

Multiple Regression Analysis: Inference. Wooldridge (2013), Chapter 4 and Chapter 6 (section 6.4)

- The t Test
- Hypothesis testing - one-sided alternatives
- Hypothesis testing - two-sided alternatives
- Confidence Intervals
- Testing a Linear Combination
- Multiple Linear Restrictions
- Testing Exclusion Restrictions
- The F statistic
- Overall Significance
- Prediction for the conditional mean of y
- Prediction for y
- Predicting y in a log model

Multiple Regression Analysis: Inference

The t Test

Under the CLM assumptions

$$\frac{\hat{\beta}_j - \beta_j}{se(\hat{\beta}_j)} \sim t(n - k - 1),$$

where

$$se(\hat{\beta}_j) = \frac{\hat{\sigma}}{\sqrt{SST_j(1 - R_j^2)}}.$$

- Note this is a *t* – *student* distribution because we estimate $sd(\hat{\beta}_j)$ by the standard error of $\hat{\beta}_j$, $se(\hat{\beta}_j)$,
- Note the degrees of freedom: $n - k - 1$ (sample size-number of parameters of the model).
- In the simple regression model $k = 1$.

Multiple Regression Analysis: Inference

The t Test (cont)

- Knowing the sampling distribution for the standardized estimator allows us to carry out hypothesis tests.
- Start with a null hypothesis $H_0 : \beta_j = b_j$, where b_j is a particular value.
- For example, $H_0 : \beta_j = 0$. If do not reject null, then x_j has no effect on the conditional mean of y , controlling for other x 's.

Multiple Regression Analysis: Inference

The t Test (cont)

- To perform our test we first need to form the statistic : $t_j = \frac{\hat{\beta}_j - b_j}{se(\hat{\beta}_j)}$.
- Besides our null, H_0 , we need an alternative hypothesis, H_1 , and a significance level α .

Alternatives:

- $H_1 : \beta_j > b_j$ and $H_1 : \beta_j < b_j$ are one-sided.
- $H_1 : \beta_j \neq b_j$ is a two-sided alternative.

Multiple Regression Analysis: Inference

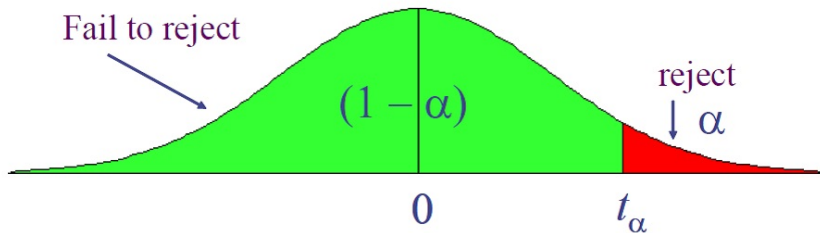
One-Sided Alternatives (cont)

$$y_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_k x_{ik} + u_i$$

- $H_0 : \beta_j = b_j$ vs $H_1 : \beta_j > b_j$.

Critical Value: t_α is defined as the constant that satisfies $\mathcal{P}(t_j > t_\alpha) = \alpha$, where t_j has the $t(n - k - 1)$ distribution. Equivalently $\mathcal{P}(t_j < t_\alpha) = 1 - \alpha$.

Rejection rule: Reject H_0 if the value of the t-statistic $> t_\alpha$.



Multiple Regression Analysis: Inference

One-Sided Alternatives (cont)

- **Example:** Consider the following regression where the standard errors are in brackets:

$$\widehat{\log(wages)} = \underset{(0.104)}{0.284} + \underset{(0.007)}{0.092educ} + \underset{(0.0017)}{0.0041exper} + \underset{(0.003)}{0.022tenure},$$
$$n = 526, R^2 = 0.316$$

Test whether, after controlling for education and tenure, higher work experience leads to higher hourly wages. Use the 5% and the 1% significance levels.

- $H_0 : \beta_2 = 0$ vs $H_1 : \beta_2 > 0$.
- $t_2^{act} = \frac{0.0041}{0.0017} = 2.41176$
- $df = n - k - 1 = 526 - 4 = 522$.
- $t_{0.05} = 1.645$.
- $t_{0.01} = 2.326$.
- Since $2.41176 > 1.645$, we reject H_0 in favour of H_1 at 5% level.
- Since $2.41176 > 2.326$, we reject H_0 in favour of H_1 at 1% level.
- Hence there is statistical evidence at both 5% and 1% that higher work experience leads to higher hourly wages.

