



```

name: <unnamed>
log: C:\Users\Conor\Documents\Conor\Grad School\TA Work\Econ 103 - Econometric
> s\STATA Work\Week 2\wk2_section_log.smcl
log type: smcl
opened on: 16 Jan 2018, 14:08:05

1 .
2 . // Demonstration STATA code for week 2
3 . // Principles of Econometrics 4th Edition
4 .
5 . set more off

6 . clear all

7 . cd "C:\Users\Conor\Documents\Conor\Grad School\TA Work\Econ 103 - Econometrics\STATA
> Work\Week 2"
C:\Users\Conor\Documents\Conor\Grad School\TA Work\Econ 103 - Econometrics\STATA Work\
> Week 2

8 .
9 . //////////////////////////////////////
> ////////////////////////////////// Question 3.8 //////////////////////////////////
> //////////////////////////////////////
>

10. * Analyze data on home sales in Baton Rouge, Louisiana in mid-2005
11. use br2.dta, clear

12.
13. // Notice that throughout this question, we are asked to use alpha = 0.01 as our
14. // critical value. Rather than hard-code this value (i.e. type 0.01), we can
15. // create a variable that we can reuse. If we need to change alpha later on, it
16. // will be much easier to just adjust this variable instead of finding 0.01
17. // (and making sure 0.01 means alpha in whatever context)
18. // Question 3.12 will use a different alpha, so to be clear in this file, I
19. // define my variable as alpha_38 (for 3.8). Other scalar variables for this
20. // section will also have a 38.
21. scalar alpha_38 = 0.01

22.
23. *****
24. *Part A: For the traditional-style houses, estimate the linear regression
25. *model PRICE = beta1 + beta2*SQFT + e. Test the null hypothesis that the slope
26. *is zero against the alternative that it is positive, using the alpha = 0.01
27. *level of significance. Follow and show all the test steps described in
28. *Chapter 3.4
29. *****
30.
31. // Note that all the exercises in this question ask that we look only at
32. // traditional-style houses. We could put " if trad == 1" into all of our
33. // reg commands, as we did last week with female-only or black-only regressions.
34. reg price sqft if trad == 1

```

Source	SS	df	MS	Number of obs	=	582
Model	2.4362e+12	1	2.4362e+12	F(1, 580)	=	1027.92
Residual	1.3746e+12	580	2.3700e+09	Prob > F	=	0.0000
				R-squared	=	0.6393
				Adj R-squared	=	0.6387
Total	3.8108e+12	581	6.5591e+09	Root MSE	=	48683

price	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]	
sqft	73.77195	2.30097	32.06	0.000	69.2527	78.2912
_cons	-28407.56	5728.161	-4.96	0.000	-39658.02	-17157.09

```

35.
36. // Alternatively, we can remove all the non-traditional observations from the
37. // data set (since we don't save over br2.dta, we can always reload to recover
38. // those observations if we need to).
39. drop if trad != 1
    (498 observations deleted)

```

```

40. reg price sqft

```

Source	SS	df	MS	Number of obs	=	582
Model	<b>2.4362e+12</b>	<b>1</b>	<b>2.4362e+12</b>	F(1, 580)	=	<b>1027.92</b>
Residual	<b>1.3746e+12</b>	<b>580</b>	<b>2.3700e+09</b>	Prob > F	=	<b>0.0000</b>
				R-squared	=	<b>0.6393</b>
				Adj R-squared	=	<b>0.6387</b>
Total	<b>3.8108e+12</b>	<b>581</b>	<b>6.5591e+09</b>	Root MSE	=	<b>48683</b>

  

price	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]	
sqft	<b>73.77195</b>	<b>2.30097</b>	<b>32.06</b>	<b>0.000</b>	<b>69.2527</b>	<b>78.2912</b>
_cons	<b>-28407.56</b>	<b>5728.161</b>	<b>-4.96</b>	<b>0.000</b>	<b>-39658.02</b>	<b>-17157.09</b>

```

41.
42. // Steps for setting up a hypothesis test:
43. // (1) Determine the null and alternative hypothesis
44. // --> Null (H0): beta2 = 0      Alternative (H1): beta2 > 0
45. //
46. // (2) Test statistic and its distribution
47. // --> t-stat: b2/se(b2)      Distribution: t with 580 degrees of freedom
48. //      for the degrees of freedom, see the residual df in the regression output,
49. //      also stored in e(df_r)
50. //
51. // (3) Select alpha and determine the rejection region
52. // --> Problem tells us to use alpha = 0.01.
53. //      Comment on STATA's invttail function:
54. // The function has inputs df and q, where df = degrees of freedom (i.e. which
55. // t distribution we're using) and q, with answer = invttail(df, q) matches t(df, an
> swer) = 1-q/
56. // Alternatively, we could use the invt function, which would give us:
57. // answer = invt(df, p) matches t(df, answer) = p
58. // where t(df, x) is the CDF for the t distribution with df degrees of freedom
59. // --> Rejection region for a right-side (beta2 > x) test is invt(1-alpha,df)
60. //      Let's call this critical value tc
61.
62. scalar tc_38_lside = invttail(e(df_r),alpha_38) // = -1*invt(e(df_r), alpha_38)

63. disp tc_38_lside // should be positive since we're doing a right-sided (beta2>0) tes
> t
2.3327943

64.
65. // (4) Calculate the sample value of the test statistic
66.
67. scalar tstat_38a = _b[sqft]/_se[sqft]

68. disp tstat_38a
32.061242

69.

