Answers

## Assignment - 2: BUAN 6312 Harikrishna Dev HXD220000

## **Answers**

- 1. Use the data in APPLE to answer this question.
- Define a binary variable as ecobuy = 1 if ecolbs > 0 and ecobuy = 0 if ecolbs = 0. In other words, ecobuy indicates whether, at the prices given, a family would buy any ecologically friendly apples. What fraction of families claim they would buy ecolabeled apples?

The fraction of families claim they would buy ecolabeled apples are 62.42%

- . gen ecobuy = 0
- . replace ecobuy = 1 if ecolbs > 0
  (412 real changes made)
- . tabulate ecobuy

ecobuy	Freq.	Percent	Cum.
0	248 412	37.58 62.42	37.58 100.00
Total	660	100.00	

• Estimate the linear probability model below and and report the results in the usual form. Carefully interpret the coefficients on the price variables (*ecoprc* and *regprc*).

$$ecobuy = \beta_0 + \beta_1 ecoprc + \beta_2 regprc + \beta_3 faminc + \beta_4 hhsize + \beta_5 educ + \beta_6 age + u$$

We get the LRM equation as follows:

$$ecobuy = 0.4236865 + -0.8026219 \times ecoprc + 0.7192675 \times regprc + 0.0005518 \times faminc + 0.0238227 \times hhsize + 0.023827 \times hhsize + 0.02$$

From the following equation, we can see that coefficients of *ecoprc* and *regprc* are 0.803 and 0.719. The p-values of these coefficients are less than 0.05, therefore they are statistically significant. We can also conclude that

$$\frac{\Delta ecoprc}{\Delta ecobuy} = -0.8026219$$

i.e. One dollar increase in price of ecolabeled apples results in a decrease in probablity of a purchase of ecobuy apples by 0.80

$$\frac{\Delta regprc}{\Delta ecobuy} = 0.7192675$$

i.e. One dollar increase in price of regular apples results in a increase in probablity of a purchase of ecobuy apples by 0.71

. reg ecobuy ec	coprc regprc f	aminc hhs	ize educ age			
Source	SS	df	MS	Number of obs	=	660
 Model	17.0019785		2.83366308	F(6, 653) Prob > F		13.43 0.0000
Residual	137.810143	653	.211041566	R-squared Adj R-squared	= =	0.1098 0.1016

interval]	[95% conf.	P> t	t	Std. err.	Coefficient	obuy
5877963	-1.017447	0.000	-7.34	.1094037	8026219	ecoprc
.9777543	4607808	0.000	5.46	.131639	.7192675	egprc
.0015916	000488	0.298	1.04	.0005295	.0005518	faminc
.0484193	0007739	0.058	1.90	.0125262	.0238227	hhsize
.0412287	.008341	0.003	2.96	.0083743	.0247849	educ
.0019536	0029551	0.689	-0.40	.0012499	0005008	age
.747617	.099756	0.010	2.57	.1649674	.4236865	cons

• Are the nonprice variables jointly significant in the LPM? (Use the usual F statistic, even though it is not valid when there is heteroskedasticity.) Which explanatory variable other than the price variables seems to have the most significant effect on the decision to buy ecolabeled apples? Does this make sense to you?

We can see that we conduct a hypothesis tests on the non price variables gives us a  $p\_value < 0.05$ . Therefore, we can reject the null hypothesis i.e. non-price variables are jointly significant. As t(educ) = 2.96 is the highest t statistic value among the non price variable, we can conclude that education makes most significant effect on purchase of eco-labeled apples. This makes sense that educated customers would prefer ecolabeled apples as they would be more well equipped in understanding the benefit of the consumption of them.

