

# **Econ 422 – Lecture Notes**

## **Part IV**

**(These notes are slightly modified versions of lecture notes provided by Stock and Watson, 2007. They are for instructional purposes only and are not to be distributed outside of the classroom.)**

# Hypothesis Tests and Confidence Intervals in Multiple Regression

## Outline

1. Hypothesis tests and confidence intervals for a single coefficient
2. Joint hypothesis tests on multiple coefficients
3. Other types of hypotheses involving multiple coefficients
4. How to decide what variables to include in a regression model?

# Hypothesis Tests and Confidence Intervals for a Single Coefficient in Multiple Regression

- $\frac{\hat{\beta}_1 - E(\hat{\beta}_1)}{\sqrt{\text{var}(\hat{\beta}_1)}}$  is approximately distributed  $N(0,1)$  (CLT).
- Thus hypotheses on  $\beta_1$  can be tested using the usual  $t$ -statistic, and confidence intervals are constructed as  $\{\hat{\beta}_1 \pm 1.96 \times SE(\hat{\beta}_1)\}$ .
- So too for  $\beta_2, \dots, \beta_k$ .
- $\hat{\beta}_1$  and  $\hat{\beta}_2$  are generally not independently distributed – so neither are their  $t$ -statistics (more on this later).

**Example:** The California class size data

$$(1) \quad \overline{TestScore} = 698.9 - 2.28 \times STR$$

(10.4) (0.52)

$$(2) \quad \overline{TestScore} = 686.0 - 1.10 \times STR - 0.650 PctEL$$

(8.7) (0.43) (0.031)

- The coefficient on  $STR$  in (2) is the effect on  $TestScores$  of a unit change in  $STR$ , holding constant the percentage of English Learners in the district
- The coefficient on  $STR$  falls by one-half
- The 95% confidence interval for coefficient on  $STR$  in (2) is  $\{-1.10 \pm 1.96 \times 0.43\} = (-1.95, -0.26)$
- The  $t$ -statistic testing  $\beta_{STR} = 0$  is  $t = -1.10/0.43 = -2.54$ , so we reject the hypothesis at the 5% significance level

# Standard errors in multiple regression in STATA

```
reg testscr str pctel, robust;
```

Regression with robust standard errors

```
Number of obs =      420
F(   2,   417) =    223.82
Prob > F       =     0.0000
R-squared      =     0.4264
Root MSE      =    14.464
```

testscr	Coef.	Robust Std. Err.	t	P> t	[95% Conf. Interval]	
str	-1.101296	.4328472	-2.54	0.011	-1.95213	-.2504616
pctel	-.6497768	.0310318	-20.94	0.000	-.710775	-.5887786
_cons	686.0322	8.728224	78.60	0.000	668.8754	703.189

$$\overline{TestScore} = 686.0 - 1.10 \times STR - 0.650 PctEL$$

(8.7)    (0.43)            (0.031)

We use **heteroskedasticity-robust standard errors** – for exactly the same reason as in the case of a single regressor.

# Tests of Joint Hypotheses

Let  $Expn$  = expenditures per pupil and consider the population regression model:

$$TestScore_i = \beta_0 + \beta_1 STR_i + \beta_2 Expn_i + \beta_3 PctEL_i + u_i$$

The null hypothesis that “school resources don’t matter,” and the alternative that they do, corresponds to:

$$H_0: \beta_1 = 0 \text{ and } \beta_2 = 0$$

vs.  $H_1$ : ***either***  $\beta_1 \neq 0$  ***or***  $\beta_2 \neq 0$  ***or both***

$$TestScore_i = \beta_0 + \beta_1 STR_i + \beta_2 Expn_i + \beta_3 PctEL_i + u_i$$

## *Tests of joint hypotheses, ctd.*

$$H_0: \beta_1 = 0 \text{ and } \beta_2 = 0$$

vs.  $H_1$ : *either*  $\beta_1 \neq 0$  *or*  $\beta_2 \neq 0$  *or both*

- A *joint hypothesis* specifies a value for two or more coefficients, that is, it imposes a restriction on two or more coefficients.
- In general, a joint hypothesis will involve  $q$  restrictions. In the example above,  $q = 2$ , and the two restrictions are  $\beta_1 = 0$  and  $\beta_2 = 0$ .
- A “common sense” idea is to reject if either of the individual  $t$ -statistics exceeds 1.96 in absolute value.
- But this “one at a time” test isn’t valid: the resulting test rejects too often under the null hypothesis (more than 5%)!