```

```

70. // --> Note that this t-stat value matches what was reported in the regression output
> t
71. //      since the null is the same (beta2 = 0)
72. //
73. // (5) State your conclusion
74. // --> Given the t-statistic of about 32 compared to a critical value of about 2.3
75. //      we reject the null that beta2 = 0
76.
77. *****
78. *Part B: Using the linear model in (a), test the null hypothesis (H0) that the
79. *expected price of a house of 2000 square feet is equal to, or less than,
80. *$120,000. What is the appropriate alternative hypothesis? Use the alpha = 0.01
81. *level of significance. Obtain the p-value of the test and show its value on a
82. *sketch. What is your conclusion?
83. *****
84.
85. // Null (H0): yhat(2000) = beta1 + 2000*beta2 = 120,000
86. // Alternative (H1): yhat(2000) > 120,000
87. // Again, we'll use a t test for this exercise.
88.
89. // First calculate the point estimate for yhat(2000)
90. scalar point_yhat_2000 = _b[_cons] + 2000*_b[sqft]
91. disp point_yhat_2000
119136.34
92.
93. // Notice that the point estimate is less than our null of 120000. Already,
94. // this tells us that we will fail to reject if our alternative is that
95. // yhat(2000) > 120000 since the rejection region will only cover positive
96. // t statistics, which will require yhat(2000) > 120000
97.
98. // Next, we need to calculate the standard error of our estimate.
99. // Since yhat is a linear combination of b1 and b2, we need information on
100 // the full variance-covariance matrix from our OLS regression
101 // We can view the variance-covariance matrix using the following command:
102 estat vce

Covariance matrix of coefficients of regress model

      e(V) |      sqft      _cons
-----|-----
      sqft | 5.2944617
      _cons | -12335.341  32811823

103 // Note that the diagonal terms (sqft, sqft and _cons, _cons) are equal to
104 // the square of the Std. Err. reported in the output table
105
106 // Next, rather than code the values from variance-covariance matrix terms by
107 // hand, we can use STATA's stored values to enter these numbers
108 // Just as we've used scalars before, now we need to define a matrix object,
109 // and use it to store a saved value from the regression.
110 matrix define varmat_38a = e(V)

111 // Next, to extract elements from a matrix, we use the syntax matname[i, j]
112 // where we want to find the row i and column j value of matrix matname
113 // From here, we can calculate the standard error of our estimate of yhat
114 // as follows: se_yhat_2000 = sqrt( var(_cons) + (2000^2)*var(sqft) + 2*2000*cov(_co
> ns,sqft))
115 scalar se_yhat_2000 = sqrt(varmat_38a[2,2] + (2000^2)*varmat_38a[1,1] + 2*2000*varma
> t_38a[2,1])

```

```

116 disp se_yhat_2000
2155.9934

117
118 // Calculate the t stat for the null of yhat(2000) = 120000
119 scalar tstat_38b = (point_yhat_2000 - 120000)/se_yhat_2000

120
121 // Compare to the critical value for the 1-sided t-test
122 // This is the same critical value tc that we had in part (a)
123 disp tstat_38b
-.4005866

124 disp tc_38_lside
2.3327943

125
126 // Given this estimate, we fail to reject the null that yhat(2000) <= 120000
127
128 // We still need to calculate the p-value for this test. To be clear about the
129 // p-value that I'm calculating, I note in the name that we are looking at a test
130 // where rejecting the null requires values to fall on the right side of the
131 // distribution
132 scalar pval_yhat_2000_rt = 1 - t(e(df_r), tstat_38b) // could also do ttail(e(df_r),
> tstat_38b)

133 disp pval_yhat_2000_rt
.65556399

134
135 // Comparing this p-value to our alpha, we once again fail to reject the null
136 // that yhat(2000) <= 120000.
137
138 ////////////////////////////////// STATA Command Alternatives //////////////////////////////////
>
139 // For comparison, I also calculate the p-value for the 2-sided test, which
140 // corresponds to the automatic output of the STATA commands discussed below.
141 // Note for the 2-sided test, to get the correct value we need to make sure that
142 // the input for t CDF is a negative value. I use the if, else syntax for STATA
143 // to assign the value of check_sign to make sure I get the correct result.
144 if tstat_38b < 0 {
145     scalar check_sign = 1
146 }

147     else {
148         scalar check_sign = -1
149 }

150
151 scalar pval_yhat_2000_2side = 2*t(e(df_r), check_sign*tstat_38b)

152
153 // In addition to calculating the test statistics by hand, we could have STATA
154 // do a bunch of the work for us using the following commands: lincom or test
155 // (1) lincom calculates the point estimate and standard error for a linear
156 // combination of our beta estimates. It also reports a confidence interval,
157 // together with the t-test and p-value for the null that the calculated
158 // value is equal to zero.
159 // syntax: lincome exp
160 // where exp = c1*var1 + c2*var2 or c1*var2 - c2*var2

```

```

161 // note that var1 and var2 are the names of RHS variables used in the most recent
162 // regression, as reported in the STATA output. lincom stores results in return list
163 lincom _cons + 2000*sqft

```

```
( 1) 2000*sqft + _cons = 0
```

price	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]	
(1)	<b>119136.3</b>	<b>2155.993</b>	<b>55.26</b>	<b>0.000</b>	<b>114901.8</b>	<b>123370.8</b>

```

164
165 // We could also feed in _cons + 2000*sqft - 120000 to lincom. Notice that this
166 // won't affect the estimate of the standard error, but rather than give us the
167 // point estimate for yhat(2000) it will give us the the numerator in our
168 // t-statistic, as well as giving us the t-stat and p-value for the 2-sided t-test
169 // that yhat(2000) = 120000 (i.e. yhat(2000) - 120000 = 0)
170 lincom _cons + 2000*sqft - 120000

