

**Economics 102\_A: Analysis of Economic Data**  
**Cameron Winter 2019**  
**Department of Economics, U.C.-Davis**

**Final Exam (A) Thursday March 21**

Compulsory. Closed book. Total of 62 points and worth 45% of course grade.

Read question carefully so you answer the question.

		<b>Question scores</b>																	
Question		1a	1b	1c	1d	1e	2a	2b	2c	2d	2e	2f	3a	3b	3c	3d	3e	3f	
Points		1	3	2	2	2	1	2	3	1	1	1	1	2	1	1	2	1	
Question		4a	4b	4c	4d	4e	4f	5a	5b	6a	6b	6c	6d	7a	7b	7c	7d	<i>Mult Choice</i>	
Points		2	1	2	2	1	1	2	2	1	2	2	2	1	2	1	1	10	

**Questions 1-5**

Consider data on annual outpatient health expenditures for individuals in the U.S.

Outpatient spending is for health services that did not require admission to a hospital.

The data come from the Rand health insurance experiment where different individuals were assigned to one of four different levels of health insurance - see variables `coins0`, `coins25`, `coins50` and `coins95` below.

**Dependent Variable**

`outspend` = Annual outpatient health expenditures (paid through insurance or out-of-pocket)

`lnout` = Natural logarithm of variable `outspend`

**Regressors**

`age` = age in years

`lnage` = Natural logarithm of variable `age`

`education` = years of schooling

`coins0` = 1 if individual's share of health spending is 0% (free) and = 0 otherwise

`coins25` = 1 if individuals share of health spending is 25% and = 0 otherwise

`coins50` = 1 if individuals share of health spending is 50% and = 0 otherwise

`coins95` = 1 if individuals share of health spending is 95% and = 0 otherwise

**Note:** People have exactly one of 0%, 25%, 50% and 95% coinsurance.

**Use the two pages of output provided at the end of this exam on:**

1. Various t critical values.
2. Various descriptive statistics output and correlations for all variables.
3. Three regressions and a test.

**Part of the following questions involves deciding which output to use.**

**You can use the output that gets the correct answer in the quickest possible way.**

1.(a) Give a 95% confidence interval for the population mean outpatient expenditures.

(b) Perform a test at significance level .05 of the claim that the population mean of outpatient expenditures is less than \$1,700.

**State clearly the null and alternative hypotheses of your test, and your conclusion.**

(c) Without using any of the regression output, which two variables (out of `age`, `education`, `coins0`, `coins25`, `coins50` and `coins95`) do you think will best explain outpatient spending?

**Explain your answer.**

(d) Suppose we OLS regress `outspend` on an intercept and `coins0`. Give the intercept and slope coefficient for this regression. **Explain your answer.**

(e) Provide an approximate graph of the distribution of the natural logarithm of outpatient spending. **Note - this is for the natural logarithm.**

**Your graph should include an appropriate scale on the horizontal axis.**

2. In this question the regression studied is a linear regression of **outspend** on **age**.

(a) According to the regression results, by how much does outpatient spending increase in response to ten years of aging?

(b) Give a 99 percent confidence interval for the population slope parameter.

(c) Test the hypothesis at significance level 10% that outpatient spending increases by \$15 with each year of aging. State clearly the null and alternative hypothesis in terms of population parameters and your conclusion.

(d) What is the sample correlation coefficient between age and outpatient spending? **Explain your answer.**

(e) Do you think model errors are likely to be heteroskedastic in this application? **Explain your answer.**

(f) Give the exact Stata command that would provide heteroskedastic-robust standard errors for this model.

**3.** In this question the regression studied is a linear regression of **outspend** on **age**.

(a) Predict the actual outpatient spending for someone aged 50 years.

(b) Give a 95 percent confidence interval for the conditional mean of outpatient spending for someone aged 50 years. Note that  $\sum_{i=1}^n (x_i - \bar{x})^2 = 248,054$  for the variable **age**.

**Give your answer as an expression involving numbers only - you do not need to complete all calculations.**

(c) Using relevant output show that  $\sum_{i=1}^n (x_i - \bar{x})^2 = 248,054$  for **age**.