• In the model from part (ii), replace faminc with log(faminc). Given the  $R^2$ , which model fits the data better? How many estimated probabilities are negative? How many are bigger than one? Should you be concerned? [Hint: Use command predict y to generate fitted values.]

```
. gen lfaminc = ln(faminc)
. reg ecobuy ecoprc regprc lfaminc hhsize educ age
                                 F(6, 653) = 660

6 2.8797815 Prob > F = 0.0000

653 .210617813 R-squared = 0.1116
                     SS
      Source |
      Model | 17.278689
   Residual | 137.533432
                                                     Adj R-squared =
       Total | 154.812121
                                  659 .234919759 Root MSE
                                                                            .45893
     ecobuy | Coefficient Std. err. t P>|t| [95% conf. interval]
      ecoprc | -.8006664 .1092981 -7.33 0.000 -1.015285 -.5860482
                 .721377 .1315196 5.48 0.000
                                                           .4631247
                                                                        .9796294
      reapro I

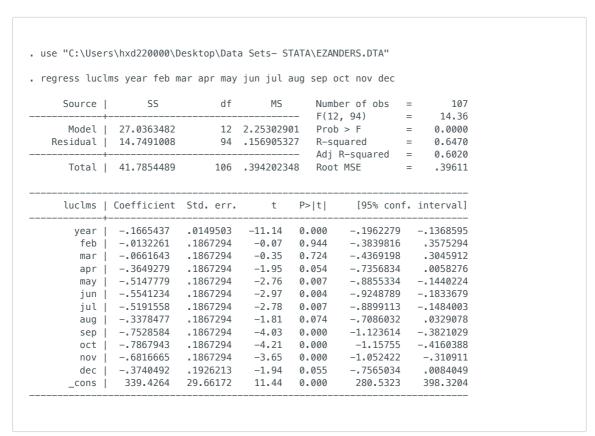
    lfaminc
    .0445162
    .0287239
    1.55
    0.122
    -.0118861
    .1009185

    hhsize
    .0227002
    .012543
    1.81
    0.071
    -.0019294
    .0473297

                                         1.81 0.071 -.0019294
2.73 0.006 .006400
      hhsize | .0227002
                 .023093 .0084508
                                                            .006499
       educ I
                                                                          .039687
        age | -.0003865 .0012517
                                         -0.31 0.758
                                                           -.0028444
                                                                          .0020713
       _cons | .3037519 .1789605
                                        1.70 0.090
                                                           -.0476555
                                                                         .6551593
```

We see that the *Adj-R sqr* of the second model is greater in the first model. This indicates that the second model fits better. In the second model, there are two fitted probabilities are above 1 and in the range of 0.185 to 1.051. The two values aren't of concern as the source has 660 observations and the values are very close to 1. There are no negative probabilities.

- 2. Use the data in EZANDERS for this exercise. The data are on monthly unemployment claims in Anderson Township in Indiana, from January 1980 through November 1988. In 1984, an enterprise zone (EZ) was located in Anderson (as well as other cities in Indiana).
- Regress log(uclms) on a monthly linear time trend and 11 monthly dummy variables. [Hint: Use jan as the
  base month for the monthly dummy variables.] What was the overall trend in unemployment claims over
  this period? (Interpret the coefficient on the time trend.) Is there evidence of seasonality in unemployment
  claims?



We see that coefficient of **YEAR** is -0.1665. This implies that the overall trend of unemployment claims decreases by 16.65% per year. As the p-value < threshold value, we can conclude that the yearly trend is significant.

We can see that some of the monthly dummy variables are significant at a 5% level of significance, whereas some are not significant at the same threshold. This helps us understand that there is a presence of seasonal factors behind unemplyment claims.

To confirm the joint significance, we perform the Wald test on the 11 monthly dummy variables.

$$H_0: feb - dec = 0$$
  
 $H_1: feb - dec \equiv 0$ 

```
. test feb mar apr may jun jul aug sep oct nov dec
(1) feb = 0
(2) mar = 0
(3)
      apr = 0
(4)
     may = 0
(5) jun = 0
(6) jul = 0
(7)
     aug = 0
(8) sep = 0
(9) oct = 0
(10) nov = 0
(11) dec = 0
      F( 11,
               94) =
                        4.32
           Prob > F =
                        0.0000
```

As the p-value < threshold, we can reject the null hypothesis. Therefore, we can conclude that the monthly dummy variables are jointly significant.