```

```
( 1) 2000*sqft + _cons = 120000
```

price	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]	
(1)	<b>-863.6621</b>	<b>2155.993</b>	<b>-0.40</b>	<b>0.689</b>	<b>-5098.168</b>	<b>3370.844</b>

```

171 disp tstat_38b
    -.4005866

```

```

172 disp pval_yhat_2000_2side
    .68887203

```

```

173
174 // Notice that the t-statistic reported here matches the t-statistic we found
175 // by hand earlier. The p-value, however, matches that for the 2-sided test
176 // and not the right-side test.
177
178 // (2) test allows you to test (potentially many) linear combinations of the
179 //     beta estimates against values of your choosing. Although test reports an
180 //     F-statistic, in the case of testing just a single linear combination, we
181 //     can recover the t-stat by noting that in the single-restriction case
182 //     t-stat = sqrt(F-stat). This still poses the issue of what the sign of the
183 //     t-stat is though, but we know that its sign will match the sign of
184 //     the point estimate minus its null value.
185 test _cons + 2000*sqft = 120000

```

```
( 1) 2000*sqft + _cons = 120000
```

```

      F( 1, 580) =    0.16
      Prob > F =    0.6889

```

```

186 disp sqrt(r(F))
    .4005866

```

```

187 disp tstat_38b
    -.4005866

```

```

188 disp pval_yhat_2000_2side
    .68887203

```

```

189
190 // We can see that the tstat we calculated earlier matches the absolute value//
191 // of the square root of the F-stat here. As with lincom above, the automatically
192 // reported p-value corresponds to the one for the 2-sided test.
193
194 *****
195 *Part C: Based on the estimated results from part (a), construct a 95% interval
196 *estimate of the expected price of a house of 2000 square feet.
197 *****
198
199 // To construct the 2-sided t-test, we need 3 elements
200 // (1) The point estimate for our variable of interest
201 // (2) The standard error for our variable of interest
202 // (3) The critical value for the t distribution associated with a 2-sided test
203 // at our desired level of significance.
204 // From these values, we can then calculate the low- and high-points of the
205 // confidence interval as:
206 // LOW = point estimate - tc_lvl_2side * se
207 // HIGH = point estimate + tc_lvl_2side * se
208
209 scalar tc_95_2side = -1*invttail(e(df_r), 0.05/2) //alpha = 0.05, could also do invttail
    > (e(df_r),0.05/2)

210 scalar yhat_2000_ci95low = point_yhat_2000-tc_95_2side*se_yhat_2000

211 scalar yhat_2000_ci95high = point_yhat_2000+tc_95_2side*se_yhat_2000

212 disp yhat_2000_ci95low
    114901.83

213 disp yhat_2000_ci95high
    123370.84

214
215 // Compare these by-hand estimates to what STATA generated using lincom
216 lincom _cons + 2000*sqft

    ( 1)  2000*sqft + _cons = 0

```

price	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]	
(1)	119136.3	2155.993	55.26	0.000	114901.8	123370.8

```

217
218 // To adjust the confidence level in the lincom command, use the level option
219 // For example, to construct the 99% confidence interval, you would say:
220 // lincom _cons + 2000*sqft, level(99)
221
222 *****
223 *Part D: For the traditional-style houses, estimate the quadratic regression
224 *model PRICE = alpha1 + alpha2*SQFT^2 + e. Test the null hypothesis that the
225 *marginal effect of an additional square foot of living area in a home with
226 *2000 square feet of living space is $75 against the alternative that the effect
227 *is less than $75. Use the alpha = 0.01 level of significance. Repeat the same
228 *test for a home of 4000 square feet of living space. Discuss your conclusions.
229 *****
230
231 gen sqft_sqr = sqft^2

```

232 reg price sqft\_sqr

Source	SS	df	MS	Number of obs	=	582
Model	<b>2.5786e+12</b>	<b>1</b>	<b>2.5786e+12</b>	F(1, 580)	=	<b>1213.74</b>
Residual	<b>1.2322e+12</b>	<b>580</b>	<b>2.1245e+09</b>	Prob > F	=	<b>0.0000</b>
				R-squared	=	<b>0.6767</b>
				Adj R-squared	=	<b>0.6761</b>
Total	<b>3.8108e+12</b>	<b>581</b>	<b>6.5591e+09</b>	Root MSE	=	<b>46093</b>

  

