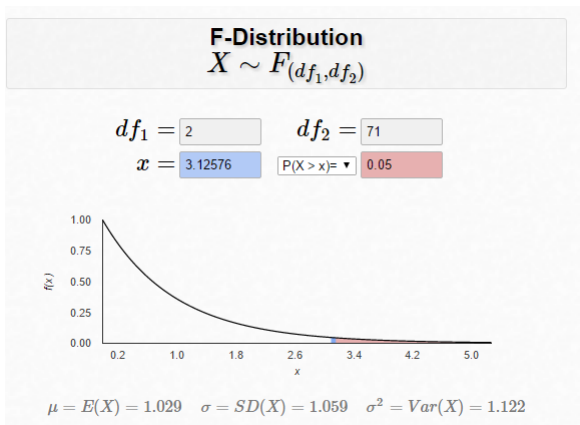
F distribution with  $J$  degree of freedom in the numerator and  $n - K$  degree of freedom in the denominator;

- Next, given significance level, we will reject  $H_0$  only when F-statistic is “too large”, i.e.  $SSE_U$  is “too smaller” than  $SSE_R$ , which means unrestricted model is much better than restricted model.

# Joint Hypothesis Testing

- **Third step:** given significance level  $\alpha$

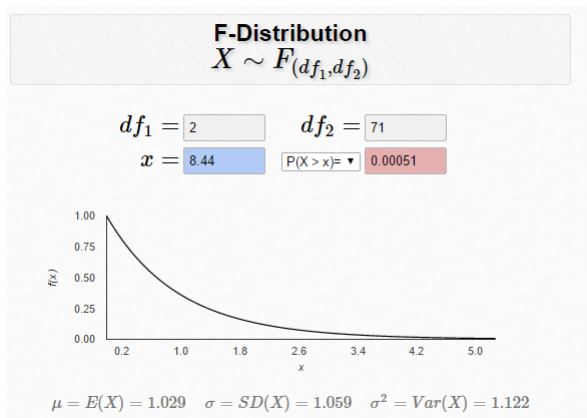
**Option 1:** find critical value with right tail probability equal to  $\alpha$  and set rejection region;



# Joint Hypothesis Testing

- **Third step:** given significance level  $\alpha$

**Option 2:** use F-statistic calculated in the first step to calculate corresponding p-value  $p = P(F_{(J,n-K)} > F)$ , then compare p-value with  $\alpha$ .





# Joint Hypothesis Testing

- State the conclusion: suppose  $F = 8.44$ ,  $J = 2$ ,  $n - K = 71$ ,  $\alpha = 0.05$  then  $F_c = 3.126$ , since

$$8.44 = F > F_c = F_{(1-\alpha, J, n-K)} = F_{(0.95, 2, 71)} = 3.126 \quad (11)$$

we reject the null hypothesis that both  $\beta_3 = 0$  and  $\beta_4 = 0$ , and conclude that at least one of them is not equal to zero.

- Then go back to the example, we conclude that advertising does have a significant effect on sales revenue.

# Joint Hypothesis Testing

In STATA output table, consider the general multiple regression model with  $K - 1$  explanatory variables and  $K$  unknown parameters,

$$y = \beta_1 + \beta_2 x_2 + \beta_3 x_3 + \cdots + \beta_K x_K + e \quad (12)$$

- **To examine whether we have a viable model**, STATA automatically does the following hypothesis testing:

$$H_0 : \beta_2 = \beta_3 = \cdots = \beta_K = 0 \quad (13)$$

$$H_1 : \text{at least one } \beta_k \neq 0, k = 2, 3, \dots, K \quad (14)$$

- This is referred to as a test of overall significance of the regression model;
- We use F test to test the above  $H_0$  against  $H_1$ .