(d) Suppose we instead want a 95 percent confidence interval for the actual value of outpatient spending for someone aged 50 years. Will the width of this confidence interval be more or less than \$4,000?

**Give a simple answer that does not require complicated calculations.**

(e) What are the four key assumptions on the model for the output of this question to have OLS coefficient estimates that are unbiased and standard error estimates that are unbiased?

[Half point off for each assumption missing].

(f) In addition to these four key assumptions, do you think we need to additionally assume that the errors are normally distributed in order for your answer in part (b) to be exactly correct?

A simple yes or no will do.

4. In this question both regressions where `outspend` is the dependent variable are relevant.

(a) In the second model, based on its coefficient (and not on statistical significance) is education an important determinant of outpatient spending? **Explain your answer.**

(b) Are age, education and the health insurance indicator variables jointly statistically significant at 5 percent? **Explain your answer.**

(c) Are the four levels of health insurance jointly statistically significant at the 5% level? If there is insufficient information to answer this question then say so. **Explain your answer.**

(d) Suppose we give command `generate total = coins0+coins25+coins50+coins95`  
Give the sample mean and standard deviation of variable `total`.

(e) Suppose we give command `regress outspend age education coins0 coins25 coins50`  
Do you expect the coefficient of `age` to be 10.36179? **Explain your answer.**

(f) Using an appropriate measure of goodness-of-fit, which model explains the data better - the second model (with five regressors) or the first model (with one regressor)? **Explain your answer.**

5. In this question consider the regression where  $\ln\text{out}$  is the dependent variable. For each of the following regressors provide a meaningful interpretation of the various slope coefficients **in terms of impacts on outpatient spending** (rather than on  $\ln\text{out}$ ).

(a) Provide a meaningful interpretation of the estimated coefficient for **education**.

(b) Provide a meaningful interpretation of the estimated coefficient for **lnage**.

6.(a) Calculate  $\sum_{i=1}^n z_i$  for  $z_i = 3 + 2i$  and  $n = 4$ .

(b) Suppose for  $X \sim (500, 80^2)$  we form 900 samples of size 400 and obtain 900 sample means  $\bar{x}$ . What approximately do you expect the average of the  $\bar{x}$  to equal? **Explain.**  
What approximately do you expect the standard deviation of the  $\bar{x}$  to equal? **Explain.**

(c) Let  $X$  be the number of students who miss a midterm exam due to illness. Suppose  $X = 1$  with probability 0.4,  $X = 2$  with probability 0.5 and  $X = 6$  with probability 0.1. What is the variance of  $X$ ? **Show all workings.**

(d) Consider a simple random sample of size 3 with values 5, 10, 30. Compute the sample standard deviation. **Show all workings.**

7. You are given the following partial Stata output

```
. regress y x z
```

Source	SS	df	MS			
Model	270			Number of obs =	21	
Residual	90			F( 2, 18) =	(B)	
				Prob > F =		
				R-squared =		
				Adj R-squared =		
				Root MSE =	(A)	

  

	y	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]
	x	6	(C)	1.5		
	z	4	2.0			
	_cons	-8	2.0			

(a) Calculate missing entry (A).

(b) Calculate missing entry (B).

(c) Calculate missing entry (C).

(d) Suppose we perform an F test of  $H_0 : \beta_z = 0$  against  $H_a : \beta_z \neq 0$ . What will the value of the F statistic be?

**Multiple choice questions (1 point each)**

1. Suppose an investment yields a return of four percent a year compounded. Then the investment will have doubled in

- a. between 0 and 10 years
- b. between 10 and 20 years
- c. between 20 and 30 years
- d. more than 30 years.