price	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]	
sqft_sqr	<b>.0120632</b>	<b>.0003463</b>	<b>34.84</b>	<b>0.000</b>	<b>.0113832</b>	<b>.0127433</b>
_cons	<b>68710.05</b>	<b>2873.195</b>	<b>23.91</b>	<b>0.000</b>	<b>63066.91</b>	<b>74353.18</b>

```

233 // Marginal effect is 2*x*b2 where x = 2000 or 4000
234
235 // Calculate by hand: sqft = 2000
236 scalar margeff_quad_2000 = 2*_b[sqft_sqr]*2000

237 scalar se_margeff_quad_2000 = 2*2000*_se[sqft_sqr]

238 scalar tstat_38d_2000 = (margeff_quad_2000 - 75)/se_margeff_quad_2000

239 // Compare tstat to the critical value for a left-hand side test at alpha = 0.01
240 disp tstat_38d_2000
-19.311502

241 disp -1*tc_38_lside
-2.3327943

242 // Conclusion: since our t-statistic is less than the critical value for the
243 // LHS test, we reject the null that marginal effect is $75 at sqft=2000 in favor
244 // of the alternative that the marginal effect is less than $75
245
246 // Calculate by hand: sqft = 4000
247 scalar margeff_quad_4000 = 2*_b[sqft_sqr]*4000

248 scalar se_margeff_quad_4000 = 2*4000*_se[sqft_sqr]

249 scalar tstat_38d_4000 = (margeff_quad_4000 - 75)/se_margeff_quad_4000

250 // Compare tstat to the critical value for a left-hand side test at alpha = 0.01
251 disp tstat_38d_4000
7.7636657

252 disp -1*tc_38_lside
-2.3327943

253 // Conclusion: fail to reject the null that the marginal effect is $75 at
254 // sqft = 4000
255
256
257 ////////////////////////////////////// STATA Command Alternatives //////////////////////////////////////
>
258 // Marginal effect is 2*b2*x where x = 2000 or 4000 and null is = 75
259 // Compare the t-stat to what we calculated above
260 lincom 2*sqft_sqr*2000 - 75

```

( 1) **4000\*sqft\_sqr = 75**

price	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]	
(1)	<b>-26.74707</b>	<b>1.385033</b>	<b>-19.31</b>	<b>0.000</b>	<b>-29.46736</b>	<b>-24.02678</b>

```
261 lincom 2*sqft_sqr*4000 - 75
```

```
( 1) 8000*sqft_sqr = 75
```

price	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]	
(1)	<b>21.50587</b>	<b>2.770066</b>	<b>7.76</b>	<b>0.000</b>	<b>16.06528</b>	<b>26.94645</b>

```
262
```

```
263 // Using the # syntax in regress together with the margins command
```

```
264 reg price c.sqft#c.sqft
```

Source	SS	df	MS	Number of obs	=	582
Model	<b>2.5786e+12</b>	<b>1</b>	<b>2.5786e+12</b>	F(1, 580)	=	<b>1213.74</b>
Residual	<b>1.2322e+12</b>	<b>580</b>	<b>2.1245e+09</b>	Prob > F	=	<b>0.0000</b>
				R-squared	=	<b>0.6767</b>
				Adj R-squared	=	<b>0.6761</b>
Total	<b>3.8108e+12</b>	<b>581</b>	<b>6.5591e+09</b>	Root MSE	=	<b>46093</b>

price	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]	
c.sqft#c.sqft	<b>.0120632</b>	<b>.0003463</b>	<b>34.84</b>	<b>0.000</b>	<b>.0113832</b>	<b>.0127433</b>
_cons	<b>68710.05</b>	<b>2873.195</b>	<b>23.91</b>	<b>0.000</b>	<b>63066.91</b>	<b>74353.18</b>

```
265 margins, dydx(sqft) at(sqft=(2000 4000))
```

```
Conditional marginal effects      Number of obs      =      582
Model VCE      : OLS
```

```
Expression      : Linear prediction, predict()
dy/dx w.r.t.    : sqft
```

```
1._at          : sqft          =      2000
```

```
2._at          : sqft          =      4000
```

	Delta-method					
	dy/dx	Std. Err.	t	P> t	[95% Conf. Interval]	
sqft						
at						
1	48.25293	1.385033	34.84	0.000	45.53264	50.97322
2	96.50586	2.770066	34.84	0.000	91.06528	101.9464

```
266
```

```
267 // Compare point estimates and standard errors to what we calculated earlier
```

```
268 matrix define margins_est = r(b) // r(b) is a matrix of point estimates, in order
```

```
269
```

```
> gins command
```

```
270 matrix define var_margins_est = r(V) // r(V) is the variance-covariance matrix
```

```
271
```

```
> reported point estimates, so that
```

```
// for the
```



```

272                                     // the vari
> ance of an individual estimate
273                                     // is on th
> e diagonal of the r(V) matrix
274
275 disp (margins_est[1,1]-75)/sqrt(var_margins_est[1,1]) // t-stat at sqft=2000
-19.311503

276 disp (margins_est[1,2]-75)/sqrt(var_margins_est[2,2]) // t-stat at sqft=4000
7.7636653

277
278 *****
279 *Part E: For the traditional-style houses, estimate the log-linear regression
280 *model ln(PRICE) = gamma1 + gamma2*SQFT + e. Test the null hypothesis that the
281 *marginal effect of an additional square foot of living area in a home with
282 *2000 square feet of living space if $75 against the alternative that the effect
283 *is less than $75. Use the alpha = 0.01 level of significance. Repeat the same
284 *test for a home of 4000 square feet of living space. Discuss your conclusions.
285 *****
286
287 gen ln_price = log(price) // generate variable for log price

288 reg ln_price sqft // log-linear regression

```

Source	SS	df	MS	Number of obs	=	582
Model	76.4414878	1	76.4414878	F(1, 580)	=	880.41
Residual	50.3587111	580	.086825364	Prob > F	=	0.0000
				R-squared	=	0.6028
				Adj R-squared	=	0.6022
Total	126.800199	581	.218244749	Root MSE	=	.29466

  

ln_price	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]
sqft	.0004132	.0000139	29.67	0.000	.0003859 .0004406
_cons	10.79894	.0346705	311.47	0.000	10.73084 10.86703

```

289
290 // Calculate marginal effect by hand using formula exp(b1 + b2*x)*b2
291 // where x = 2000 or 4000
292 scalar margeff_log_2000 = exp(_b[_cons]+_b[sqft]*2000)*_b[sqft]

293 scalar margeff_log_4000 = exp(_b[_cons]+_b[sqft]*4000)*_b[sqft]

294 matrix define varmat_loglin = e(V)

295
296 // Since we're using a non-linear transformation of the beta values, we need to
297 // use some different techniques, which are incorporated into the following
298 // functions: nlcom and testnl
299 // nlcom is analogous to lincom, while testnl is analogous to test
300 //
301 // You will notice below that in the output for nlcom STATA refers to a z value
302 // which reminds us that (1) this result is relying on completely asymptotic
303 // results, rather than making a finite sample adjustment, and (2) that we
304 // should compare that value to a normal distribution. Similarly, testnl reports
305 // results for a chi2 statistic, rather than an F-distribution.
306 // Once again, the reported p-values correspond to a 2-sided test.