• Add ez, a dummy variable equal to one in the months Anderson had an EZ, to the regression in part (i). Does having the enterprise zone seem to decrease unemployment claims? By how much?

Source	l SS	df	MS	Numh	per of obs	= 107	
	55 +				3, 93)		
Model	28.7422487	13	2.2109422		) > F :		
Residual					uared :		
	+				R-squared :		
Total	41.7854489	106		_		3745	
luclms	Coefficient	Std. err.	t	P> t	[95% conf	. interval]	
year	0811489				1372918		
feb	0132261	.1765405	-0.07	0.940	3638005	.3373484	
mar	0661643	.1765405	-0.37	0.709	4167388	.2844101	
apr	3649279	.1765405	-2.07	0.042	7155023	0143534	
may	5147779	.1765405	-2.92	0.004	8653523	1642034	
jun	5541234	.1765405	-3.14	0.002	9046978	203549	
jul	5191558	.1765405	-2.94	0.004	8697303		
aug	3378477	.1765405	-1.91	0.059	6884222	.0127267	
sep	7528584	.1765405	-4.26	0.000	-1.103433	4022839	
oct	7867943	.1765405	-4.46	0.000	-1.137369	4362198	
nov	6816665	.1765405	-3.86	0.000	-1.032241	3310921	
dec			-1.97		7213057		
ez					7972917		
_cons	170.2854	56.02201	3.04	0.003	59.03674	281.534	

When ez is added to the regression, its coefficient is about -.508 (se  $\approx .146$ ). EZ decreases the unemplyment claims by:

$$100(1 - e^{-0.508}) = 39.82\%$$

- 3. Use the data in HSEINV for this exercise.
- Find the first order autocorrelation in log(invpc) and log(price) respectively. Which of the two series may have a unit root?

The first order autocorrelation for *log(invpc)* is 0.6391.

The first order autocorrelation for log(price) is 0.9492.

As the correlation coefficient is high, we can assume they both have a unit root.

Based on your findings in part (i), estimate the equation below and report the results in standard form.
 Interpret the coefficient β<sub>-</sub>1 and determine whether it is statistically significant.

$$log(invpc_t) = \beta_0 + \beta_1 \times \Delta log(price_t) + \beta_2 \times t + u_t$$

Source	l SS	df	MS	Numl	ber of obs	=	41	
	+			- F(2	<b>,</b> 38)	=	19.77	
Model	575457228	2	.28772861	4 Prol	b > F	=	0.0000	
Residual	.553167094	38	.01455702	9 R-s	quared	=	0.5099	
	+			- Adj	R-squared	=	0.4841	
Total	1.12862432	40	.02821560	8 Roo	t MSE	=	.12065	
linvpc	Coefficient		t 		[95% co	nf.	interval]	
gprice	3.878646				1.939282	2	5.81801	
t	.008037	.0015952	5.04	0.000	.004807	7	.0112664	
_cons	8532545	.040291	-21.18	0.000	9348193	3	7716897	

We can see that the co-efficient of gprice is statistically significant. This implies that 1% growth of price results in 3.87% increase in per capita in the housing investment above it mean value.

 Now use Δlog(invpc\_t) as the dependent variable. Re-run the equation and report the results in standard form. How do your results of the coefficient βˆ\_1 change from part (ii)? Is the time trend still significant? Why or why not?

Source	l SS	df	MS	Numb	er of obs	=	41	
	+				38)			
	.039000234				) > F			
kesidual	.782237921	38			quared .			
	•			_	R-squared			
Total	.821238155	40	.02053095	4 K001	MSE	=	. 14348	
ginvpc	Coefficient	Std. err.				f.	interval]	
gprice	1.566526						3.872745	
t	.000037	.001897	0.02	0.985	0038032		.0038772	
_cons	.0059315	.0479125	0.12	0.902	0910623		.1029253	

We see that the co-efficient is 1.567 and is not statistically significant. The time trend is not significant at 5% level of significance as the p value is 0.902.