Multiple Regression Analysis: Inference

One-Sided Alternatives (cont)

$$y_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_k x_{ik} + u_i$$

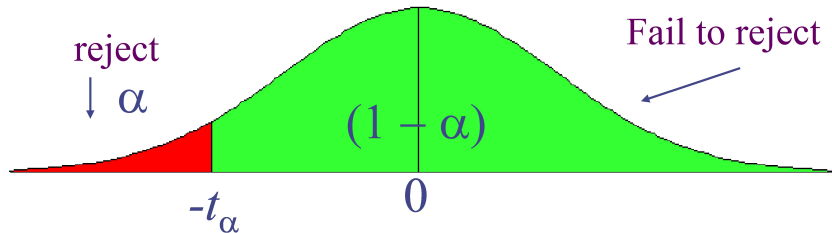
- $H_0 : \beta_j = b_j$ vs $H_1 : \beta_j < b_j$.

Critical Value: $-t_\alpha$ that is the constant that satisfies

$\mathcal{P}(t_j < -t_\alpha) = \alpha$, where t_j has the $t(n - k - 1)$ distribution.

Equivalently $\mathcal{P}(t_j > -t_\alpha) = 1 - \alpha$.

Rejection rule: Reject H_0 if the value of the t-statistic $< -t_\alpha$.



Multiple Regression Analysis: Inference

One-Sided Alternatives (cont)

Example: Student performance and school size

- Consider the following regression

$$\widehat{math10} = 2.274 + 0.00046totcomp + 0.048staff - 0.0002enroll,$$

(6.113) (0.0001) (0.04) (0.00022)

$$n = 408, R^2 = 0.0541$$

where

<i>math10</i>	-percentage of students passing math test
<i>totcomp</i>	-average annual teacher compensation
<i>staff</i>	-staff per one thousand students
<i>enroll</i>	-School enrollment=school size

Test whether smaller school size leads to better student performance at 5% level and 10% level.

Multiple Regression Analysis: Inference

One-Sided Alternatives (cont)

- $H_0 : \beta_3 = 0$ vs $H_1 : \beta_3 < 0$.
- $t_3^{act} = \frac{-0.0002}{0.00022} = -0.90909$
- $df = n - k - 1 = 408 - 4 = 404$
- $-t_{0.05} = -1.645$.
- $-t_{0.1} = -1.282$.
- Given that $-1.645 < -0.90909$ we fail to reject H_0 in favour of H_1 at 5% level.
- Given that $-1.282 < -0.90909$ we fail to reject H_0 in favour of H_1 at 10% level.
- Therefore, there is no statistical evidence (at 5% and 10% levels) that smaller school size leads to better student performance.

Multiple Regression Analysis: Inference

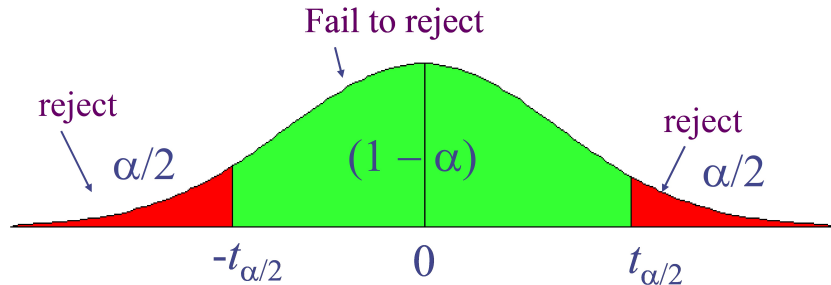
Two-Sided Alternatives

$$y_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_k x_{ik} + u_i$$

- $H_0 : \beta_j = b_j$ vs $H_1 : \beta_j \neq b_j$.

Critical Value: $t_{\alpha/2}$ is defined as the constant that satisfies $\mathcal{P}(t_j > t_{\alpha/2}) = \alpha/2$, where t_j has the $t(n - k - 1)$ distribution.

Rejection rule: Reject H_0 if the *absolute value* of the t-statistic $> t_{\alpha/2}$.



Multiple Regression Analysis: Inference

Two-Sided Alternatives

Example: Campus crime and enrollment

An interesting hypothesis is whether crime increases by one percent if enrollment is increased by one percent

$$\widehat{\log(\text{crime})} = \frac{-6.63}{(1.03)} + \frac{1.27}{(0.11)} \log(\text{enroll}),$$
$$n = 97, R^2 = 0.0541$$

The estimate 1.27 is different from one but is this difference statistically significant? (use the 5% significance level)?

- $H_0 : \beta_1 = 1$ vs $H_1 : \beta_1 \neq 1$.
- $t_1^{act} = \frac{1.27-1}{0.11} = 2.4545$
- $df = n - k - 1 = 97 - 2 = 95$
- $t_{0.025} = 1.985$
- In Wooldridge(2013)' book you can only find the critical values for the t-student distribution with $df = 90$. In this case $t_{0.025} = 1.987$.
- Given that $|2.4545| = 2.4545 > 1.985$ (also $|2.4545| > 1.987$) we reject H_0 in favour of H_1 at 5% level

Multiple Regression Analysis: Inference

Two-Sided Alternatives

Remarks on $H_0 : \beta_j = 0$ vs $H_1 : \beta_j \neq 0$

- The quantity $t_j = \frac{\hat{\beta}_j}{se(\hat{\beta}_j)}$ is called the t-ratio.
- If we reject the null, we typically say “ x_j is statistically significant at the α level”.
- If we fail to reject the null, we typically say “ x_j is statistically insignificant at the α level”.
- If asked to test whether a regressor is statistical significant, the alternative is assumed to be two-sided.