```

```

307 //
308 // The underlying method for working with non-linear functions of the coefficients
309 // is called the delta method, and it is discussed briefly in section 5.6.3 of
310 // the textbook. The underlying result is that we can estimate the variance for
311 // a non-linear function of the data as follows:
312 //  $\text{Lambda} = f(b_1, b_2)$ 
313 //  $\text{Var}(\text{Lambda}) = (d\text{Lambda}/db_1)^2 * \text{var}(b_1) + (d\text{Lambda}/db_2)^2 * \text{var}(b_2)$ 
314 //  $+ 2 * (d\text{Lambda}/db_1) * (d\text{Lambda}/db_2) * \text{cov}(b_1, b_2)$ 
315 //
316 // Below we show how to do this by hand, and then use nlcom and testnl:
317
318 // Calculating variance of marginal effect estimate by hand
319 scalar partial_b1_2000 = exp(_b[_cons]+_b[sqft]*2000)*_b[sqft]
320 scalar partial_b2_2000 = exp(_b[_cons]+_b[sqft]*2000)*(1+_b[sqft]*2000)
321 scalar var_margeff_log_2000 = (partial_b1_2000^2)*varmat_loglin[2,2] + ///
> (partial_b2_2000^2)*
> varmat_loglin[1,1] + ///
> (2*partial_b1_2000*p
> artial_b2_2000)*varmat_loglin[2,1]
322
323 scalar partial_b1_4000 = exp(_b[_cons]+_b[sqft]*4000)*_b[sqft]
324 scalar partial_b2_4000 = exp(_b[_cons]+_b[sqft]*4000)*(1+_b[sqft]*4000)
325 scalar var_margeff_log_2000 = (partial_b1_4000^2)*varmat_loglin[2,2] + ///
> (partial_b2_4000^2)*
> varmat_loglin[1,1] + ///
> (2*partial_b1_4000*p
> artial_b2_4000)*varmat_loglin[2,1]
326
327 // Pick the critical value for our test. Should be negative since we are doing
328 // a left-side test.
329 scalar zc_95_lt = invnormal(0.01)
330
331 // Calculate the point estimates using nlcom
332 nlcom exp(_b[_cons]+_b[sqft]*2000)*_b[sqft]

```

```
_nl_1: exp(_b[_cons]+_b[sqft]*2000)*_b[sqft]
```

ln_price	Coef.	Std. Err.	z	P> z	[95% Conf. Interval]	
_nl_1	<b>46.2433</b>	<b>1.459765</b>	<b>31.68</b>	<b>0.000</b>	<b>43.38221</b>	<b>49.10439</b>

```
333 nlcom exp(_b[_cons]+_b[sqft]*4000)*_b[sqft]
```

```
_nl_1: exp(_b[_cons]+_b[sqft]*4000)*_b[sqft]
```

ln_price	Coef.	Std. Err.	z	P> z	[95% Conf. Interval]	
_nl_1	<b>105.677</b>	<b>3.663876</b>	<b>28.84</b>	<b>0.000</b>	<b>98.49592</b>	<b>112.858</b>

```

334
335 // Calculate the t value for null = 75 using nlcom
336 nlcom exp(_b[_cons]+_b[sqft]*2000)*_b[sqft]-75

```

```

      _nl_1:  exp(_b[_cons]+_b[sqft]*2000)*_b[sqft]-75

```

ln_price	Coef.	Std. Err.	z	P> z	[95% Conf. Interval]	
_nl_1	<b>-28.7567</b>	<b>1.459765</b>	<b>-19.70</b>	<b>0.000</b>	<b>-31.61779</b>	<b>-25.89561</b>

```

337 nlcom exp(_b[_cons]+_b[sqft]*4000)*_b[sqft]-75

```

```

      _nl_1:  exp(_b[_cons]+_b[sqft]*4000)*_b[sqft]-75

```

ln_price	Coef.	Std. Err.	z	P> z	[95% Conf. Interval]	
_nl_1	<b>30.67698</b>	<b>3.663876</b>	<b>8.37</b>	<b>0.000</b>	<b>23.49592</b>	<b>37.85805</b>

```

338
339 // Compare the z-scores here to the critical value:
340 disp zc_95_lt
      -2.3263479

```

```

341
342 // Do the same exercise using testnl. Compare the sqrt of the chi2 stat to the
343 // t stat we calculated earlier.
344 testnl exp(_b[_cons]+_b[sqft]*2000)*_b[sqft]=75 // sqft = 2000

```

```

      (1)  exp(_b[_cons]+_b[sqft]*2000)*_b[sqft] = 75

```

```

              chi2(1) =      388.07
      Prob > chi2 =      0.0000

```

```

345 disp sqrt(r(chi2))
      19.699539

```

```

346 testnl exp(_b[_cons]+_b[sqft]*4000)*_b[sqft]=75 // sqft = 4000

```

```

      (1)  exp(_b[_cons]+_b[sqft]*4000)*_b[sqft] = 75

```

```

      warning: derivative with respect to sqft coefficient is near zero,
               derivative treated as zero