- 4. Recall that in the example of testing Efficient Markets Hypothesis, it may be that the expected value of the return at time t, given past returns, is a quadratic function of  $return_{t-1}$ .
- To check this possibility, use the data in NYSE to estimate

$$return_t = \beta_0 + \beta_1 return_{t-1} + \beta_2 return_{t-1}^2 + u_t$$

report the results in standard form.

Total	3070.42479	688	4.4628267	73 Root	: MSE =	2.109	
return	Coefficient				[95% conf.	_	
return_1	.0485723	.0387224			- <b>.</b> 0274563		
return_1_2	009735	.0070296	-1.38	0.167	023537	.004067	
_cons	.2255486	.087234	2.59	0.010	.0542708	.3968263	

We can see both estimates are not statistically significant at 5%.

• State and test the null hypothesis that E(return\_t | return\_(t-1)) does not depend on returnt-1. [Hint: There are two restrictions to test here.] What do you conclude?

$$H_0: \beta_1 = 0 \qquad \beta_2 = 0$$

We need to satisfy the above null for our hypothesis to be satisfied. As the p value > 0.05, we cannot reject the nnull hypothesis.

• Drop  $return_{t-1}^2$  from the model, but add the interaction term  $return_{t-1} \times return_{t-2}$ . Now test the efficient markets hypothesis. [Hint: stata can create lag (or lead) variables using subscripts conveniently. For example, you can use the command gen return\_2 = return[\_n-2] to create  $return_{t-2}$  fast.]

```
. gen return_2 = return[_n-2]
(3 missing values generated)
. gen return_2_1 = return_1*return_2
(3 missing values generated)
. reg return return_1 return_2_1
                                               Number of obs =
                                                                     688
     Source |
                                               F(2, 685) =
                                                                    1.80
      Model |
              16.0639248
                                2 8.03196242
                                               Prob > F
                                                                  0.1658
                                               R-squared
                              685 4.45747442
   Residual | 3053.36998
                                                                  0.0052
                                               Adj R-squared =
                                                                   0.0023
               3069.4339
                              687
                                    4.4678805 Root MSE
                                                                   2.1113
      Total |
     return | Coefficient Std. err.
                                      t
                                          P>|t|
                                                    [95% conf. interval]
   return_1 | .0687116 .0392472
                                   1.75
                                            0.080
                                                    -.0083476
                                                                .1457709
                         .0100134
               .0113384
                                                                 .030999
 return_2_1 |
                                            0.258
                                                    -.0083222
                                     1.13
      _cons | .1731605
                         .0809626
                                      2.14
                                            0.033
                                                     .0141959
                                                                 .3321251
. test return_1 return_2_1
(1) return_1 = 0
(2) return_2_1 = 0
      F(2, 685) =
                       1.80
          Prob > F =
                       0.1658
```

$$H_0: \beta_1 = 0 \qquad \beta_2 = 0$$

What do you conclude about predicting weekly stock returns based on past stock returns?

As both models look very weak when we look at the R sqr and summary statistics, we cannot predict weekly stock returns from our models.

- 5. Use the data in KIELMC for this exercise.
- The variable dist is the distance from each home to the incinerator site, in feet. Consider the model

$$log(price) = \beta_0 + \delta_0 y_{81} + \beta_1 log(dist) + \delta_1 y_{81} \cdot log(dist) + u.$$

If building the incinerator reduces the value of homes closer to the site, what is the sign of  $\delta$ 1? What does it mean if  $\beta$ 1 > 0?

Assuming all the other variables remain constant, we can conlcude that cost of home is positively correlated to the distance from the incinerator. Therefore,

$$\delta_1 > 0$$

Assuming  $\beta 1 > 0$ , We can assume the distance between the expensive houses and the incinerator is large.