# Joint Hypothesis Testing

- The unrestricted model is equation (12).
- Assuming  $H_0$  is true, the restricted model becomes:

$$y = \beta_1 + e \quad (15)$$

then the OLS estimator of  $\beta_1$  in the restricted model is:

$$b_1^* = \frac{1}{n} \sum_{i=1}^n y_i = \bar{y} \quad (16)$$

and

$$SSE_R = \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \sum_{i=1}^n (y_i - b_1^*)^2 = \sum_{i=1}^n (y_i - \bar{y})^2 = SST \quad (17)$$

- Thus to test the overall significance of a model (not in general, only for multiple regression model), the F-statistic can be modified and written as:

$$F = \frac{(SSE_R - SSE_U)/(K - 1)}{SSE_U/(n - K)} = \frac{(SST - SSE)/(K - 1)}{SSE/(n - K)} \quad (18)$$

# Joint Hypothesis Testing

$$FOODEXP = \beta_1 + \beta_2 INCOME + e \quad (19)$$

For joint hypothesis testing to test the overall significance of the model, what is  $J$ ?

```
. reg food_exp income
```

Source	SS	df	MS	Number of obs = 40		
Model	190626.984	1	190626.984	F( 1, 38) = 23.79		
Residual	304505.176	38	8013.2941	Prob > F = 0.0000		
Total	495132.16	39	12695.6964	R-squared = 0.3850		
				Adj R-squared = 0.3688		
				Root MSE = 89.517		

food_exp	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]	
income	10.20964	2.093264	4.88	0.000	5.972052	14.44723
_cons	83.416	43.41016	1.92	0.062	-4.463279	171.2953

$$SST = 495132.16, SSE = 304505.176$$

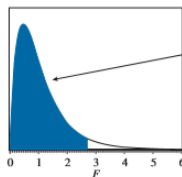
# Joint Hypothesis Testing

- F-statistic:

$$F = \frac{(SST - SSE)/(K - 1)}{SSE/(n - K)} = \frac{(495132.16 - 304505.176)/(2 - 1)}{304505.176/(40 - 2)} = 23.7888 \quad (20)$$

- **Option 2:** The STATA calculates p-value as very close to zero, given significance level  $\alpha = 0.05$ , then we reject  $H_0$ .
- **Option 1:** How can we check  $F_c = F_{(0.95, 1, 38)}$ ?

# Joint Hypothesis Testing



Example:  
 $P(F_{(4,30)} \leq 2.69) = 0.95$   
 $P(F_{(4,30)} > 2.69) = 0.05$

**Table 4** 95th Percentile for the  $F$ -distribution

$v_2/v_1$	1	2	3	4	5	6	7
1	161.45	199.50	215.71	224.58	230.16	233.99	236.77
2	18.51	19.00	19.16	19.25	19.30	19.33	19.35
3	10.13	9.55	9.28	9.12	9.01	8.94	8.89
4	7.71	6.94	6.59	6.39	6.26	6.16	6.09
5	6.61	5.79	5.41	5.19	5.05	4.95	4.88
6	5.99	5.14	4.76	4.53	4.39	4.28	4.21
7	5.59	4.74	4.35	4.12	3.97	3.87	3.79
8	5.32	4.46	4.07	3.84	3.69	3.58	3.50
9	5.12	4.26	3.86	3.63	3.48	3.37	3.29
10	4.96	4.10	3.71	3.48	3.33	3.22	3.14
15	4.54	3.68	3.29	3.06	2.90	2.79	2.71
20	4.35	3.49	3.10	2.87	2.71	2.60	2.51
25	4.24	3.39	2.99	2.76	2.60	2.49	2.40
30	4.17	3.32	2.92	2.69	2.53	2.42	2.33
35	4.12	3.27	2.87	2.64	2.49	2.37	2.29
40	4.08	3.23	2.84	2.61	2.45	2.34	2.25
--	---	---	---	---	---	---	---

# Joint Hypothesis Testing

**Connection between F test and t test:** when testing single hypothesis (focus on the significance of only one parameter), F test and t test are equivalent.