## *Why can't we just test the coefficients one at a time?*

Because the rejection rate under the null isn't 5%. We'll calculate the probability of incorrectly rejecting the null using the “common sense” test based on the two individual  $t$ -statistics. To simplify the calculation, suppose that  $\hat{\beta}_1$  and  $\hat{\beta}_2$  are independently distributed. Let  $t_1$  and  $t_2$  be the  $t$ -statistics:

$$t_1 = \frac{\hat{\beta}_1 - 0}{SE(\hat{\beta}_1)} \text{ and } t_2 = \frac{\hat{\beta}_2 - 0}{SE(\hat{\beta}_2)}$$

The “one at time” test is:

reject  $H_0: \beta_1 = \beta_2 = 0$  if  $|t_1| > 1.96$  and/or  $|t_2| > 1.96$

What is the probability that this “one at a time” test rejects  $H_0$ , when  $H_0$  is actually true? (We *want* it to be 5%.)



***Suppose  $t_1$  and  $t_2$  are independent (for this calculation).***

The probability of incorrectly rejecting the null hypothesis using the “one at a time” test

$$\begin{aligned} &= \Pr_{H_0} [|t_1| > 1.96 \text{ and/or } |t_2| > 1.96] \\ &= \Pr_{H_0} [|t_1| > 1.96, |t_2| > 1.96] + \Pr_{H_0} [|t_1| > 1.96, |t_2| \leq 1.96] \\ &\quad + \Pr_{H_0} [|t_1| \leq 1.96, |t_2| > 1.96] \quad (\text{disjoint events}) \\ &= \Pr_{H_0} [|t_1| > 1.96] \times \Pr_{H_0} [|t_2| > 1.96] \\ &\quad + \Pr_{H_0} [|t_1| > 1.96] \times \Pr_{H_0} [|t_2| \leq 1.96] \\ &\quad + \Pr_{H_0} [|t_1| \leq 1.96] \times \Pr_{H_0} [|t_2| > 1.96] \\ &\quad (t_1, t_2 \text{ are independent by assumption}) \\ &= .05 \times .05 + .05 \times .95 + .95 \times .05 \\ &= .0975 = 9.75\% - \text{which is } \textit{not} \text{ the desired } 5\%!! \end{aligned}$$

The *size* of a test is the actual rejection rate under the null hypothesis.

- The size of the “common sense” test isn’t 5%!
- In fact, its size depends on the correlation between  $t_1$  and  $t_2$  (and thus on the correlation between  $\hat{\beta}_1$  and  $\hat{\beta}_2$ ).

### **Two Solutions:**

- Use a different critical value in this procedure – not 1.96 (this is the “Bonferroni method) (this method is rarely used in practice however)
- Use a different test statistic that test both  $\beta_1$  and  $\beta_2$  at once: the  $F$ -statistic (this is common practice)

## The $F$ -statistic

The  $F$ -statistic tests all parts of a joint hypothesis at once.

Formula for the special case of the joint hypothesis  $\beta_1 = \beta_{1,0}$  and  $\beta_2 = \beta_{2,0}$  in a regression with two regressors:

$$F = \frac{1}{2} \left( \frac{t_1^2 + t_2^2 - 2\hat{\rho}_{t_1, t_2} t_1 t_2}{1 - \hat{\rho}_{t_1, t_2}^2} \right)$$

where  $\hat{\rho}_{t_1, t_2}$  estimates the correlation between  $t_1$  and  $t_2$ .

Reject when  $F$  is large (how large?)

The  $F$ -statistic testing  $\beta_1$  and  $\beta_2$ :

$$F = \frac{1}{2} \left( \frac{t_1^2 + t_2^2 - 2\hat{\rho}_{t_1, t_2} t_1 t_2}{1 - \hat{\rho}_{t_1, t_2}^2} \right)$$

- The  $F$ -statistic is large when  $t_1$  and/or  $t_2$  is large
- The  $F$ -statistic corrects (in just the right way) for the correlation between  $t_1$  and  $t_2$ .
- The formula for more than two  $\beta$ 's is nasty unless you use matrix algebra.
- This gives the  $F$ -statistic a nice large-sample approximate distribution, which is...

## Large-sample distribution of the $F$ -statistic

Consider *special case* that  $t_1$  and  $t_2$  are independent, so  $\hat{\rho}_{t_1, t_2}$

$\xrightarrow{p} 0$ ; in large samples the formula becomes

$$F = \frac{1}{2} \left( \frac{t_1^2 + t_2^2 - 2\hat{\rho}_{t_1, t_2} t_1 t_2}{1 - \hat{\rho}_{t_1, t_2}^2} \right) \approx \frac{1}{2} (t_1^2 + t_2^2)$$

- Under the null,  $t_1$  and  $t_2$  have standard normal distributions that, in this special case, are independent
- The large-sample distribution of the  $F$ -statistic is the distribution of the average of two independently distributed squared standard normal random variables.

