

**PRACTICE Midterm**

**1 True/False/Uncertain (30 points, 5 questions \* 6 each)**

Indicate whether the following are true, false, or uncertain, with a brief explanation.

1. Suppose both  $X_i$  and  $Y_i$  are dummy variables. The slope from a bivariate regression of  $Y_i$  on  $X_i$  is a difference in conditional probabilities.
2. Consider the residual from a bivariate regression:  $e_i = Y_i - \hat{\alpha} - \hat{\beta}X_i$ , where hats denote OLS estimates. The slope coefficient from a regression of  $e_i$  on  $X_i$  depends on the value of  $\hat{\beta}$ .
3. Regression estimates are made less precise by highly variable regressors.
4. Consider the long regression  $Y_i = \alpha + \beta_1 X_{1i} + \beta_2 X_{2i} + \varepsilon_i$ , and the short regression  $Y_i = \alpha^* + \beta_1^* X_{1i} + \nu_i$ . When  $\beta_2 \neq 0$ , the short coefficient,  $\beta_1^*$ , surely differs from the long coefficient,  $\beta_1$ .
5. Studies of the RAND Health Insurance Experiment show a strong relationship between health insurance coverage and health.

**2 Short Answer (30 points, 10 each)**

1. A survey of 100 male and 400 female MIT grads reveals that 60 of the men and 200 of the women work in a Big Tech firm like Google or Facebook. Show how to use this information to compute the standard error for the male-female difference in Big Tech employment rates by sex. Call your estimate  $se$ ; show how to use  $se$  to construct an approximate 95% confidence interval for the difference.
2. Many MIT students are torn between an Economics (Course 14) major and a Computer Science (Course 6) major.
  - (a) Let  $Y_{1i}$  represent the potential post-graduation earnings of MIT students who major in Course 14 and let  $Y_{0i}$  represent the potential post-graduation earnings of MIT students who major in Course 6. Let  $D_i$  be a dummy that indicates the Course 14 majors in a population limited to Course 14 and Course 6 majors. Assuming  $Y_{1i}$  and  $Y_{0i}$  are defined for everyone regardless of major choice, use this notation to define the average causal effect of doing a Course 14 major on Course 14 majors' earnings.
  - (b) Let  $\beta = E[Y_i|D_i = 1] - E[Y_i|D_i = 0]$ . Show that  $\beta$  is likely to be a biased measure of the causal effect of majoring in 14.
3. Use the fact that residuals have expectation zero and are uncorrelated with regressors to show that the slope in a regression of  $Y_i$  on  $X_i$  is  $\frac{C(Y_i, X_i)}{V(X_i)}$ . Suggestion: Work out your answer on a separate piece of paper. Try to keep it brief. Then, writing in Gradescope, use either latex for mathematical symbols or the symbols  $a$  for intercept,  $b$  for slope, and  $Ybar$  and  $Xbar$  for means of  $Y_i$  and  $X_i$ .

### 3 Longer Questions (40 points, 20 each)

1. This question refers to analysis of data on married women from the 1980 Census. The women in this sample all have at least two children; some have 3 or more. The descriptive table at the top of the log below describes the following variables: weeks worked last year (*weeksm1*), number of children (*kidcount*), a dummy indicating women with 3 or more children (*morekids*), years of schooling (*educm*), a dummy for women with a college degree or higher (*colgrad\_c*). The rest of the log reports three bivariate regression results.