Multiple Regression Analysis: Inference

Two-Sided Alternatives

Example: Consider the following regression where the standard errors are in brackets:

$$\log(\widehat{wages}) = \underset{(0.104)}{0.284} + \underset{(0.007)}{0.092}educ + \underset{(0.0017)}{0.0041}exper + \underset{(0.003)}{0.022}tenure,$$
$$n = 526, R^2 = 0.316$$

Test whether, after controlling for experience and tenure, education is statistically significant at 5% and the 1% significance levels.

- $H_0 : \beta_1 = 0$ vs $H_1 : \beta_1 \neq 0$.
- $t_1^{act} = \frac{0.092}{0.007} = 13.1429$
- $df = 526 - 4 = 522$
- $t_{0.025} = 1.96$.
- $t_{0.005} = 2.576$.
- Since $|13.1429| = 13.1429 > 1.96$, we reject H_0 in favour of H_1 at 5% level.
- Since $|13.1429| = 13.1429 > 2.576$, we reject H_0 in favour of H_1 at 1% level.
- Therefore Education is statistically significant at 5% and 1% levels.

Multiple Regression Analysis: Inference

p-values

- The smallest significance level at which the null hypothesis is still rejected, is called the *p-value* of the hypothesis test.
- A small p-value is evidence against the null hypothesis because one would reject the null hypothesis even at small significance levels.
- A large p-value is evidence in favor of the null hypothesis

Multiple Regression Analysis: Inference

Computing p-values for t tests

- Let t_j^{act} be the actual value of the t-statistic in the sample.
- If the alternative hypothesis is $H_1 : \beta_j > b_j$,

$$p - value = \mathcal{P} \left(t_j > t_j^{act} \right).$$

- If the alternative hypothesis is $H_1 : \beta_j < b_j$,

$$p - value = \mathcal{P} \left(t_j < t_j^{act} \right).$$

- If the alternative hypothesis is $H_1 : \beta_j \neq b_j$

$$p - value = \mathcal{P} \left(|t_j| > |t_j^{act}| \right).$$

- **Rejection rule:** If $p - value < \alpha$, we reject the null hypothesis.

Multiple Regression Analysis: Inference

Example: We are studying the returns to education at junior colleges and four year colleges (universities) and we have the model

$$\log(wages) = \beta_0 + \beta_1 jc + \beta_2 univ + \beta_3 exper + u,$$

where:

- jc = number of years attending a two year college
- $univ$ = number of years at a four year college
- $exper$ = months in workforce
- Data set taken from Kane and Rouse, 1995, "Labor Market Returns to Two- and Four-Year College", American Economic Review 85, 600-614. Sample size $n = 6,763$.

Multiple Regression Analysis: Inference

Running a regression of $\log(wages)$ on jc , $univ$ and $exper$ we obtain:

	Estimate	Std. Err.	t-Ratio	p-Value
<i>Intercept</i>	1.47233	0.02106	69.911	0
<i>exper</i>	0.00494	0.00016	30.901	0
<i>jc</i>	0.0667	0.00683	9.765	0
<i>univ</i>	0.07688	0.00231	33.28	0

$$n = 6763, R^2 = 0.2224.$$

This is the typical output of a software in a regression model. The p-value computed in this table is for the hypothesis $H_0 : \beta_j = 0$ vs $H_1 : \beta_j \neq 0$.

Given that for $j = 1, 2, 3$, $p - value < 0.05$ we reject H_0 in favour of H_1 at 5% level.

The regressors *exper*, *jc* and *univ* are individually significant at 5% level.

Multiple Regression Analysis: Inference

Confidence Intervals

- Another way to use classical statistical testing is to construct a confidence interval using the same critical value as was used for a two-sided test.
- Using

$$\frac{\hat{\beta}_j - \beta_j}{se(\hat{\beta}_j)} \sim t(n - k - 1),$$

we have

$$\mathcal{P}(-t_{\alpha/2} < \frac{\hat{\beta}_j - \beta_j}{se(\hat{\beta}_j)} < t_{\alpha/2}) = 1 - \alpha,$$

where $t_{\alpha/2}$ the constant that satisfies $\mathcal{P}(t_j < -t_{\alpha/2}) = \alpha/2$, where t_j is a random variable with distribution $t(n - k - 1)$. Equivalently $\mathcal{P}(t_j > t_{\alpha/2}) = \alpha/2$.