**Suppose now we want to test whether PRICE affects SALES**, that is for:

$$SALES = \beta_1 + \beta_2 PRICE + \beta_3 ADVERT + \beta_4 ADVERT^2 + e \quad (21)$$

we want to test

$$H_0 : \beta_2 = 0 \quad (22)$$

$$H_1 : \beta_2 \neq 0 \quad (23)$$

or restricted model is:

$$SALES = \beta_1 + \beta_3 ADVERT + \beta_4 ADVERT^2 + e \quad (24)$$

with F-statistic equal to 53.355.

# Joint Hypothesis Testing

**Fitted model:**

$$\widehat{SALES}_{(Se)} = 109.72 - \frac{7.64}{(1.046)} PRICE + \frac{12.15}{(5.556)} ADVERT - \frac{2.77}{(0.941)} ADVERT^2 \quad (25)$$

Use t test:

$$t = \frac{-7.64}{1.046} \quad (26)$$

and

$$t^2 = \left( \frac{-7.64}{1.046} \right)^2 = 53.35 = F \quad (27)$$

Then

$$t \text{ is either too large or too small} \iff F \text{ is too large} \quad (28)$$

So for single hypothesis testing, **two-tail** t test and F test are consistent. Actually in lecture 1, we have:  $F_{(1,m)} = t_m^2$ , for  $\forall m$ . In this case,  $m = n - 2$ .



# Joint Hypothesis Testing

**How about the single hypothesis test on linear combination of parameters?**

- Consider testing the following claim: the marginal sales of advertising when advertising expenditure is \$1900/month is equal to \$1 (assume unit of ADVERT is \$1000 and marginal cost of advertising is \$1), which means \$1900/month is the optimal advertising expenditure

⇒

$$H_0 : \beta_3 + 2\beta_4 ADVERT|_{ADVERT=1.9} = 1 \quad (29)$$

⇒

$$H_0 : \beta_3 + 3.8\beta_4 = 1 \quad (30)$$

$$H_1 : \beta_3 + 3.8\beta_4 \neq 1 \quad (31)$$

- We already learned how to use **t test** to test the above  $H_0$  against  $H_1$ .

# Joint Hypothesis Testing

- We can also use **F test** to test the above  $H_0$  (restricted model) against  $H_1$  (unrestricted model).
- **When  $H_0$  applies, the restricted model is:**

$$SALES = \beta_1 + \beta_2 PRICE + (1 - 3.8\beta_4) ADVERT + \beta_4 ADVERT^2 + e \quad (32)$$

$\implies$

$$SALES - ADVERT = \beta_1 + \beta_2 PRICE + \beta_4 (ADVERT^2 - 3.8 ADVERT) + e \quad (33)$$

- The F-statistic is:

$$F = \frac{(SSE_R - SSE_U)/1}{SSE_U/(n - K)} = 0.9362 \quad (34)$$

suppose  $\alpha = 0.05$ ,  $F_c = 3.976$ , then  $F < F_c$ , we cannot reject  $H_0$ .

- We conclude an advertising expenditure of \$1900/month is optimal.

# Joint Hypothesis Testing

- Suppose now we have the following hypothesis:

$$H_0 : \beta_3 + 3.8\beta_4 \leq 1 \quad (35)$$

$$H_1 : \beta_3 + 3.8\beta_4 > 1 \quad (36)$$

- In this case, we can no longer use F test.
- Because  $F_{(1,n-K)} = t_{n-K}^2$  cannot distinguish between the left and right tails as needed for a one-tail test.
- When we have alternative hypothesis  $H_1$  containing inequality signs  $\leq$ ,  $\geq$ , we restrict to t-test. (test using t statistic and t distribution)

# Model Specification

In any econometric investigation, choice of the model is one of the first steps

- What are the important considerations when choosing a model?
- What are the consequences of choosing the wrong model?
- Are there ways of assessing whether a model is adequate?

We have already learned some ways to evaluate a model: significance separately or jointly for parameters,  $R^2$  to measure the goodness-of-fit.

