

**Econ 3040 – Final Exam, Dec. 20<sup>th</sup>, 2020**

**Professor: Ryan Godwin**

- The exam consists of 100 marks. Answer all questions.
- You must submit your exam by 3:45 pm on December 20<sup>th</sup> in the UM Learn “Final Exam” Dropbox, under the assignment tab. If you are unable to submit on UM Learn, email me your exam directly.
- You do not need to use R Studio for the final exam, but you can. You should at least have a calculator available.
- The exam is open book / open notes, however, you must provide answers to the questions **in your own words**. You may not plagiarize, or copy-paste, responses to questions from the internet, or from anywhere.
- You must complete the exam individually. Please be aware that there are ways to detect cheating, even electronically.
- A table of standard normal probabilities is provided at the back of the exam.

**Part A – Short Answer – 8 marks each, except Q6**

1) Why are estimators random? Why are the least squares estimators,  $b_0, b_1, \dots, b_k$ , considered to be “good” estimators for  $\beta_0, \beta_1, \dots, \beta_k$ ?

2) Consider the population model:

$$Y = \beta_0 + \beta_1 X + \epsilon$$

Suppose that the sample correlation between  $X$  and  $Y$  is exactly 0. What is the estimated value  $b_1$ ? Explain your answer carefully.

Hints: the formula for calculating  $b_1$  is

$$b_1 = \frac{\sum_{i=1}^n [(X_i - \bar{X})(Y_i - \bar{Y})]}{\sum_{i=1}^n [(X_i - \bar{X})^2]}$$

and the sample covariance between  $X$  and  $Y$  is

$$s_{xy} = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})$$

3) Suppose that your data is: *wage* – the hourly wage of a worker; *male* – a dummy variable; and *educ* – a categorical variable which equals “high” if the worker has obtained a high-school degree, “ugrad” if a undergraduate university degree, “MA” if a Master’s degree, and “PhD” if a PhD. Explain how this categorical variable would be used in an analysis of the effects of education on wage.

4) Why do you reject a null hypothesis when the p-value is small?

5) Explain why it is important to use adjusted-R-square ( $\bar{R}^2$ ) instead of R-square ( $R^2$ ) in a multiple regression model. How does  $\bar{R}^2$  fix the problem that occurs with  $R^2$ ?

6) [10 marks] The following question uses data from the U.S. The two variables are *marriages* – the number of marriages per 1000 people in Kentucky, and *deaths* – the number of people who drowned after falling out of a fishing boat in the U.S. The data was collected annually from the year 1999 to 2010.

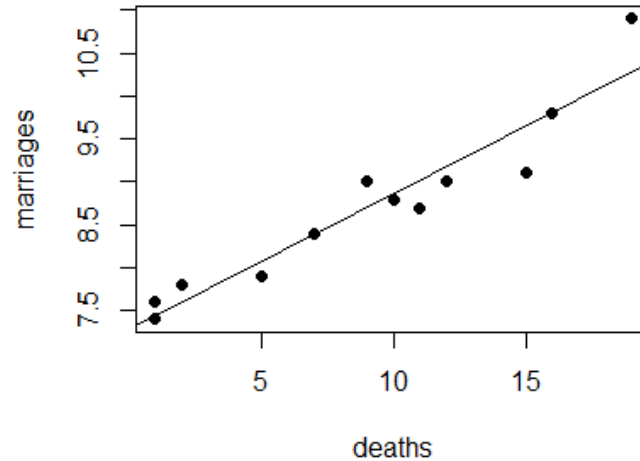
	Year											
	1999	2000	2001	2002	2003	2004	2005	2006	2007	2008	2009	2010
<i>marriages</i>	10.9	9.8	9	9	9.1	8.8	8.7	8.4	7.8	7.9	7.6	7.4
<i>deaths</i>	19	16	9	12	15	10	11	7	2	5	1	1

The model estimated by OLS, is:

$$\widehat{marriages} = 7.28 + 0.16 \times deaths, \quad R^2 = 0.907$$

(0.17) (0.02)

Below is a plot of the data, where the line indicates the LS fitted model:



Test the null hypothesis that *deaths* has **no effect** on marriages. Using the result of this test, the above figure, and the  $R^2$ , comment on whether drowning deaths on fishing boats causes the number of marriages in Kentucky.