Multiple Regression Analysis: Inference

Confidence Intervals

- Now notice that

$$\begin{aligned}\mathcal{P}(-t_{\alpha/2} < \frac{\hat{\beta}_j - \beta_j}{se(\hat{\beta}_j)} < t_{\alpha/2}) &= \mathcal{P}(-t_{\alpha/2}se(\hat{\beta}_j) < \hat{\beta}_j - \beta_j < t_{\alpha/2}se(\hat{\beta}_j)), \\ &= \mathcal{P}(-\hat{\beta}_j - t_{\alpha/2}se(\hat{\beta}_j) < -\beta_j < -\hat{\beta}_j + t_{\alpha/2}se(\hat{\beta}_j)) \\ &= \mathcal{P}(\hat{\beta}_j - t_{\alpha/2}se(\hat{\beta}_j) < \beta_j < \hat{\beta}_j + t_{\alpha/2}se(\hat{\beta}_j))\end{aligned}$$

Therefore

$$\mathcal{P}(\hat{\beta}_j - t_{\alpha/2}se(\hat{\beta}_j) < \beta_j < \hat{\beta}_j + t_{\alpha/2}se(\hat{\beta}_j)) = 1 - \alpha$$

Multiple Regression Analysis: Inference

Confidence Intervals

- Hence a $100(1 - \alpha)\%$ confidence interval is defined as

$$(\hat{\beta}_j - t_{\alpha/2}se(\hat{\beta}_j), \hat{\beta}_j + t_{\alpha/2}se(\hat{\beta}_j)),$$

- In repeated samples, the interval that is constructed in the above way will cover the population regression coefficient in $100(1 - \alpha)\%$ of the cases. The interval that we compute with the actual sample is one of these intervals
- Relationship between confidence interval and hypotheses tests:

$b_j \notin \text{conf. interval} \Rightarrow \text{reject } H_0 : \beta_j = b_j \text{ in favour of } H_1 : \beta_j \neq b_j$
at $100\alpha\%$ level.

Multiple Regression Analysis: Inference

Confidence Intervals

Example: Running a regression of $\log(wages)$ on jc , $univ$ and $exper$ we obtain:

	Estimate	Std. Err.	t-Ratio	p-Value
<i>Intercept</i>	1.47233	0.02106	69.911	0
<i>exper</i>	0.00494	0.00016	30.901	0
<i>jc</i>	0.0667	0.00683	9.765	0
<i>univ</i>	0.07688	0.00231	33.28	0

$$n = 6763, R^2 = 0.2224.$$

- Construct a 90% confidence interval for the coefficient of the variable *exper*.
 - $df = n - k - 1 = 6763 - 4 = 6759$
 - $t_{0.05} = 1.645$.
 - $(0.00494 - 1.645 \times 0.00016, 0.00494 + 1.645 \times 0.00016)$
 - $(0.0046768, 0.0052032)$

Multiple Regression Analysis: Inference

Confidence Intervals

- Construct a 95% confidence interval for the coefficient of the variable jc .
- $t_{0.025} = 1.96$.
- $(0.0667 - 1.96 \times 0.00683, 0.0667 + 1.96 \times 0.00683)$
- $(0.0533132, 0.0800868)$
- Construct a 99% confidence interval for the coefficient of the variable $univ$.
- $t_{0.005} = 2.576$.
- $(0.07688 - 2.576 \times 0.00231, 0.07688 + 2.576 \times 0.00231)$
- $(0.07092944, 0.08283056)$

Multiple Regression Analysis: Inference

Testing a Linear Combination

- Suppose instead of testing whether β_1 is equal to a constant, you want to test if it is equal to another parameter, that is $H_0 : \beta_1 = \beta_2$.
- Use same basic procedure for forming a t statistic

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

Multiple Regression Analysis: Inference

Testing Linear Combination (cont)

Notice that the standard error of $\hat{\beta}_1 - \hat{\beta}_2$, $se(\hat{\beta}_1 - \hat{\beta}_2)$, is an estimator of the standard deviation of $\hat{\beta}_1 - \hat{\beta}_2$:

$$\sqrt{\text{Var}(\hat{\beta}_1 - \hat{\beta}_2)}$$

Since

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

an estimator for $\sqrt{\text{Var}(\hat{\beta}_1 - \hat{\beta}_2)}$ is given by

$$se(\hat{\beta}_1 - \hat{\beta}_2) = \sqrt{se(\hat{\beta}_1)^2 + se(\hat{\beta}_2)^2 - 2s_{12}},$$

where s_{12} is an estimate of $\text{Cov}(\hat{\beta}_1, \hat{\beta}_2)$.

Multiple Regression Analysis: Inference

Testing a Linear Combination (cont)

In some cases you can always restate the problem to get the test you want.

Example: Consider the model on the returns to education at junior colleges and four year colleges

$$\log(\text{wages}) = \beta_0 + \beta_1jc + \beta_2univ + \beta_3exper + u,$$

- We would like to test whether one year at a junior college is worth one year at a university, that is $H_0 : \beta_1 = \beta_2$.
- The alternative hypothesis is that a year at junior college is worth less than a year at a university. That is $H_1 : \beta_1 < \beta_2$.
- One can test H_0 by using the approach described before.
- However there is an easier way.