## Omitted Variable Bias

- It is possible that a chosen model may have important variables omitted, possibly because the economic theory has overlooked a variable or the lack of data makes us drop a variable even when it is prescribed by economic theory.

$$SALES = \beta_1 + \beta_2 PRICE + e \quad (37)$$

# Model Specification

Consider the following model:

$$FAMINC = \beta_1 + \beta_2 HEDUC + \beta_3 WEDUC + e \quad (38)$$

where FAMINC is family income, HEDUC is husband's education, WEDUC is wife's education.

- The estimated model is:

$$\widehat{FAMINC} = \underset{(Se)}{-5534} + \underset{(11230)}{3132} HEDUC + \underset{(1066)}{4523} WEDUC \quad (39)$$

- If we incorrectly omit WEDUC,

$$\widehat{FAMINC} = \underset{(Se)}{-26191} + \underset{(8541)}{5155} HEDUC \quad (40)$$

- Omitting WEDUC leads us to overstate the effect of an extra year of husband's education on family income by about \$2000**
- Then omission of a relevant variable leads to an estimator that is biased, which is called **omitted-variable bias**.

# Model Specification

More generally, write a general model as:

$$y = \beta_1 + \beta_2 x_2 + \beta_3 x_3 + e \quad (41)$$

- Omitting  $x_3$  is equivalent to imposing restriction  $\beta_3 = 0$ ;
- It can be viewed as an example of imposing an incorrect constraint on the parameters;
- Suppose  $b_2^*$  is estimator of  $\beta_2$  in the following model:

$$y = \beta_1 + \beta_2 x_2 + e \quad (42)$$

then we analyze the bias of  $b_2^*$ .

- The bias is:

$$\text{bias}(b_2^*) = E(b_2^*) - \beta_2 = \beta_3 \frac{\widehat{Cov}(x_2, x_3)}{\widehat{Var}(x_2)} \quad (43)$$

why  $E(b_2^*) \neq \beta_2$  in this case?

# Model Specification

Given

$$\text{bias}(b_2^*) = E(b_2^*) - \beta_2 = \beta_3 \frac{\widehat{Cov}(x_2, x_3)}{\widehat{Var}(x_2)} \quad (44)$$

- If wife's education has positive effect on family income:  $\beta_3 > 0$ ;
- If wife's education is positively correlated with husband's education:  $\widehat{Cov}(x_2, x_3) > 0$ ;
- Then we can conclude that the bias is positive, in words, the second estimated regression attributes too much to the husbands education because of the omission of the wife's education.



# Model Specification

**Table 6.1** Correlation Matrix for Variables Used in Family Income Example

	<i>FAMINC</i>	<i>HEDU</i>	<i>WEDU</i>	<i>KL6</i>	$X_5$	$X_6$
<i>FAMINC</i>	1.000					
<i>HEDU</i>	0.355	1.000				
<i>WEDU</i>	0.362	0.594	1.000			
<i>KL6</i>	-0.072	0.105	0.129	1.000		
$X_5$	0.290	0.836	0.518	0.149	1.000	
$X_6$	0.351	0.821	0.799	0.160	0.900	1.000

- $KL6$ : number of kids lower than six years old;
- $X_5$  and  $X_6$  are just another two economic variables possibly affecting family income and highly correlated with  $HEDU$  and  $WEDU$ .

# Model Specification

Now consider the model:

$$\widehat{FAMINC} = -7755 + 3212HEDUC + 4777WEDUC - 14311KL6 \quad (45)$$

(Se)                    (11163)                    (797)                    (1061)                    (5004)

- Note that in this example the coefficient estimators for HEDUC and WEDUC have not changed too much, because *KL6* is not highly correlated with those two education variables.

## The more explanatory variables, the better?

- The presence of many explanatory variables may inflate the variances of the estimators because of multi-collinearity, remember

$$\text{Var}(b_2) = \frac{\sigma^2}{(1-r_{23}^2) \sum_{i=1}^n (x_{2i} - \bar{x}_2)^2}.$$

- Consider the following fitted model:

$$\widehat{FAMINC} = \underset{(Se)}{-7755} + \underset{(11195)}{3340} HEDUC + \underset{(1250)}{5869} WEDUC - \underset{(5044)}{14200} KL6 \quad (46)$$
$$- \underset{(2242)}{889} X_5 + \underset{(1982)}{1067} X_6$$

The inclusion of irrelevant variables ( $X_5$  and  $X_6$ ) has reduced the precision of the coefficient estimators for other explanatory variables in the regression.