```

```

              chi2(1) =      70.10
      Prob > chi2 =      0.0000

```

```

347 disp sqrt(r(chi2))
      8.3728216

```

```

348
349
350 /* Discussion:
> The point estimate for the log-linear regression is quite similar to the
> point estimate in the quadratic case for sqft=2000. Unsurprisingly, in both cases
> we reject the null that the marginal effect is $75 in favor of the alternative that
> the effect is less than $75. The difference between the two approaches is larger
> in the sqft=4000 case, but in the log-linear case we once again fail to reject
> the null given that the point estimate is above $75. This discussion would be
> largely the same if we implemented the adjustment to the point estimate for the
> log-linear regression.
> */

```

```

351
352 //////////////////////////////////////
> ////////////////////////////////////// Question 3.12 //////////////////////////////////////
> //////////////////////////////////////
>
353 * How does the relationship between experience and wages change over a lifetime?
354 * How does sample size affect inference in OLS?
355 clear all

356 use cps4_small.dta, clear

357
358 scalar alpha_312 = 0.05 // set the alpha for this section at 0.05

359 *****
360 *Part A: Create a new variable called EXPER30 = EXPER - 30. Construct a scatter
361 *diagram with WAGE on the vertical axis and EXPER30 on the horizontal axis. Are
362 *any patterns evident?
363 *****
364 gen exper30 = exper - 30

365 twoway scatter wage exper30

366 graph export "Q 3-12 Wage Exper30 Scatter.pdf", replace
    (file Q 3-12 Wage Exper30 Scatter.pdf written in PDF format)

367
368 /* Discussion:
> For all levels of experience, there is a large mass in the 0-20 range for wages
> but the upper level of wages seems to have a parabolic shape (i.e. rising,
> and then falling).
> */
369
370 *****
371 *Part B: Estimate by least squares the quadratic model
372 *WAGE = gamma1+gamma2*(EXPER30)^2 + e. Are the coefficient estimates
373 *statistically significant? Test the null that gamma2 >= 0 against the
374 *alternative that gamma2 < 0 at the alpha = 0.05 level of significance. What
375 *conclusion do you draw?
376 *****
377 gen exper30_sqr = exper30^2

378 reg wage exper30_sqr

```

Source	SS	df	MS	Number of obs	=	1,000
Model	<b>7845.5768</b>	<b>1</b>	<b>7845.5768</b>	F(1, 998)	=	<b>49.96</b>
Residual	<b>156719.851</b>	<b>998</b>	<b>157.033919</b>	Prob > F	=	<b>0.0000</b>
				R-squared	=	<b>0.0477</b>
				Adj R-squared	=	<b>0.0467</b>
Total	<b>164565.428</b>	<b>999</b>	<b>164.730158</b>	Root MSE	=	<b>12.531</b>

  

wage	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]
exper30_sqr	<b>-.0138283</b>	<b>.0019564</b>	<b>-7.07</b>	<b>0.000</b>	<b>-.0176674    -.0099892</b>
_cons	<b>23.06694</b>	<b>.5265962</b>	<b>43.80</b>	<b>0.000</b>	<b>22.03358    24.10031</b>

```

379
380 scalar tc_312b_1side = invt(e(df_r),alpha_312) // critical value for left-hand test

```

```

381 // Since we are working with the null that gamma2 = 0, we can use the t-stat
382 // stata automatically reported with the regression.
383 disp tc_312b_lside
-1.6463819

```

```

384
385 /* Discussion - the estimate for the beta on exper30_sqr is statistically
> different from zero at the 95% confidence level. We can see this from the the
> fact that the P value is approximately 0 and that 0 is not in the 95% confidence
> interval. Since the 2-sided test is more aggressive on either side than the 1-sided
> test, passing the 2-sided guarantees that we will pass the 1-sided test. To be
> certain, we compare our t-stat to the critical value for the 1-sided test.
> */

```

```

386
387 *****
388 *Part C: Using the estimation in part (b), compute the estimated marginal effect
389 *of experience upon wage for a person with 10 years' experience, 30 years'
390 *experience and 50 years' experience. Are these slopes significantly different
391 *from zero at the alpha = 0.05 level of significance?
392 *****
393
394 lincom 2*exper30_sqr*(10-30) //when exper=30, exper30 = -20

```

```
( 1) - 40*exper30_sqr = 0
```

wage	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]	
(1)	<b>.5531314</b>	<b>.0782551</b>	<b>7.07</b>	<b>0.000</b>	<b>.399568</b>	<b>.7066948</b>

```

395 lincom 2*exper30_sqr*(30-30) //when exper=30, exper30 = 0

```

```
( 1) = 0
```

wage	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]	
(1)	<b>0</b>	(omitted)				

```

396 lincom 2*exper30_sqr*(50-30) //when exper=50, exper30 = 20