1 .	sum weeksml kidcount morekids educm colgrad					
	Variable	Obs	Mean	Std. Dev.	Min	Max
	weeksm1	254,654	19.01833	21.86728	0	52
	kidcount	254,654	2.507799	.7693323	2	12
	morekids	254,654	.3805634	.4855263	0	1
	educm	254,654	12.38582	2.450936	0	20
	colgrad_c	254,654	.1316806	.3381439	0	1
2 .						
3 .	regress weeksml morekids					
	Source	SS	df	MS	Number of obs = 254,654	
	Model	1742078.14	1	1742078.14	F(1, 254652) = 3696.02	
	Residual	120027337	254,652	471.338679	Prob > F = 0.0000	
	Total	121769415	254,653	478.177816	R-squared = 0.0143	
					Adj R-squared = 0.0143	
					Root MSE = 21.71	
4 .	regress weeksml colgrad_c					
	Source	SS	df	MS	Number of obs = 254,654	
	Model	32787.929	1	32787.929	F(1, 254652) = 68.59	
	Residual	121736627	254,652	478.050938	Prob > F = 0.0000	
	Total	121769415	254,653	478.177816	R-squared = 0.0003	
					Adj R-squared = 0.0003	
					Root MSE = 21.864	
	weeksm1	Coef.	Std. Err.	t	P> t  [95% Conf. Interval]	
	morekids	-5.386996	.0886093	-60.79	0.000 -5.560667 -5.213324	
	_cons	21.06843	.0546629	385.42	0.000 20.96129 21.17557	
5 .	regress weeksml educm					
	Source	SS	df	MS	Number of obs = 254,654	
	Model	525110.698	1	525110.698	F(1, 254652) = 1102.90	
	Residual	121244305	254,652	476.117622	Prob > F = 0.0000	
	Total	121769415	254,653	478.177816	R-squared = 0.0043	
					Adj R-squared = 0.0043	
					Root MSE = 21.82	
	weeksm1	Coef.	Std. Err.	t	P> t  [95% Conf. Interval]	
	educm	.5858939	.0176421	33.21	0.000 .5513158 .620472	
	_cons	11.76156	.2227491	52.80	0.000 11.32498 12.19814	

- (a) What's the probability a woman with 2 or more children has more than 2?

- (b) Interpret (one sentence for each) the slope coefficient magnitudes in these regression results (ignore intercepts).
- (c) Are any of these regressions likely to capture a causal relationship?
- (d) Regressions 3 and 4 below shows results from a multivariate regression of weeks worked on *educm* and *morekids*, followed by a bivariate regression of *morekids* on *educm*. Use the omitted variables bias formula to explain the change in coefficients from the regression of *weeksm1* on *educm* with and without *morekids* control.

2 .						
3 . regress <i>weeksm1</i> <i>educm</i> <i>morekids</i>						
Source	SS	df	MS	Number of obs	=	<b>254,654</b>
Model	<b>2045544.77</b>	2	<b>1022772.39</b>	F(2, 254651)	=	<b>2175.42</b>
Residual	<b>119723871</b>	<b>254,651</b>	<b>470.148834</b>	Prob > F	=	<b>0.0000</b>
Total	<b>121769415</b>	<b>254,653</b>	<b>478.177816</b>	R-squared	=	<b>0.0168</b>
				Adj R-squared	=	<b>0.0168</b>
				Root MSE	=	<b>21.683</b>
<i>weeksm1</i>	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]	
educm	<b>.449545</b>	<b>.0176944</b>	<b>25.41</b>	<b>0.000</b>	<b>.4148645</b>	<b>.4842255</b>
morekids	<b>-5.079497</b>	<b>.0893212</b>	<b>-56.87</b>	<b>0.000</b>	<b>-5.254565</b>	<b>-4.90443</b>
_cons	<b>15.38342</b>	<b>.230329</b>	<b>66.79</b>	<b>0.000</b>	<b>14.93198</b>	<b>15.83486</b>

4 .						
regress <i>morekids</i> <i>educm</i>						
Source	SS	df	MS	Number of obs	=	<b>254,654</b>
Model	<b>1102.23662</b>	1	<b>1102.23662</b>	F(1, 254652)	=	<b>4763.17</b>
Residual	<b>58928.6001</b>	<b>254,652</b>	<b>.231408354</b>	Prob > F	=	<b>0.0000</b>
Total	<b>60030.8368</b>	<b>254,653</b>	<b>.235735832</b>	R-squared	=	<b>0.0184</b>
				Adj R-squared	=	<b>0.0184</b>
				Root MSE	=	<b>.48105</b>
<i>morekids</i>	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]	
educm	<b>-.026843</b>	<b>.0003889</b>	<b>-69.02</b>	<b>0.000</b>	<b>-.0276053</b>	<b>-.0260807</b>
_cons	<b>.7130359</b>	<b>.0049108</b>	<b>145.20</b>	<b>0.000</b>	<b>.7034109</b>	<b>.7226608</b>