2. Suppose we want to transform a random variable  $X \sim (\mu, \sigma)$  to a variable  $Y$  with mean 0 and standard deviation one. Then

- a.  $Y = (X - \mu)/\sigma$
- b.  $Y = (X - \mu)/\sigma^2$
- c.  $Y = (X - \mu)^2/\sigma$
- d.  $Y = (X - \mu)^2/\sigma^2$

3. The Stata command `use mydata` is used to

- a. read in data from a text spreadsheet file `mydata.csv`
- b. read in data from an Excel formatted spreadsheet file `mydata.xlsx`
- c. read in data from a Stata dataset `mydata.dta`
- d. none of the above.

4. The Stata command to find the p-value for a two-tailed t-test in a regression with three regressors (including the intercept) when  $t = 1.57$  and there are 40 observations is

- a. `2*invttail(1.57,37)`
- b. `2*invttail(37,1.57)`
- c. `2*ttail(37,1.57)`
- d. `2*ttail(1.57,37)`

5. Suppose we estimate a quadratic model and find  $\hat{y}_i = 3 + 2x_i + x_i^2$ . Then the marginal effect of a change in  $x_i$  equals

- a.  $2x_i + x_i^2$
- b.  $2 + 2x_i$
- c. 2
- d. none of the above.



6. A linear regression of the academic performance index ( $y$ ) for California high schools on various regressors finds that the most important explanators are
- socioeconomic background measures such as fraction of students eligible for free meals
  - teacher quality measures such as whether they have teaching credentials
  - average educational attainment of parents
  - none of the above.
7. Suppose  $\sum_{i=1}^n (x_i - \bar{x})^2 = 4$ ,  $\sum_{i=1}^n (y_i - \bar{y})^2 = 10$  and  $\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) = 20$ . Then the slope coefficient from OLS regression equals
- 0.5
  - 2
  - 2.5
  - 5
  - none of the above.
8. We obtain OLS estimate  $\hat{y} = 2 + 5d$  where  $d$  is an indicator variable taking values 0 or 1. Then
- The mean of  $y = 7$  for those observations with  $d = 0$
  - The mean of  $y = 2$  for those observations with  $d = 1$
  - both a. and b.
  - neither a. nor b.
9. A type 1 error of a statistical test of  $H_0$  against  $H_a$  occurs if we
- reject  $H_0$  given  $H_0$  true
  - reject  $H_0$  given  $H_a$  true
  - do not reject  $H_0$  given  $H_0$  true
  - do not reject  $H_0$  given  $H_a$  true.
10. Multicollinearity is a problem that arises when
- the dependent variable is highly correlated with the regressors
  - the error is highly correlated with the regressors
  - the error is highly correlated with the dependent variable
  - the regressors are highly correlated with each other

SOME USEFUL FORMULAS

Univariate Data

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i \quad \text{and} \quad s_x^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$$

$$\bar{x} \pm t_{\alpha/2; n-1} \times (s_x / \sqrt{n}) \quad \text{and} \quad t = \frac{\bar{x} - \mu_0}{s / \sqrt{n}}$$

ttail(df, t) = Pr[T > t] where T ~ t(df)

t<sub>α/2</sub> such that Pr[|T| > t<sub>α/2</sub>] = α is calculated using invttail(df, α/2).

Bivariate Data

$$r_{xy} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2 \times \sum_{i=1}^n (y_i - \bar{y})^2}} = \frac{s_{xy}}{s_x \times s_y} \quad [\text{Here } s_{xx} = s_x^2 \text{ and } s_{yy} = s_y^2].$$

$$\hat{y} = b_1 + b_2 x_i \quad b_2 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} \quad b_1 = \bar{y} - b_2 \bar{x}$$

TSS =  $\sum_{i=1}^n (y_i - \bar{y})^2$  ResidualSS =  $\sum_{i=1}^n (y_i - \hat{y}_i)^2$  Explained SS = TSS - Residual SS

$$R^2 = 1 - \text{ResidualSS}/\text{TSS}$$

$$b_2 \pm t_{\alpha/2; n-2} \times s_{b_2}$$

$$t = \frac{b_2 - \beta_2^*}{s_{b_2}} \quad s_{b_2}^2 = \frac{s_e^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \quad s_e^2 = \frac{1}{n-2} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

$$y|x = x^* \in b_1 + b_2 x^* \pm t_{\alpha/2; n-2} \times s_e \times \sqrt{\frac{1}{n} + \frac{(x^* - \bar{x})^2}{\sum_i (x_i - \bar{x})^2} + 1}$$

$$E[y|x = x^*] \in b_1 + b_2 x^* \pm t_{\alpha/2; n-2} \times s_e \times \sqrt{\frac{1}{n} + \frac{(x^* - \bar{x})^2}{\sum_i (x_i - \bar{x})^2}}$$