The *chi-squared* distribution with  $q$  degrees of freedom ( $\chi_q^2$ ) is defined to be the distribution of the sum of  $q$  independent squared standard normal random variables.

**In large samples,  $F$  is distributed as  $\chi_q^2/q$ .**

**Selected large-sample critical values of  $\chi_q^2/q$**

<u><math>q</math></u>	<u>5% critical value</u>	
1	3.84	( <i>why?</i> )
2	3.00	(the case $q=2$ above)
3	2.60	
4	2.37	
5	2.21	

### ***Computing the $p$ -value using the $F$ -statistic:***

$p$ -value = tail probability of the  $\chi^2_q/q$  distribution  
beyond the  $F$ -statistic actually computed.

### **Implementation in STATA**

Use the “test” command after the regression

*Example:* Test the joint hypothesis that the population coefficients on  $STR$  and expenditures per pupil ( $expn\_stu$ ) are both zero, against the alternative that at least one of the population coefficients is nonzero.

# F-test example, California class size data:

```
reg testscr str expn_stu pctel, r;
```

Regression with robust standard errors

Number of obs = 420  
 F( 3, 416) = 147.20  
 Prob > F = 0.0000  
 R-squared = 0.4366  
 Root MSE = 14.353

testscr	Coef.	Robust Std. Err.	t	P> t	[95% Conf. Interval]	
str	-.2863992	.4820728	-0.59	0.553	-1.234001	.661203
expn_stu	.0038679	.0015807	2.45	0.015	.0007607	.0069751
pctel	-.6560227	.0317844	-20.64	0.000	-.7185008	-.5935446
_cons	649.5779	15.45834	42.02	0.000	619.1917	679.9641

## NOTE

```
test str expn_stu;
```

*The test command follows the regression*

```
( 1) str = 0.0
```

*There are q=2 restrictions being tested*

```
( 2) expn_stu = 0.0
```

F( 2, 416) = 5.43  
 Prob > F = 0.0047

*The 5% critical value for q=2 is 3.00  
 Stata computes the p-value for you*



**More on  $F$ -statistics:** *a simple  $F$ -statistic formula that is easy to understand (it is only valid if the errors are homoskedastic, but it might help intuition).*

### **The homoskedasticity-only $F$ -statistic**

When the errors are homoskedastic, there is a simple formula for computing the “homoskedasticity-only”  $F$ -statistic:

- Run two regressions, one under the null hypothesis (the “restricted” regression) and one under the alternative hypothesis (the “unrestricted” regression).
- Compare the fits of the regressions – the  $R^2$ 's – if the “unrestricted” model fits sufficiently better, reject the null

## *The “restricted” and “unrestricted” regressions*

*Example:* are the coefficients on *STR* and *Expn* zero?

Unrestricted population regression (under  $H_1$ ):

$$TestScore_i = \beta_0 + \beta_1 STR_i + \beta_2 Expn_i + \beta_3 PctEL_i + u_i$$

Restricted population regression (that is, under  $H_0$ ):

$$TestScore_i = \beta_0 + \beta_3 PctEL_i + u_i \quad (why?)$$

- The number of restrictions under  $H_0$  is  $q = 2$  (why?).
- The fit will be better ( $R^2$  will be higher) in the unrestricted regression (why?)

By how much must the  $R^2$  increase for the coefficients on *Expn* and *PctEL* to be judged statistically significant?

## *Simple formula for the homoskedasticity-only F-statistic:*

$$F = \frac{(R_{unrestricted}^2 - R_{restricted}^2) / q}{(1 - R_{unrestricted}^2) / (n - k_{unrestricted} - 1)}$$

where:

$R_{restricted}^2$  = the  $R^2$  for the restricted regression

$R_{unrestricted}^2$  = the  $R^2$  for the unrestricted regression

$q$  = the number of restrictions under the null

$k_{unrestricted}$  = the number of regressors in the  
unrestricted regression.

- The bigger the difference between the restricted and unrestricted  $R^2$ 's – the greater the improvement in fit by adding the variables in question – the larger is the homoskedasticity-only  $F$ .

**Example:**

Restricted regression:

$$\overline{TestScore} = 644.7 - 0.671PctEL, \quad R^2_{restricted} = 0.4149$$
$$(1.0) \quad (0.032)$$

Unrestricted regression:

$$\overline{TestScore} = 649.6 - 0.29STR + 3.87Expn - 0.656PctEL$$
$$(15.5) \quad (0.48) \quad (1.59) \quad (0.032)$$

$$R^2_{unrestricted} = 0.4366, \quad k_{unrestricted} = 3, \quad q = 2$$

so

$$F = \frac{(R^2_{unrestricted} - R^2_{restricted}) / q}{(1 - R^2_{unrestricted}) / (n - k_{unrestricted} - 1)}$$
$$= \frac{(.4366 - .4149) / 2}{(1 - .4366) / (420 - 3 - 1)} = \mathbf{8.01}$$