```

```
( 1) 40*exper30_sqr = 0
```

wage	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]	
(1)	<b>-.5531314</b>	<b>.0782551</b>	<b>-7.07</b>	<b>0.000</b>	<b>-.7066948</b>	<b>-.399568</b>

```

397
398 /* Discussion:
> Notice that for exper = 10 and exper = 50, the t-stat has the same magnitude
> as the beta coefficient itself. This follows from the fact that (1) we are only
> evaluating a scalar multiple of the beta coefficient, and (2) that the fixed
> null for the 3 tests is 0. In addition, notice that given the definition of
> exper30 and our choice of regression, we've basically assumed that the marginal
> value at exper=30 will be zero. Given this, we would fail to reject the null
> that the marginal effect is equal to zero for any finite variance.
>
> Comment: why is (absolute value of) the t-stat not affected across the three
> cases b2 = 0, 2*(-20)*b2, 2*20*b2? Recall that the t-stat is given as:
> t-stat = (b2-b2 null)/se(b2). In addition, for any constant c, se(c*b2) =
> abs(c)*se(b2) where abs(c) is the absolute value of c. Then, in the 3 cases, we have
> :
> t-stat = b2/se(b2) = -c*b2/(c*se(b2)) and c*b2/se(b2)
> */

```

```

399
400 *****
401 *Part D: Construct 95% confidence interval estimates of each of the slopes in
402 *part (c) How precisely are we estimating these values?
403 *****
404
405 // The lincom command already generated the confidence intervals for us, but
406 // here we recreate the estimates by hand. For comparison, I also show the
407 // calculations for the confidence interval for the beta coefficient itself
408
409 scalar tc_312d_2side = invttail(e(df_r),0.05/2) // = (-1)*invt(e(df_r),0.05/2)
410
411 scalar beta2 = _b[exper30_sqr] // useful for comparison later
412
413 scalar beta2_cilow = _b[exper30_sqr] - tc_312d_2side*_se[exper30_sqr]
414 scalar margeff_10_cilow = -40*_b[exper30_sqr] - tc_312d_2side*_se[exper30_sqr]*40
415 scalar margeff_50_cilow = 40*_b[exper30_sqr] - tc_312d_2side*_se[exper30_sqr]*40
416
417 scalar beta2_cihigh = _b[exper30_sqr] + tc_312d_2side*_se[exper30_sqr]
418 scalar margeff_10_cihigh = -40*_b[exper30_sqr] + tc_312d_2side*_se[exper30_sqr]*40
419 scalar margeff_50_cihigh = 40*_b[exper30_sqr] + tc_312d_2side*_se[exper30_sqr]*40
420
421 disp "Confidence Interval for effect at exper = 10: [" margeff_10_cilow ", " margeff
> _10_cihigh "]"
Confidence Interval for effect at exper = 10: [.39956799, .70669478]
422 disp "Confidence Interval for effect at exper = 50: [" margeff_50_cilow ", " margeff
> _50_cihigh "]"
Confidence Interval for effect at exper = 50: [-.70669478, -.39956799]
423
424 /* Discussion:
> Overall, we have a fairly tight estimate of the confidence intervals.
> */
425
426 *****
427 *Part E: Using the estimation result from part (b) create the fitted values
428 *WAGE_hat = gamma1_hat + gamma2_hat*(EXPER30)^2 where _hat denotes the least
429 *squares estimates. Plot these fitted values and WAGE on the vertical axis of
430 *the same graph against EXPER30 on the horizontal axis. Are the estimates in
431 *part (c) consistent with the graph?
432 *****
433 predict wage_hat, xb
434 twoway (scatter wage exper30) (line wage_hat exper30, sort)
435 graph export "Q 3-12E Fitted Values.pdf", replace
(file Q 3-12E Fitted Values.pdf written in PDF format)
436
437 /* Discussion:
> Not surprisingly, given that everything is based on the same underlying regression,
> the results in part (c) are consistent with this graph. The curve of fitted
> values is symmetric around zero, with wages tending to rise before 30 (exper30 = 0),
> and falling after 30 (exper30 = 0).
> */

```

```

438
439 *****
440 *Part F: Estimate the linear regression WAGE = beta1+beta2*EXPER30 + e and the
441 *linear regression WAGE = alpha1 + alpha2*EXPER + e. What differences do you
442 *observe between these regressions and why do they occur? What is the estimated
443 *marginal effect of experience on wage from these regressions? Based on our work
444 *in parts (b)-(d), is the assumption of constant slope in this model a good one?
445 *Explain.
446 *****
447
448 reg wage exper30

```

Source	SS	df	MS	Number of obs	=	1,000
Model	1306.16677	1	1306.16677	F(1, 998)	=	7.98
Residual	163259.261	998	163.586434	Prob > F	=	0.0048
				R-squared	=	0.0079
				Adj R-squared	=	0.0069
Total	164565.428	999	164.730158	Root MSE	=	12.79

  

wage	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]	
exper30	.0889534	.0314802	2.83	0.005	.0271785	.1507283
_cons	20.92629	.419131	49.93	0.000	20.10381	21.74876

```

449 reg wage exper

```

Source	SS	df	MS	Number of obs	=	1,000
Model	1306.16677	1	1306.16677	F(1, 998)	=	7.98
Residual	163259.261	998	163.586434	Prob > F	=	0.0048
				R-squared	=	0.0079
				Adj R-squared	=	0.0069
Total	164565.428	999	164.730158	Root MSE	=	12.79

  

wage	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]	
exper	.0889534	.0314802	2.83	0.005	.0271785	.1507283
_cons	18.25768	.9273279	19.69	0.000	16.43795	20.07742

```

450
451 /* Discussion:
452 > The only difference between the two sets of output are in the estimates tied
453 > to the constant term. With exper30, the constant is larger but with a smaller
454 > standard error. If we think about how OLS works, this should make sense since
455 > adding or subtracting a constant to x in the true model is equivalent to moving
456 > the constant term around. Mechanically, for OLS, the b2 estimate only cares about
457 > deviations in x from its sample mean, so that adding/subtracting a constant gets
458 > stripped out. Similarly, the b1 estimate moves around to ensure that the point
459 > (xbar, ybar) is on the best-fit line, so adding/subtracting from xbar just moves
460 > the b1 estimate around so that we continue to satisfy b1 = ybar - b2*xbar
461 > */
462
463 *****
464 *Part G: Use the larger data cps4.dta (4838 observations) to repeat parts (b),
465 *(c), and (d). How much has the larger sample improved the precision of the
466 *interval estimates in part (d)?
467 *****