Multiple Regression Analysis: Inference

Testing a Linear Combination (cont)

Define a new parameter $\theta = \beta_1 - \beta_2$. Hence the null hypothesis becomes

$$H_0 : \theta = 0$$

and the alternative hypothesis becomes:

$$H_1 : \theta < 0,$$

We can always write the model in terms of θ . Under H_0 , the model is equivalent to

$$\log(wages) = \beta_0 + \theta jc + \beta_2 totcoll + \beta_3 exper + u,$$

where $totcoll = jc + univ$.

This model is linear in the parameters so one can use the usual tests on hypothesis for single parameters described before.

Multiple Regression Analysis: Inference

Testing a Linear Combination (cont)

- Running the regression of $\log(wages)$ on $exper$, jc and $totcoll$ we obtain:

	Estimate	Std.Err.	t-ratio
<i>Intercept</i>	1.47233	0.02106	69.911
<i>exper</i>	0.00494	0.00016	30.901
<i>jc</i>	-0.01018	0.00694	-1.467
<i>totcoll</i>	0.07688	0.00231	33.28

$$n = 6763, R^2 = 0.2224$$

Test $H_0 : \theta = 0$ vs $H_1 : \theta < 0$ (use the 5% significance level).

Multiple Regression Analysis: Inference

Testing a Linear Combinations (cont)

Example (cont):

- This is the same model as originally, but now you get a standard error for $\hat{\beta}_1 - \hat{\beta}_2$ directly from the basic regression
- Any linear combination of parameters could be tested in a similar manner
- Other examples of hypotheses about a single linear combination of parameters: $\beta_1 = 1 + \beta_2$; $\beta_1 = 5\beta_2$; $\beta_1 = -(1/2)\beta_2$; *etc.*

Multiple Regression Analysis: Inference

Multiple Linear Restrictions

- Everything we've done so far has involved testing a single linear restriction, (e.g. $\beta_1 = 0$ or $\beta_1 = \beta_2$)
- However, we may want to jointly test multiple hypotheses about our parameters.
- A typical example is testing “exclusion restrictions” – we want to know if a group of parameters are all equal to zero.

Multiple Regression Analysis: Inference

Testing Exclusion Restrictions

- Now the null hypothesis might be something like $H_0 : \beta_1 = 0, \dots, \beta_q = 0$ in the model

$$y = \beta_0 + \beta_1 x_1 + \dots + \beta_q x_q + \dots + \beta_k x_k + u.$$

That is, we want to test whether the parameters of the first q regressors (x_1 to x_q) are equal to zero.

- The alternative is just H_1 : At least one of the $\beta_j \neq 0, j = 1, \dots, q$.
- Can't just check each t statistic separately, because we want to know if the q parameters are jointly significant at a given level – it is possible for none to be individually significant at that level.

Multiple Regression Analysis: Inference

Exclusion Restrictions (cont)

- To do the test we need to estimate the “restricted model” without x_1, \dots, x_q included, as well as the “unrestricted model” with all x 's included and compute

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

where SSR_r is the sum of squared residuals of the restricted model and SSR_{ur} is the sum of squared residuals of the unrestricted model.

- Intuitively, we want to know if the change in SSR is big enough to warrant inclusion of x_1, \dots, x_q .

Multiple Regression Analysis: Inference

The F statistic

- The F statistic is always positive, since the SSR from the restricted model can't be less than the SSR from the unrestricted.
- Essentially the F statistic is measuring the relative increase in SSR when moving from the unrestricted to restricted model.
- q = number of restrictions, or $df_r - df_{ur}$.
- $n - k - 1 = df_{ur}$.
- $n - k - 1 + q = df_r$.

Multiple Regression Analysis: Inference

The F statistic (cont)

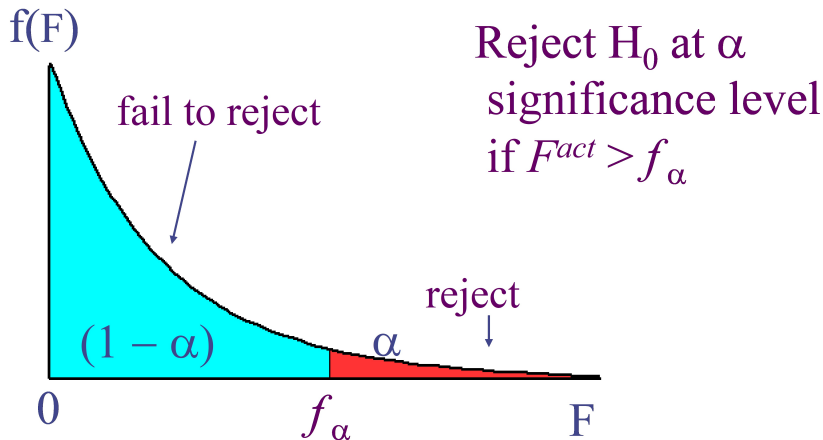
- To decide if the increase in SSR when we move to a restricted model is “big enough” to reject the exclusions, we need to know about the sampling distribution of our F statistic.
- $F \sim F(q, n - k - 1)$, where q is referred to as the numerator degrees of freedom and $n - k - 1$ as the denominator degrees of freedom.
- Denote F^{act} the actual value of the statistic in a given sample.
- The critical value is denoted as f_α and corresponds to the constant that satisfies

$$\mathcal{P}(F > f_\alpha) = \alpha.$$

Multiple Regression Analysis: Inference

The F statistic (cont)

- **Rejection rule:** Reject H_0 if $F^{act} > f_\alpha$.