**Note:** Heteroskedasticity-robust  $F = \mathbf{5.43...}$

## *The homoskedasticity-only $F$ -statistic – summary*

$$F = \frac{(R_{unrestricted}^2 - R_{restricted}^2) / q}{(1 - R_{unrestricted}^2) / (n - k_{unrestricted} - 1)}$$

- The homoskedasticity-only  $F$ -statistic rejects when adding the two variables increased the  $R^2$  by “enough” – that is, when adding the two variables improves the fit of the regression by “enough”
- If the errors are homoskedastic, then the homoskedasticity-only  $F$ -statistic has a large-sample distribution that is  $\chi_q^2/q$ .
- But if the errors are heteroskedastic, the large-sample distribution is a mess and is not  $\chi_q^2/q$

## Digression: The $F$ distribution

Your regression printouts might refer to the “ $F$ ” distribution.

If the four multiple regression LS assumptions hold *and*:

5.  $u_i$  is homoskedastic, that is,  $\text{var}(u|X_1, \dots, X_k)$  does not depend on  $X$ 's
6.  $u_1, \dots, u_n$  are normally distributed

then the homoskedasticity-only  $F$ -statistic has the “ $F_{q, n-k-1}$ ” distribution, where  $q$  = the number of restrictions and  $k$  = the number of regressors under the alternative (the unrestricted model).

- The  $F$  distribution is to the  $\chi^2_q/q$  distribution what the  $t_{n-1}$  distribution is to the  $N(0,1)$  distribution

## ***The $F_{q,n-k-1}$ distribution:***

- The  $F$  distribution is tabulated many places
- As  $n \rightarrow \infty$ , the  $F_{q,n-k-1}$  distribution asymptotes to the  $\chi_q^2/q$  distribution:

**The  $F_{q,\infty}$  and  $\chi_q^2/q$  distributions are the same.**

- For  $q$  not too big and  $n \geq 100$ , the  $F_{q,n-k-1}$  distribution and the  $\chi_q^2/q$  distribution are essentially identical.
- Many regression packages (including STATA) compute  $p$ -values of  $F$ -statistics using the  $F$  distribution
- You will encounter the  $F$  distribution in published empirical work.

## Another digression: A little history of statistics...

- The theory of the homoskedasticity-only  $F$ -statistic and the  $F_{q,n-k-1}$  distributions rests on implausibly strong assumptions (are earnings normally distributed?)
- These statistics dates to the early 20<sup>th</sup> century... back in the days when data sets were small and computers were people...
- The  $F$ -statistic and  $F_{q,n-k-1}$  distribution were major breakthroughs: an easily computed formula; a single set of tables that could be published once, then applied in many settings; and a precise, mathematically elegant justification.



## *A little history of statistics, ctd...*

- The strong assumptions seemed a minor price for this breakthrough.
- But with modern computers and large samples we can use the heteroskedasticity-robust  $F$ -statistic and the  $F_{q,i}$  distribution, which only require the four least squares assumptions (not assumptions #5 and #6)
- This historical legacy persists in modern software, in which homoskedasticity-only standard errors (and  $F$ -statistics) are the default, and in which  $p$ -values are computed using the  $F_{q,n-k-1}$  distribution.

## Summary: the homoskedasticity-only $F$ -statistic and the $F$ distribution

- These are justified only under very strong conditions – stronger than are realistic in practice.
- Yet, they are widely used.
- *You* should use the heteroskedasticity-robust  $F$ -statistic, with  $\chi_q^2/q$  (that is,  $F_{q,\infty}$ ) critical values.
- For  $n \geq 100$ , the  $F$ -distribution essentially is the  $\chi_q^2/q$  distribution.
- For small  $n$ , sometimes researchers use the  $F$  distribution because it has larger critical values and in this sense is more conservative.

## Summary: testing joint hypotheses

- The “one at a time” approach of rejecting if either of the  $t$ -statistics exceeds 1.96 rejects more than 5% of the time under the null (the size exceeds the desired significance level)
- The heteroskedasticity-robust  $F$ -statistic is built in to STATA (“test” command); this tests all  $q$  restrictions at once.
- For  $n$  large, the  $F$ -statistic is distributed  $\chi_q^2/q (= F_{q,\infty})$
- The homoskedasticity-only  $F$ -statistic is important historically (and thus in practice), and can help intuition, but isn’t valid when there is heteroskedasticity

# Testing Single Restrictions on Multiple Coefficients

$$Y_i = \beta_0 + \beta_1 X_{1i} + \beta_2 X_{2i} + u_i, \quad i = 1, \dots, n$$

Consider the null and alternative hypothesis,

$$H_0: \beta_1 = \beta_2 \quad \text{vs.} \quad H_1: \beta_1 \neq \beta_2$$

This null imposes a *single* restriction ( $q = 1$ ) on *multiple* coefficients – it is not a joint hypothesis with multiple restrictions (compare with  $\beta_1 = 0$  and  $\beta_2 = 0$ ).