• Estimate the model from part (i) and report the results in the usual form. Interpret the coefficient on  $y_81 \cdot log(dist)$ . What do you conclude?

ea Inrice y	y81 ldist y81l	dict						
eg thire	yor turst yort	uist						
Source	SS	df	MS					
	+			- F(3	, 317)	=	69.22	
	24.3172548				b > F			
Residual	37.1217306				quared			
				_	R-squared			
Total	61.4389853	320	.191996829	9 Roo	t MSE	=	.3422	
lprice	Coefficient	Std. err.	t	P> t	[95% cor	nf.	interval]	
v81	0113101	.8050622	-0.01	0.989	-1.59525	5	1.57263	
,	316689							
y81ldist	.0481862	.0817929	0.59	0.556	1127394	ļ	.2091117	
cons	8.058468	.5084358	15.85	0.000	7.058133	3	9.058803	

From our analysis, we get the following equation:

$$lp\hat{r}ice = 8.06 - 0.0113 \times y81 + 0.317 ldist + 0.0481 \times y81 \times ldist$$
  
 $n = 321, R^2 = 0.3958, AdjR^2 = 0.3901$ 

We see that  $\delta 1 = 0.0481862$ , but the p-value > 0.05. So, it is not statistically significant.

• Add  $age, age^2, rooms, baths, log(intst), log(land), andlog(area)$  to the equation. Now, what do you conclude about the effect of the incinerator on housing values?

Source		df	MS		ber of obs		
+ Model I			4.8353762		0, 310) b > F		
Residual	13.0852234			R-s	quared	=	0.7870
Total			.191996829		R-squared t MSE		0.7802 .20545
lprice	Coefficient	Std. err.	t	P> t	 [95% cor	 nf.	interval]
y81	2254466	.4946914	-0.46	0.649	-1.198824	1	.7479309
ldist	.0009226 .0624668		0.02 1.24		0868674 036464		.0887125

rooms   .0461389 .0173442 2.66 0.008 .0120117 .0807 baths   .1010478 .0278224 3.63 0.000 .0463032 .155	q   .0000357 8.71e-06 4.10 0.000 .0000186 .000052
baths   .1010478 .0278224 3.63 0.000 .0463032 .1553	
	s   .0461389 .0173442 2.66 0.008 .0120117 .080266
lintst  0599757 .0317217 -1.89 0.0601223929 .0024	s   .1010478 .0278224 3.63 0.000 .0463032 .155792
	t  0599757 .0317217 -1.89 0.0601223929 .00244
lland   .0953425 .0247252 3.86 0.000 .046692 .143	d   .0953425 .0247252 3.86 0.000 .046692 .14399
larea   .3507429 .0519485 6.75 0.000 .2485266 .4529	a   .3507429 .0519485 6.75 0.000 .2485266 .452959
_cons   7.673854 .5015718 15.30 0.000 6.686938 8.660	s   7.673854 .5015718 15.30 0.000 6.686938 8.66076

We can see that  $\delta 1 = 0.0624668$  with a p-value = 0.215. As the summary of the regression output conducts a two-tailed test, we can assume for the one tailed test

$$H_0: \delta_1 = 0$$
 
$$H_1: \delta_1 > 0$$
 
$$p-value_{one-tailed} = \frac{p-value_{two-tailed}}{2} = \frac{0.215}{2} = 0.107$$

As the p-value > 0.05, we can conclude that the distance from the incinerator is not affecting the price of the houses.

Why is the coefficient on log(dist) positive and statistically significant in part (ii) but not in part (iii)? What
does this say about the controls used in part (iii)?

We can see that in the first model, the coefficient of dist is statistically significant, where it is insignificant in the second model. This is due to the absense of these additional factor. To ensure they are jointly significant, we can perform the Wald's test.

```
. test age agesq rooms baths lintst lland larea

( 1) age = 0
( 2) agesq = 0
( 3) rooms = 0
( 4) baths = 0
( 5) lintst = 0
( 6) lland = 0
( 7) larea = 0

F( 7, 310) = 81.35
Prob > F = 0.0000
```

As the p-value of the test is lesser than the threshold, we can conclude they are jointly significant.