Multiple Regression

$$\hat{y} = b_1 + b_2 x_{2i} + \dots + b_k x_{ki}$$

$$R^2 = 1 - \text{ResidualSS}/\text{TSS} \quad \bar{R}^2 = R^2 - \frac{k-1}{n-k} (1 - R^2)$$

$$b_j \pm t_{\alpha/2; n-k} \times s_{b_j} \quad \text{and} \quad t = \frac{b_j - \beta_j^*}{s_{b_j}}$$

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

$$F = \frac{(\text{ResSS}_r - \text{ResSS}_u) / (k-g)}{\text{ResSS}_u / (n-k)} \sim F(k-g, n-k)$$

Ftail(df1, df2, f) = Pr[F > f] where F is F(df1, df2) distributed.

F<sub>α</sub> such that Pr[F > f<sub>α</sub>] = α is calculated using invFtail(df1, df2, α).



. regress outspend age

Source	SS	df	MS	Number of obs	=	1,950
Model	16055842	1	16055842	F(1, 1948)	=	3.16
Residual	9.9112e+09	1,948	5087864.54	Prob > F	=	0.0758
Total	9.9272e+09	1,949	5093492.03	R-squared	=	0.0016
				Adj R-squared	=	0.0011
				Root MSE	=	2255.6

outspend	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]	
age	8.045322	4.52892	1.78	0.076	-.8367169	16.92736
_cons	1325.305	178.5266	7.42	0.000	975.1813	1675.428

. regress outspend age education coins25 coins50 coins95

Source	SS	df	MS	Number of obs	=	1,950
Model	157414231	5	31482846.2	F(5, 1944)	=	6.26
Residual	9.7698e+09	1,944	5025618.18	Prob > F	=	0.0000
Total	9.9272e+09	1,949	5093492.03	R-squared	=	0.0159
				Adj R-squared	=	0.0133
				Root MSE	=	2241.8

outspend	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]	
age	10.36179	4.606929	2.25	0.025	1.32675	19.39683
education	41.37921	16.97463	2.44	0.015	8.088812	74.6696
coins25	-462.7932	126.2862	-3.66	0.000	-710.4638	-215.1226
coins50	-659.371	191.112	-3.45	0.001	-1034.177	-284.565
coins95	-428.764	134.1436	-3.20	0.001	-691.8443	-165.6837
_cons	990.4725	305.762	3.24	0.001	390.8168	1590.128

. test coins25 coins50 coins95

- ( 1) coins25 = 0
- ( 2) coins50 = 0
- ( 3) coins95 = 0

F( 3, 1944) = 7.79  
 Prob > F = 0.0000

. regress lnout lnage education coins25 coins50 coins95

Source	SS	df	MS	Number of obs	=	1,950
Model	158.317781	5	31.6635562	F(5, 1944)	=	20.59
Residual	2989.28621	1,944	1.53769867	Prob > F	=	0.0000
Total	3147.60399	1,949	1.61498409	R-squared	=	0.0503
				Adj R-squared	=	0.0479
				Root MSE	=	1.24

lnout	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]	
lnage	.4622027	.0961673	4.81	0.000	.2736008	.6508045
education	.0375138	.0093656	4.01	0.000	.0191461	.0558814
coins25	-.3907023	.0698611	-5.59	0.000	-.5277129	-.2536918
coins50	-.6245889	.1057155	-5.91	0.000	-.8319164	-.4172613
coins95	-.4994225	.0741998	-6.73	0.000	-.6449421	-.3539029
_cons	4.815274	.3863439	12.46	0.000	4.057582	5.572966