## *Testing single restrictions on multiple coefficients, ctd.*

Here are two methods for testing single restrictions on multiple coefficients:

1. *Rearrange (“transform”) the regression*

Rearrange the regressors so that the restriction becomes a restriction on a single coefficient in an equivalent regression; or,

2. *Perform the test directly*

Some software, including STATA, lets you test restrictions using multiple coefficients directly

## ***Method 1: Rearrange (“transform”) the regression***

$$Y_i = \beta_0 + \beta_1 X_{1i} + \beta_2 X_{2i} + u_i$$

$$H_0: \beta_1 = \beta_2 \quad \text{vs.} \quad H_1: \beta_1 \neq \beta_2$$

Add and subtract  $\beta_2 X_{1i}$ :

$$Y_i = \beta_0 + (\beta_1 - \beta_2) X_{1i} + \beta_2 (X_{1i} + X_{2i}) + u_i$$

or

$$Y_i = \beta_0 + \gamma_1 X_{1i} + \beta_2 W_i + u_i$$

where

$$\gamma_1 = \beta_1 - \beta_2$$

$$W_i = X_{1i} + X_{2i}$$

## ***Rearrange the regression, ctd.***

*(a) Original system:*

$$Y_i = \beta_0 + \beta_1 X_{1i} + \beta_2 X_{2i} + u_i$$

$$H_0: \beta_1 = \beta_2 \quad \text{vs.} \quad H_1: \beta_1 \neq \beta_2$$

*(b) Rearranged (“transformed”) system:*

$$Y_i = \beta_0 + \gamma_1 X_{1i} + \beta_2 W_i + u_i$$

where  $\gamma_1 = \beta_1 - \beta_2$  and  $W_i = X_{1i} + X_{2i}$

so

$$H_0: \gamma_1 = 0 \quad \text{vs.} \quad H_1: \gamma_1 \neq 0$$

The testing problem is now a simple one:

test whether  $\gamma_1 = 0$  in specification (b).

## *Method 2: Perform the test directly*

$$Y_i = \beta_0 + \beta_1 X_{1i} + \beta_2 X_{2i} + u_i$$

$$H_0: \beta_1 = \beta_2 \quad \text{vs.} \quad H_1: \beta_1 \neq \beta_2$$

*Example:*

$$\text{TestScore}_i = \beta_0 + \beta_1 \text{STR}_i + \beta_2 \text{Expn}_i + \beta_3 \text{PctEL}_i + u_i$$

In STATA, to test  $\beta_1 = \beta_2$  vs.  $\beta_1 \neq \beta_2$  (two-sided):

```
regress testscore str expn pctel, r  
test str=expn
```

The details of implementing this method are software-specific.



# Confidence Sets for Multiple Coefficients

$$Y_i = \beta_0 + \beta_1 X_{1i} + \beta_2 X_{2i} + \dots + \beta_k X_{ki} + u_i, \quad i = 1, \dots, n$$

What is a *joint* confidence set for  $\beta_1$  and  $\beta_2$ ?

A 95% *joint confidence set* is:

- A set-valued function of the data that contains the true parameter(s) in 95% of hypothetical repeated samples.
- The set of parameter values that cannot be rejected at the 5% significance level.
- You can find a 95% confidence set as the set of  $(\beta_1, \beta_2)$  that cannot be rejected at the 5% level using an  $F$ -test (*why not just combine the two 95% confidence intervals?*).

## *Joint confidence sets ctd.*

Let  $F(\beta_{1,0}, \beta_{2,0})$  be the (heteroskedasticity-robust)  $F$ -statistic testing the hypothesis that  $\beta_1 = \beta_{1,0}$  and  $\beta_2 = \beta_{2,0}$ :

95% confidence set =  $\{\beta_{1,0}, \beta_{2,0}: F(\beta_{1,0}, \beta_{2,0}) < 3.00\}$

- 3.00 is the 5% critical value of the  $F_{2,\infty}$  distribution
- This set has coverage rate 95% because the test on which it is based (the test it “inverts”) has size of 5%

*5% of the time, the test incorrectly rejects the null when the null is true, so 95% of the time it does not; therefore the confidence set constructed as the nonrejected values contains the true value 95% of the time (in 95% of all samples).*