Multiple Regression Analysis: Inference

Exclusion Restrictions (cont)

Example: Consider the following model that explains major league baseball players' salaries:

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

where

- *salary* = salary of major league baseball player
- *years* = Years in the league
- *gamesyr* = Average number of games per year
- *bavg* = Batting average
- *hrunsyr* = Home runs per year
- *rbisyr* = Runs batted in per year

We would like to test $H_0 : \beta_3 = 0, \beta_4 = 0, \beta_5 = 0$ vs $H_1 : H_0$ is not true.

Multiple Regression Analysis: Inference

Exclusion Restrictions (cont)

- Estimating the unrestricted model we obtain

$$\begin{aligned}\widehat{\log(\text{salary})} &= 11.19 + 0.0689\text{years} + 0.0126\text{gamesyr} \\ &\quad (0.29) \quad (0.0121) \quad (0.0026) \\ &\quad + 0.00098\text{bavg} + 0.0144\text{hrunsyr} + 0.0108\text{rbisyr}, \\ &\quad (0.00110) \quad (0.0161) \quad (0.0072) \\ n &= 353, SSR = 183.186, R^2 = 0.6278\end{aligned}$$

- Estimating the restricted model we obtain

$$\begin{aligned}\widehat{\log(\text{salary})} &= 11.22 + 0.0713\text{years} + 0.0202\text{gamesyr}, \\ &\quad (0.11) \quad (0.0125) \quad (0.0013) \\ n &= 353, SSR = 198.311, R^2 = 0.5971.\end{aligned}$$

- Test $H_0 : \beta_3 = 0, \beta_4 = 0, \beta_5 = 0$ vs $H_1 : H_0$ is not true at 5% level

Multiple Regression Analysis: Inference

The R^2 form of the F statistic

- Because the SSR 's may be large and unwieldy, an alternative form of the formula is useful.
- We use the fact that $SSR = SST(1 - R^2)$ for any regression, so can substitute in for SSR_r and SSR_{ur} :

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

where R_r^2 is the R^2 of the restricted model and R_{ur}^2 is the R^2 of the unrestricted model.

Example: For the baseball salary example, use (1) to obtain the F statistic.

Multiple Regression Analysis: Inference

Overall Significance

- A special case of exclusion restrictions is to test $H_0 : \beta_1 = \beta_2 = \dots = \beta_k = 0$.
- Since the R^2 from a model with only an intercept will be zero, the F statistic is simply

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

Multiple Regression Analysis: Inference

- **Example:** Consider the estimated model

$$\begin{aligned}\widehat{\log(\text{salary})} &= 11.19 + 0.0689\text{years} + 0.0126\text{gamesyr} \\ &\quad (0.29) \quad (0.0121) \quad (0.0026) \\ &\quad + 0.00098\text{bavg} + 0.0144\text{hrunsyr} + 0.0108\text{rbisyr}, \\ &\quad (0.00110) \quad (0.0161) \quad (0.0072) \\ n &= 353, SSR = 183.186, R^2 = 0.6278\end{aligned}$$

We would like to test

$$H_0 : \beta_1 = 0, \beta_2 = 0, \beta_3 = 0, \beta_4 = 0, \beta_5 = 0$$

vs

$$H_1 : H_0 \text{ not true}$$

at 5% level.

Multiple Regression Analysis: Inference

General Linear Restrictions

- The basic form of the F statistic will work for any set of linear restrictions.
- First estimate the unrestricted model obtain SSR_{ur} and then estimate the restricted model and obtain SSR_r .
- The F statistic as the usual form

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

where q is the number of restrictions being tested.

- Imposing the restrictions can be tricky – will likely have to redefine variables again.