6. Use the data in PHILLIPS for this exercise. As we mentioned in Lecture 7, instead of the static Phillips curve model, we can estimate an expectations-augmented Phillips curve of the form

$$\Delta inf_t = \beta_0 + \beta_1 unem_t + e_t$$

where  $\Delta inf_t = inf_t - inf_{t-1}$ 

 Estimate this equation by OLS and report the results in the usual form. In estimating this equation by OLS, we assumed that the supply shock, et, was uncorrelated with unemt. If this is false, what can be said about the OLS estimator of β1?

			9369398 .3714212	

We obtain the following equation by running a regression as follows:

$$\Delta inf_t = 2.83 - 0.518 \times unem_t + e_t$$

If  $e_t$  is correlated with  $unem_t$ , then the estimator for  $\beta 1$  would be biased and inconsistent.

• Suppose that et is unpredictable given all past information:  $E(e_t \mid inf_(t-1), unem_(t-1), ...) = 0$ . Explain why this makes  $unem_t - 1$  a good IV candidate for  $unem_t$ .

Assuming e\_t is unpredictable, we can choose unme\_t-1 as it correlated with the endogenous variable unem\_t, but not to e\_t. therefore, it can serve as IV for unem\_t. As it satisfies the E(et/unem\_t-1)=0, we can conclude that unem\_t-1 is not correlated to e\_t. By using unem\_t-1 as an IV for unem\_t in the regression, we can obtain consistent estimates of the causal effect of unem\_t on dinf, even if unem\_t is endogenous.

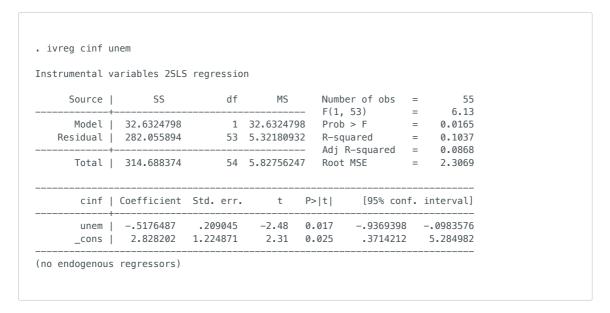
• Does  $unem_t - 1$  satisfy the instrument relevance assumption? [Hint: You need to run a regression to answer this question.]

Source	SS							
Model	+   68.9295284		68 0205284		53) > F			
	52.8515619				ared			
	•			_	-squared			
Total	121.78109	54	2.25520538	Root	MSE	=	.9986	
unem	Coefficient	Std. err.	t	P> t	 [95% con	f. i	nterval]	
unem_1	   <b>.</b> 7423824		8.31				.9214809	
_cons	1.489685	.5202033	2.86	0.006	.446289		2.53308	

As we can see that p-value of the unem\_1 is below the threshold, we can conclude that the unem\_t-1 is strongly correlated with unem\_t and satisfies the assumption.

• Estimate the expectations augmented Phillips curve by 2SLS using  $unem_t - 1$  as an IV for  $unem_t$ . Report the results in the usual form and compare them with the OLS estimates from (i).

## IV Model



## . reg cinf unem

Source	SS	df	MS		Number of obs		55
	<del> </del>			- F(1, 5	3)	=	6.13
Model	32.6324798	1	32.6324798 Prob >		F	=	0.0165
Residual	282.055894	53	5.32180932	80932 R-squared		=	0.1037
	+			- Adi R-	squared	=	0.0868
Total	314.688374	54	5.82756247	7 Root M	SE	=	2.3069
cinf	Coefficient	Std. err.	t	P> t	[95% co	 nf.	interval]
unem	+   - <b>.</b> 5176487	.209045	 -2.48	0.017	.0179369398		0983576
_cons	2.828202	1.224871	2.31	0.025	.371421	2	5.284982