```

458 reg wage exper30\_sqr // re-run regression so easier to directly compare output

Source	SS	df	MS	Number of obs	=	1,000
Model	<b>7845.5768</b>	<b>1</b>	<b>7845.5768</b>	F(1, 998)	=	<b>49.96</b>
Residual	<b>156719.851</b>	<b>998</b>	<b>157.033919</b>	Prob > F	=	<b>0.0000</b>
				R-squared	=	<b>0.0477</b>
				Adj R-squared	=	<b>0.0467</b>
Total	<b>164565.428</b>	<b>999</b>	<b>164.730158</b>	Root MSE	=	<b>12.531</b>

  

wage	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]	
exper30_sqr	<b>-.0138283</b>	<b>.0019564</b>	<b>-7.07</b>	<b>0.000</b>	<b>-.0176674</b>	<b>-.0099892</b>
_cons	<b>23.06694</b>	<b>.5265962</b>	<b>43.80</b>	<b>0.000</b>	<b>22.03358</b>	<b>24.10031</b>

459

460 use cps4, clear

461

462 gen exper30 = exper-30

463 gen exper30\_sqr = exper30^2

464 reg wage exper30\_sqr

Source	SS	df	MS	Number of obs	=	4,838
Model	<b>29995.3377</b>	<b>1</b>	<b>29995.3377</b>	F(1, 4836)	=	<b>198.74</b>
Residual	<b>729879.371</b>	<b>4,836</b>	<b>150.926255</b>	Prob > F	=	<b>0.0000</b>
				R-squared	=	<b>0.0395</b>
				Adj R-squared	=	<b>0.0393</b>
Total	<b>759874.709</b>	<b>4,837</b>	<b>157.096281</b>	Root MSE	=	<b>12.285</b>

  

wage	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]	
exper30_sqr	<b>-.0123931</b>	<b>.0008791</b>	<b>-14.10</b>	<b>0.000</b>	<b>-.0141165</b>	<b>-.0106696</b>
_cons	<b>22.35452</b>	<b>.2366069</b>	<b>94.48</b>	<b>0.000</b>	<b>21.89066</b>	<b>22.81838</b>

465

466 /\* Discussion:

```

> There are 3 changes between the two samples:
> (1) A different value for the point estimate of beta2 (associated with changes
> in sample variance/covariance). A different value of the point estimate of
> sigma_hat^2 (associated with changes in beta2 and betal).
> (2) Changes in estimate of var(b2). This can come from either changes in estimate of
> sigma_hat^2 or changes in sum (xi - xbar)^2 since
> var_hat(b2) = (sigma_hat^2)/(sum (xi -xbar)^2)
> (3) A higher degrees of freedom for the regression leads to smaller critical
> values for t tests (so smaller confidence intervals/easier to reject nulls) and
> smaller p-values (again, easier to reject null for a given alpha).
>
> While effect (1) can be important, the expected effect of these changes
> should be zero and can have positive or negative effects on point estimates and
> whether we reject certain null hypotheses. Effects (2) and (3) have a clear
> direction in which they will affect our estimates from a larger vs. a smaller
> sample. The gain from effect (3) is shrinking with the size of the sample, as
> the t-distribution approaches a normal distribution at large degrees of freedom.
> Effect (2) will tend to be lower var_hat(b2) since sum (xi-xbar)^2
> always increases with more observations, but with diminishing effects for a given
> number of new observations (i.e. the effect is larger going from 100 to 200 than
> it is going from 1000 to 1100)
> */

```



```

467
468 // Compare the t critical values for this section to those from earlier
469 scalar tc_312g_1side = invttail(e(df_r), alpha_312) // = (-1)*invt(e(df_r),0.05)

470 scalar tc_312g_2side = invttail(e(df_r), alpha_312/2) // = (-1)*invt(e(df_r),0.05/2)
    > = invt(e(df_r), 1 - 0.05/2)

471
472 scalar zc_1side = (-1)*invnormal(alpha_312)

473 scalar zc_2side = (-1)*invnormal(alpha_312/2)

474 disp "CPS_Small 1-sided critical value: " tc_312b_1side "    CPS 1-sided critical val
    > ue: " tc_312g_1side "Normal 1-sided critical value: " zc_1side
    CPS_Small 1-sided critical value: -1.6463819    CPS 1-sided critical value: 1.6451688No
    > mal 1-sided critical value: 1.6448536

475 disp "CPS_Small 2-sided critical value: " tc_312d_2side "    CPS 2-sided critical val
    > ue: " tc_312g_2side "Normal 2-sided critical value: " zc_2side
    CPS_Small 2-sided critical value: 1.9623438    CPS 2-sided critical value: 1.9604546Nor
    > mal 2-sided critical value: 1.959964

476
477 /* Discussion:
    > As suggested by the forces discussed above, the biggest change is to the standard
    > error of b2, which halved between the two cases. While there is some benefit
    > in the confidence interval from the smaller t critical values, most of the shrinking
    > comes from the smaller standard error.
    > */
478
479 //Convert log file (smcl) to pdf

```