Multiple Regression Analysis: Inference

General Linear Restrictions

Example: Test whether house price assessments are rational

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

- price = Actual house price
- assess = The assessed housing value before the house was sold
- lotsize = Size of lot (in feet)
- sqrft = Square footage
- bdrms = number of bedrooms

Multiple Regression Analysis: Inference

General Linear Restrictions

- Now, suppose we would like to test whether the assessed housing price is a rational valuation. If this is the case, then a 1% change in assess should be associated with a 1% change in price; that is, $\beta_1 = 1$. In addition, *lotsize*, *sqrft*, and *bdrms* should not help to explain $\log(\text{price})$, once the assessed value has been controlled for.
- Hence we want to test $H_0: \beta_1 = 1, \beta_2 = 0, \beta_3 = 0, \beta_4 = 0$ vs $H_1: H_0$ not true
- Sample size: 88.
- Running the regression of $\log(\text{price})$ on $\log(\text{assess})$, $\log(\text{lotsize})$, $\log(\text{sqrft})$ and *bdrms* we obtain $SSR_{ur} = 1.822$
- Imposing the restriction given by H_0 we have

$$\log(\text{price}) - \log(\text{assess}) = \beta_0 + u.$$

- Estimating the parameter of this model by OLS we obtain $SSR_r = 1.88$.
- Test $H_0: \beta_1 = 1, \beta_2 = 0, \beta_3 = 0, \beta_4 = 0$ vs $H_1: H_0$ not true at 5% level.

Prediction for the conditional mean of y

Suppose that we want an estimate of

$$E(y|x_1 = x_{1,0}, \dots, x_k = x_{k,0}) = \beta_0 + \beta_1 x_{1,0} + \dots + \beta_k x_{k,0}.$$

That is, we would like to estimate the the mean of y when the regressors are equal to known values $x_{1,0}, \dots, x_{k,0}$.

- This is easy to obtain by substituting the x 's in our estimated model with x_0 's ,

$$\hat{y}_0 = \hat{\beta}_0 + \hat{\beta}_1 x_{1,0} + \dots + \hat{\beta}_k x_{k,0}.$$

- We would like to construct confidence intervals for $E(y|x_1 = x_{1,0}, \dots, x_k = x_{k,0})$.
- But what about a standard error of \hat{y}_0 , ?
- There is general formula for this standard error in the case $k > 1$, but it requires knowledge of matrix algebra. However there is a simple way to obtain this standard error.
- Let us change notation and define $\theta = E(y|x_1 = x_{1,0}, \dots, x_k = x_{k,0})$.
- Thus now the objective becomes to construct a confidence interval for θ .
- θ is just a linear combination of the parameters.

Prediction for the conditional mean of y

- Can rewrite

$$\beta_0 + \beta_1 x_{1,0} + \dots + \beta_k x_{k,0} = \theta$$

as

$$\beta_0 = \theta - \beta_1 x_{1,0} - \dots - \beta_k x_{k,0}$$

- Substitute in

$$y = \beta_0 + \beta_1 x_1 + \dots + \beta_k x_k + u, \quad u \sim N(0, \sigma^2)$$

to obtain

$$y = \theta + \beta_1 (x_1 - x_{1,0}) + \dots + \beta_k (x_k - x_{k,0}) + u$$

- So, if you regress y on $(x_j - x_{j,0}), j = 1, \dots, k$, the intercept will give the predicted value and its standard error.
- Hence constructing a confidence interval for θ is similar to constructing a confidence interval for a parameter.
- $se(\hat{y}_0)$ is the standard error of the intercept in the regression of y on an intercept and $(x_j - x_{j,0}), j = 1, \dots, k$.

Prediction for the conditional mean of y

Remark: In the simple regression model we have

$$y = \beta_0 + \beta_1 x + u, \quad E(u|x) = 0, \quad \text{var}(u|x) = \sigma^2$$

Suppose that we would like to predict the value of

$$E(y|x = x_0) = \beta_0 + \beta_1 x_0$$

In this case

$$\text{se}(\hat{y}_0)^2 = \hat{\sigma}^2 \left[\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \right]$$

where $\hat{\sigma}^2 = \sum_{i=1}^n \hat{u}_i^2 / (n - 2)$ (recall that $k = 1$ in the simple regression model).

Prediction for the conditional mean of y in the multiple regression model

Example: Consider the following equation:

$$y_i = \beta_1 + \beta_2 x_i + u_i, i = 1, \dots, 60$$

The results from estimating this equation using 60 observations by Ordinary Least Squares were (standard errors in parentheses) are:

$$\hat{y} = \underset{(0.125)}{0.395} - \underset{(0.189)}{0.550}x,$$

$$SSR = 42.307, SSE = 6.1771,$$

$$S_x^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 = 0.34033$$

Given that $x_0 = 0.075$, the sample mean of x is 0.105 and that $u \sim N(0, \sigma^2)$, calculate the 95% confidence intervals for $E(y|x = x_0)$

Prediction for y

Suppose now that we would like to construct a confidence interval for y when the regressors are equal to known values $x_{1,0}, \dots, x_{k,0}$ and denote this value as y_0 .