# Model Specification

## Some important points for choosing a model form:

- Choose explanatory variables and a functional form based on your theoretical and general understanding of the relationship;
- If a fitted model has estimators with unexpected signs, or unrealistic magnitudes, they could be caused by a mis-specification such as the omission of an important explanatory variable;
- One method for assessing whether one or a group of explanatory variables should be included in an equation is to perform significance tests, both separately and jointly.
- We already talked about how to modify the measure of goodness-of-fit to prevent adding as many explanatory variables as possible

$$\bar{R}^2 = 1 - \frac{SSE/(n - K)}{SST/(n - 1)} \quad (47)$$

**Table 6.2** Goodness-of-Fit and Information Criteria for Family Income Example

Included Variables	$R^2$	$\bar{R}^2$	AIC	SC
<i>HEDU</i>	0.1258	0.1237	21.262	21.281
<i>HEDU, WEDU</i>	0.1613	0.1574	21.225	21.253
<i>HEDU, WEDU, KL6</i>	0.1771	0.1714	21.211	21.248
<i>HEDU, WEDU, KL6, X5, X6</i>	0.1778	0.1681	21.219	21.276

- Selecting variables to maximize  $\bar{R}^2$  can be viewed as selecting variables to minimize SSE, **subject to a penalty for introducing too many variables.**
- Both the other two information criteria: the **AIC** and the **SC(BIC)** work in a similar way, but with different penalties for introducing too many variables. (Both are positively correlated with  $SSE$  and positively correlated with  $K$ )

## When would a model be mis-specified?

- We have omitted important explanatory variables;
- Included irrelevant ones;
- Chosen a wrong functional form;
- Have a model that violates the assumptions of the multiple regression model, most usually to violate the “no-exact-collinearity” assumption.
- Poor data quality, e.g. from uncontrolled experiment which will generate economic variables that move together in a systematic way

$$\text{Var}(b_2) = \frac{\sigma^2}{(1 - r_{23}^2) \sum_{i=1}^n (x_{2i} - \bar{x}_2)^2} \quad (48)$$

# Model Specification

## Example:

- MPG=miles per gallon;
- CYL=number of cylinders;
- ENG= engine displacement in cubic inches
- WGT=vehicle weight in pounds

Regression of MPG on CYL is:

$$\begin{array}{rcc} \widehat{MPG} & = & 42.9 - 3.558CYL \\ (Se) & & (0.83) \quad (0.146) \\ (p\text{-value}) & & (0.000) \quad (0.000) \end{array} \quad (49)$$

Now add ENG and WGT:

$$\begin{array}{rccccc} \widehat{MPG} & = & 44.4 & - & 0.268CYL & - & 0.0127ENG & - & 0.0057WGT \\ (Se) & & (1.5) & & (0.413) & & (0.0083) & & (0.00071) \\ (p\text{-value}) & & (0.000) & & (0.517) & & (0.125) & & (0.000) \end{array} \quad (50)$$

## How to test collinearity?

- One simple way to detect collinear relationships is to use sample correlation coefficients between pairs of explanatory variables;
- However, in some cases, collinear relationships involve more than two of the explanatory variables, the collinearity may not be detected by examining pairwise correlations;
- Try an auxiliary model:

$$x_2 = a_1x_1 + a_3x_3 + a_4x_4 + \cdots + a_Kx_K + v \quad (51)$$

- If the  $R^2$  (or adjusted  $\bar{R}^2$ ) from this artificial model is high, e.g. above 0.8, then the implication is that a large portion of the variation in  $x_2$  is explained by variation in the other explanatory variables.