### Part B – Long Answer

7) Consider the population model:

$$Wage = \beta_0 + \beta_1 Education + \epsilon$$

where *Wage* and *Education* are the hourly wage, and the number of years of education of a worker, respectively. The model, estimated by least squares, is:

```
summary(lm(Wage ~ Education))
```

Coefficients:

	Estimate	Std. Error	t value	Pr(> t )	
(Intercept)	20.42	3.23	52.13	<2e-16	***
Education	12.63	0.79	?????	??????	***

---

Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

a) [2.5 marks] Fill in the missing t-statistic and p-value.

b) [2.5 marks] Test the hypothesis that *Education* has no effect on *Wage*.

Now, consider the population model:

$$Wage = \beta_0 + \beta_1 Education + \beta_2 IQ + \epsilon$$

where  $IQ$  is the intelligence quotient score of the worker. The model, estimated by least squares, is:

```
summary(lm(Wage ~ Education + IQ))
```

Coefficients:

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	8.23	5.00	1.65	0.0990 .
Education	2.96	3.38	?????	???????
IQ	5.09	0.23	21.68	< 2e-16 ***

---

Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

c) [5 marks] Compared to the first model, do you think the  $R^2$  for this model will be larger or smaller? Explain.

d) [2.5 marks] Calculate the missing t-statistic and p-value.

e) [2.5 marks] Test the hypothesis that *Education* has no effect on *Wage*.

f) [5 marks] Why has the estimated value for  $\beta_1$  changed, compared to the first model? Why do these two estimated models have different implications about the effect of *Education* on *Wage*?

g) [5 marks] Try to tell a story about why the inclusion of *IQ* makes the *Education* variable insignificant in the second model compared to the first (think about the Fireplaces example from chapter 6).

You will need the following table for question 8:

Table 7.1:  $\chi^2$  critical values for the  $F$ -test statistic.

$q$	5% critical value
1	3.84
2	3.00
3	2.60
4	2.37
5	2.21

8) [5 marks each] The following question uses data from a version of the CPS dataset. The dependent variable is *ahe* – the hourly earnings of a worker. *female* is a dummy variable indicating gender; *age* is self explanatory; *yrseduc* is the number of years of education of the worker. In addition, there is a *location* variable which indicates the region in which the worker lives: *northeast*, *south*, *west*, or *midwest* (the “base” group). The sample size is  $n = 61395$ .

Regressor	Model number				
	1	2	3	4	5
<i>female</i>	-4.17* (0.07)	-4.18* (0.07)	-4.19* (0.07)	-4.24* (0.25)	-4.18* (0.07)
<i>age</i>	1.89* (0.65)	1.79* (0.65)	0.15* (0.00)	0.16* (0.00)	0.98* (0.02)
<i>age</i> <sup>2</sup>	-0.03 (0.05)	-0.03 (0.03)			-0.01* (0.00)
<i>age</i> <sup>3</sup>	0.00 (0.00)	0.00 (0.00)			
<i>age</i> <sup>4</sup>	-0.00 (0.00)	-0.00 (0.00)			
<i>yrseduc</i>	0.40* (0.10)	0.38* (0.10)	0.39* (0.10)	1.74* (0.02)	0.39* (0.10)
<i>yrseduc</i> <sup>2</sup>	0.05* (0.00)	0.05* (0.00)	0.05* (0.00)		0.05* (0.00)
<i>northeast</i>	1.21* (0.11)		1.24* (0.11)	1.29* (0.11)	1.21* (0.11)
<i>south</i>	-0.02 (0.10)		-0.02 (0.10)	0.04 (0.10)	-0.01 (0.10)
<i>west</i>	0.76* (0.40)		0.74* (0.10)	0.81* (0.10)	0.76* (0.10)
<i>intercept</i>	-28.30* (6.27)	26.75* (6.27)	-1.22 (0.73)	-10.37* (0.25)	-17.03 (0.86)
$R^2$	0.2676	0.2651	0.2536	0.2514	0.2672
$\bar{R}^2$	0.2675	0.2650	0.2535	?????	0.2671

\*significant at the 1% level.

- How much more do workers in the *northeast* make compared to those in the *west*?
- Using model (1) as the unrestricted model, test the hypothesis that *age* has a linear effect on *ahe*.
- Test the hypothesis that *location* has no effect on *ahe*.
- Using model (5), interpret the effect of *yrseduc* on *ahe*. Be sure to illustrate how this effect depends on the value of *yrseduc* itself.
- Calculate the missing value for adjusted R-square.

END