- How can we construct a confidence interval for y_0 ?
- Notice that

$$y_0 = \beta_0 + \beta_1 x_{1,0} + \dots + \beta_k x_{k,0} + u_0$$

- Our best prediction for y_0 is the regression line

$$\hat{y}_0 = \hat{\beta}_0 + \hat{\beta}_1 x_{1,0} + \dots + \hat{\beta}_k x_{k,0}$$

- The prediction error is given by

$$\begin{aligned}\hat{u}_0 &= y_0 - \hat{y}_0 \\ &= \beta_0 + \beta_1 x_{1,0} + \dots + \beta_k x_{k,0} + u_0 - \hat{y}_0\end{aligned}$$

- Therefore, as u_0 and \hat{y}_0 are independent (conditional on the regressors):

$$\begin{aligned}\text{Var}(\hat{u}_0) &= \text{Var}(u_0) + \text{Var}(\hat{y}_0) \\ &= \sigma^2 + \text{Var}(\hat{y}_0).\end{aligned}$$

Prediction for y

$$\text{Var}(\hat{u}_0) = \sigma^2 + \text{Var}(\hat{y}_0).$$

- Hence an estimator for $\text{Var}(\hat{u}_0)$ is given by

$$se_0^2 = \hat{\sigma}^2 + se(\hat{y}_0)^2,$$

where $se(\hat{y}_0)$ is the standard error of the intercept in the regression of y on $(x_j - x_{j,0}), j = 1, \dots, k$, and

$$\hat{\sigma}^2 = \sum_{i=1}^n \hat{u}_i^2 / (n - k - 1).$$

- It can be shown that if $u \sim N(0, \sigma^2)$,

$$\frac{y_0 - \hat{y}_0}{se_0} \sim t(n - k - 1)$$

- Hence the $(1 - \alpha)\%$ prediction interval for y_0 is given by

$$(\hat{y}_0 - t_{\alpha/2} se_0, \hat{y}_0 + t_{\alpha/2} se_0),$$

where $t_{\alpha/2}$ is the percentile $(1 - \alpha/2)^{th}$ of the the t distribution with $n - k - 1$ df .

Example: Suppose we have the following regression model

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \beta_4 x_3^2 + u.$$

We have a sample of 4,137 observations . The estimated model is

$$\begin{aligned}\hat{y} &= 1.493 + 0.00149 x_1 - 0.01386 x_2 - 0.06088 x_3 \\ &\quad (0.075) \quad (0.00007) \quad (0.00056) \quad (0.01650) \\ &\quad + 0.00546 x_4, \\ &\quad (0.00227) \\ \hat{\sigma} &= 0.56\end{aligned}$$

Objectives:

- Construct a 95% confidence interval for the mean of y when $x_1 = 1,200$, $x_2 = 30$ and $x_3 = 5$, $x_4 = 25$.
- Construct a 95% confidence interval for y when $x_1 = 1,200$, $x_2 = 30$, $x_3 = 5$, $x_4 = 25$.
- Define a new set of regressors:
 - $x_1^* = x_1 - 1,200$.
 - $x_2^* = x_2 - 30$.
 - $x_3^* = x_3 - 5$.
 - $x_4^* = x_4 - 25$.

Running the regression of y on these new regressors we obtain

$$\hat{y} = 2.700 + 0.00149 x_1^* - 0.01386 x_2^* - 0.06088 x_3^* + 0.00546 x_4^* \\ \begin{matrix} (0.020) & (0.00007) & (0.00056) & (0.01650) \\ & & & (0.00227) \end{matrix}$$
$$\hat{\sigma} = 0.56$$

Predicting y in a log model

Suppose that we have the model

$$\log(y) = \beta_0 + \beta_1 x_1 + \dots + \beta_k x_k + u,$$

$E(u|x_1, \dots, x_k) = 0$, $Var(u|x_1, \dots, x_k) = \sigma^2$ and we would like to predict the mean of y for any value of the regressors: $E(y|x_1, \dots, x_k)$.

What can we do?

Given the OLS estimators the predicted value for the mean of $\log(y)$ for any values of the regressors is

$$\widehat{\log(y)} = \hat{\beta}_0 + \hat{\beta}_1 x_1 + \dots + \hat{\beta}_k x_k$$

Our first guess would be to exponentiate $\widehat{\log(y)}$.

However, simple exponentiation of $\widehat{\log(y)}$ will underestimate the expected value of y as $\widehat{\log(y)}$ is an estimator of $E(\log(y)|x_1, \dots, x_k)$ and it can be shown using an inequality known as *Jensen's inequality* that

$$\exp[E(\log(y)|x_1, \dots, x_k)] \leq E(y|x_1, \dots, x_k).$$

Predicting y in a log model

If $u \sim N(0, \sigma^2)$, it can be shown that

$$E(y|x_1, \dots, x_k) = \exp\left(\frac{\sigma^2}{2}\right) \exp(\beta_0 + \beta_1 x_1 + \dots + \beta_k x_k).$$

Therefore, a simple way to predict y is

$$\hat{y} = \exp\left(\frac{\hat{\sigma}^2}{2}\right) \exp(\hat{\beta}_0 + \hat{\beta}_1 x_1 + \dots + \hat{\beta}_k x_k).$$