*The confidence set based on the F-statistic is an ellipse*

$$\{\beta_1, \beta_2: F = \frac{1}{2} \left( \frac{t_1^2 + t_2^2 - 2\hat{\rho}_{t_1, t_2} t_1 t_2}{1 - \hat{\rho}_{t_1, t_2}^2} \right) \leq 3.00\}$$

Now

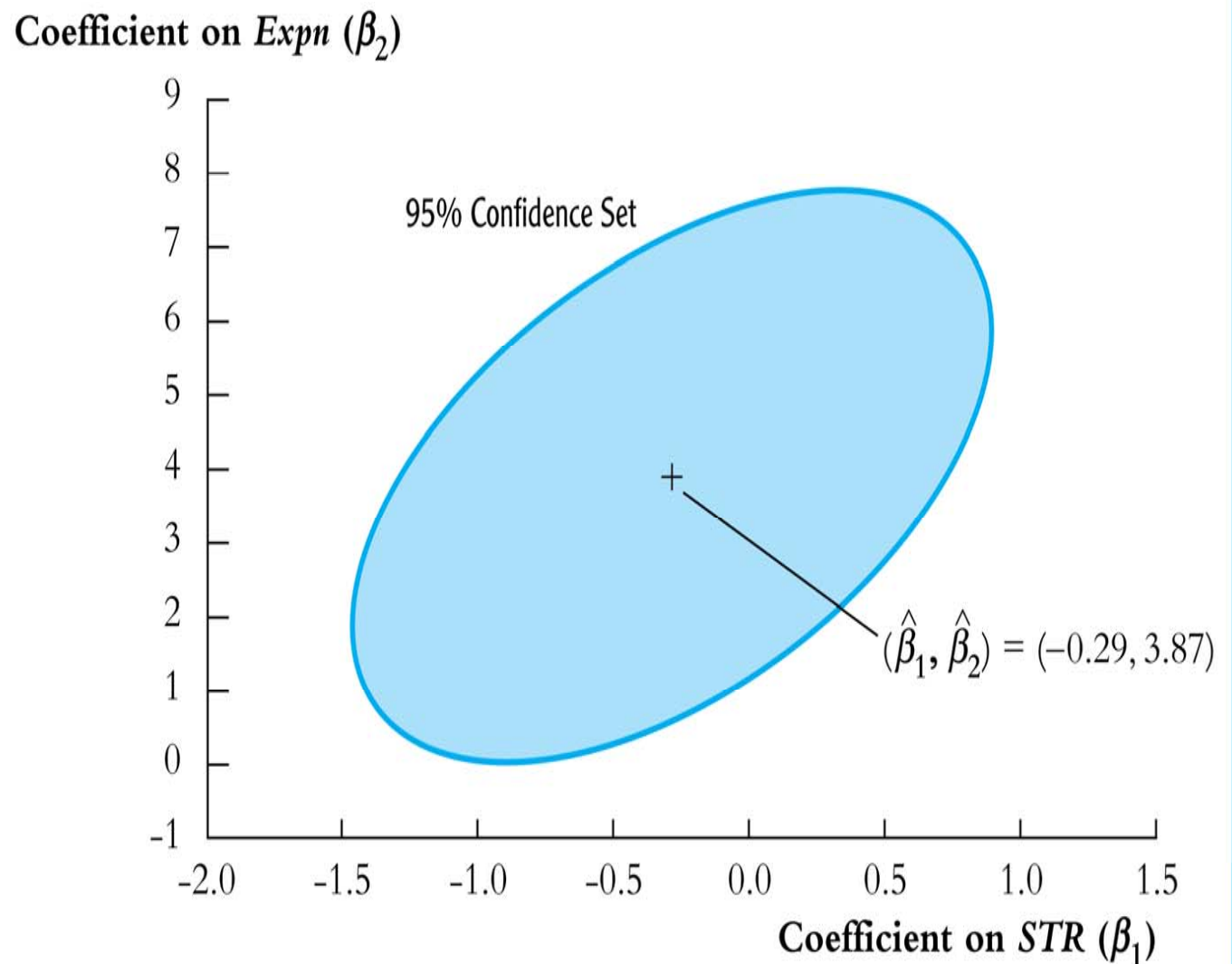
$$\begin{aligned} F &= \frac{1}{2(1 - \hat{\rho}_{t_1, t_2}^2)} \times [t_1^2 + t_2^2 - 2\hat{\rho}_{t_1, t_2} t_1 t_2] \\ &= \frac{1}{2(1 - \hat{\rho}_{t_1, t_2}^2)} \times \\ &\quad \left[ \left( \frac{\hat{\beta}_2 - \beta_{2,0}}{SE(\hat{\beta}_2)} \right)^2 + \left( \frac{\hat{\beta}_1 - \beta_{1,0}}{SE(\hat{\beta}_1)} \right)^2 + 2\hat{\rho}_{t_1, t_2} \left( \frac{\hat{\beta}_1 - \beta_{1,0}}{SE(\hat{\beta}_1)} \right) \left( \frac{\hat{\beta}_2 - \beta_{2,0}}{SE(\hat{\beta}_2)} \right) \right] \end{aligned}$$

This is a quadratic form in  $\beta_{1,0}$  and  $\beta_{2,0}$  – thus the boundary of the set  $F = 3.00$  is an ellipse.

## Confidence set based on inverting the $F$ -statistic

**FIGURE 7.1** 95% Confidence Set for Coefficients on *STR* and *Expn* from Equation (7.6)

The 95% confidence set for the coefficients on *STR* ( $\beta_1$ ) and *Expn* ( $\beta_2$ ) is an ellipse. The ellipse contains the pairs of values of  $\beta_1$  and  $\beta_2$  that cannot be rejected using the  $F$ -statistic at the 5% significance level.



*An example of a multiple regression analysis – and how to decide which variables to include in a regression...*

## **A Closer Look at the Test Score Data**

We want to get an unbiased estimate of the effect on test scores of changing class size, holding constant student and school characteristics (but not necessarily holding constant the budget (*why?*)).

To do this we need to think about what variables to include and what regressions to run – and we should do this before we actually sit down at the computer. This entails thinking beforehand about your *model specification*.

## A general approach to variable selection and “*model specification*”

- Specify a “base” or “benchmark” model.
- Specify a range of plausible alternative models, which include additional candidate variables.
- Does a candidate variable change the coefficient of interest ( $\beta_1$ )?
- Is a candidate variable statistically significant?
- Use judgment, not a mechanical recipe...
- Don’t just try to maximize  $R^2$ !

## *Digression about measures of fit...*

It is easy to fall into the trap of maximizing the  $R^2$  and  $\bar{R}^2$  – but this loses sight of our real objective, an unbiased estimator of the class size effect.

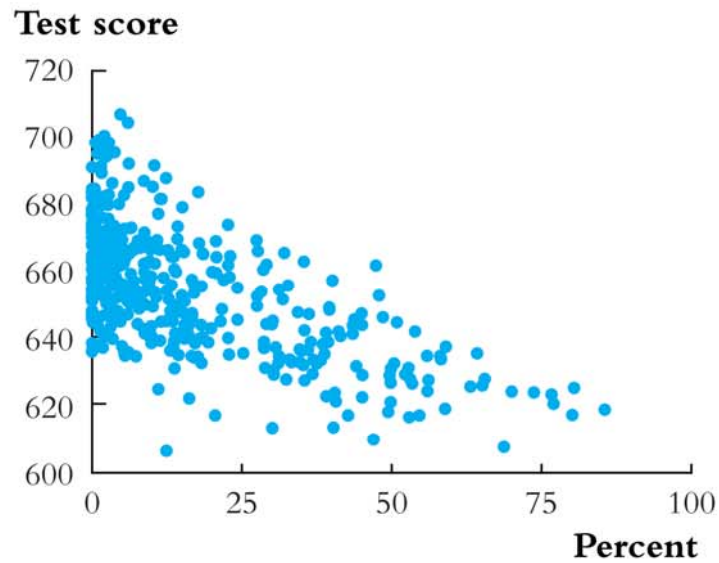
- A high  $R^2$  (or  $\bar{R}^2$ ) means that the regressors explain the variation in  $Y$ .
- A high  $R^2$  (or  $\bar{R}^2$ ) does *not* mean that you have eliminated omitted variable bias.
- A high  $R^2$  (or  $\bar{R}^2$ ) does *not* mean that you have an unbiased estimator of a causal effect ( $\beta_1$ ).
- A high  $R^2$  (or  $\bar{R}^2$ ) does *not* mean that the included variables are statistically significant – this must be determined using hypotheses tests.

## ***Back to the test score application:***

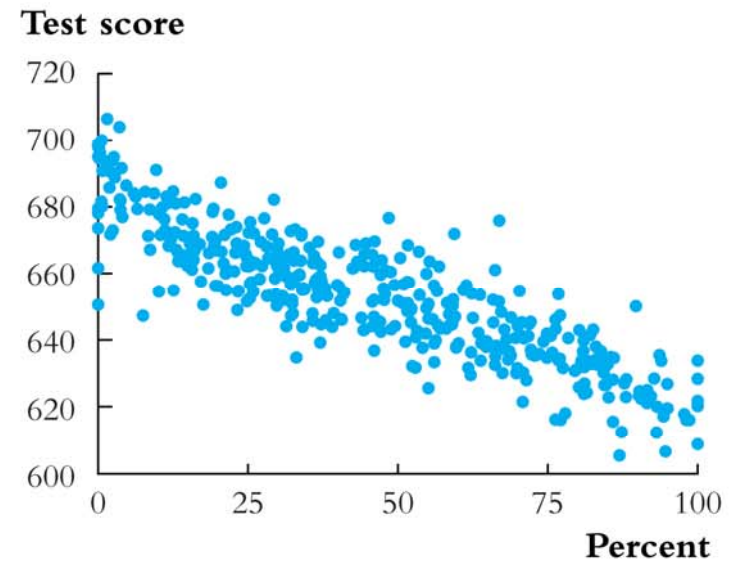
- *What variables would you want – ideally – to estimate the effect on test scores of STR using school district data?*
- *Variables actually in the California class size data set:*
  - student-teacher ratio (*STR*)
  - percent English learners in the district (*PctEL*)
  - school expenditures per pupil
  - name of the district (so we could look up average rainfall, for example)
  - percent eligible for subsidized/free lunch
  - percent on public income assistance
  - average district income
- *Which of these variables would you want to include?*



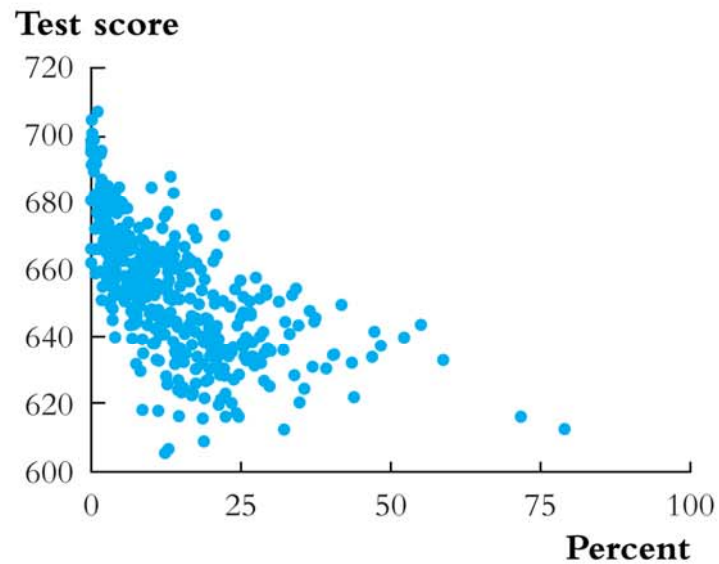
## *More California data...*



**(a)** Percentage of English language learners



**(b)** Percentage qualifying for reduced price lunch



**(c)** Percentage qualifying for income assistance

## Digression on presentation of regression results

- We have a number of regressions and we want to report them. It is awkward and difficult to read regressions written out in equation form, so instead it is conventional to report them in a table.
- A table of regression results should include:
  - estimated regression coefficients
  - standard errors
  - measures of fit
  - number of observations
  - relevant  $F$ -statistics, if any
  - Any other pertinent information.

Find this information in the following table:

**TABLE 7.1** Results of Regressions of Test Scores on the Student–Teacher Ratio and Student Characteristic Control Variables Using California Elementary School Districts

Dependent variable: average test score in the district.

Regressor	(1)	(2)	(3)	(4)	(5)
Student–teacher ratio ( $X_1$ )	−2.28** (0.52)	−1.10* (0.43)	−1.00** (0.27)	−1.31** (0.34)	−1.01** (0.27)
Percent English learners ( $X_2$ )		−0.650** (0.031)	−0.122** (0.033)	−0.488** (0.030)	−0.130** (0.036)
Percent eligible for subsidized lunch ( $X_3$ )			−0.547** (0.024)		−0.529** (0.038)
Percent on public income assistance ( $X_4$ )				−0.790** (0.068)	0.048 (0.059)
Intercept	698.9** (10.4)	686.0** (8.7)	700.2** (5.6)	698.0** (6.9)	700.4** (5.5)
<b>Summary Statistics</b>					
<i>SER</i>	18.58	14.46	9.08	11.65	9.08
$\overline{R}^2$	0.049	0.424	0.773	0.626	0.773
<i>n</i>	420	420	420	420	420

These regressions were estimated using the data on K-8 school districts in California, described in Appendix 4.1. Standard errors are given in parentheses under coefficients. The individual coefficient is statistically significant at the \*5% level or \*\*1% significance level using a two-sided test.

## Summary: Multiple Regression

- Multiple regression allows you to estimate the effect on  $Y$  of a change in  $X_1$ , holding  $X_2$  constant.
- If you can measure a variable, you can avoid omitted variable bias from that variable by including it.
- There is no simple recipe for deciding which variables belong in a regression – you must exercise judgment.
- One approach is to specify a base model – relying on *a-priori* reasoning – then explore the sensitivity of the key estimate(s) in alternative specifications.