5 .						
regress <i>weeksm1</i> <i>educm</i> <i>agem1</i> <i>kidcount</i>						
Source	SS	df	MS	Number of obs	=	<b>254,654</b>
Model	<b>3804978.76</b>	3	<b>1268326.25</b>	F(3, 254650)	=	<b>2737.94</b>
Residual	<b>117964437</b>	<b>254,650</b>	<b>463.241455</b>	Prob > F	=	<b>0.0000</b>
Total	<b>121769415</b>	<b>254,653</b>	<b>478.177816</b>	R-squared	=	<b>0.0312</b>
				Adj R-squared	=	<b>0.0312</b>
				Root MSE	=	<b>21.523</b>
<i>weeksm1</i>	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]	
educm	<b>.1453183</b>	<b>.0181865</b>	<b>7.99</b>	<b>0.000</b>	<b>.1096732</b>	<b>.1809634</b>
agem1	<b>.8018506</b>	<b>.0130619</b>	<b>61.39</b>	<b>0.000</b>	<b>.7762497</b>	<b>.8274516</b>
kidcount	<b>-3.81788</b>	<b>.0570737</b>	<b>-66.89</b>	<b>0.000</b>	<b>-3.929742</b>	<b>-3.706017</b>
_cons	<b>2.422062</b>	<b>.4274962</b>	<b>5.67</b>	<b>0.000</b>	<b>1.584181</b>	<b>3.259943</b>

- (e) Regression 5 above reports results from a longer regression on *educm* that includes controls for age of mother (*agem1*) and number of kids (*kidcount*). Briefly explain how to compute the coefficient on *educm* in this model using a two-step regression procedure in which the second step is a bivariate regression.

2. School vouchers are government-funded coupons that cover the cost of schooling. Legendary economist Milton Friedman proposed replacing traditional government-run public schools with school vouchers that families can use to attend private schools. Since Friedman's original proposal, a number of US states and some Latin American countries have experimented with vouchers. Angrist, Bettinger, and Kremer (2002) and Angrist, et al. (2006) study voucher effects on educational attainment using data from a voucher lottery in Colombia. In Colombia, voucher recipients were chosen randomly from a pool of applicants (much as with health insurance in the OregonHealth Plan we read about in class).
- (a) Consider Table 1 from the 2006 paper (reproduced below). This table reports means and differences in applicant characteristics between voucher winners and losers. Variables examined include a dummy for whether the student has a valid national ID number, age, and dummies for boys and whether the applicant has a phone.

TABLE 1—CHARACTERISTICS OF ICFES MATCHING SAMPLE BY VOUCHER STATUS

	Means		Difference by voucher status (winners vs. losers)			
	Full sample (1)	Sample w/valid age (2)	Full sample (3)	Sample w/valid age (4)	Valid ID and age (5)	Valid ID and age and has phone (6)
Won voucher	0.588	0.585				
Valid ID	0.876	0.967	-0.010 (0.010)	0.001 (0.006)	—	—
Age at time of application	12.7 (1.8)	12.7 (1.3)	-0.137 (0.064)	-0.086 (0.045)	-0.085 (0.044)	-0.091 (0.047)
Male	0.487	0.493	0.004 (0.016)	0.011 (0.017)	0.012 (0.017)	0.008 (0.018)
Phone	0.882	0.886	0.013 (0.010)	0.008 (0.011)	0.008 (0.011)	—
N	4,044	3,661	4,044	3,661	3,542	3,139

*Notes:* Robust standard errors are reported in parentheses in columns 3–6. Regression estimates of differences by voucher status in column 4 are for the sample with valid age data embedded in the national ID number. A valid age must be between 9 and 25. Column 5 reports results for a sample limited to those with a valid ID check digit and column 6 shows results for a sample further limited to those with a phone. The sample includes applicants from the 1995 lottery cohort.

- i. What proportion of the full sample has a phone? Does phone ownership differ significantly between winners and losers?
- ii. Looking at all variables and samples described in the table, how plausible is the claim that vouchers were indeed awarded randomly?

- (b) Next, consider Table 2 from the same study (reproduced below). This table reports treatment effects on the likelihood of something called an “ICFES match.” In Colombia, students who graduate high school and hope to go to college take a matriculation exam called the ICFES. The 2006 study therefore uses ICFES matching (test taking) as a proxy for high school graduation.

TABLE 2—VOUCHER STATUS AND THE PROBABILITY OF ICFES MATCH

	Exact ID match		ID and city match		ID and 7-letter name match		ID, city, and 7-letter match	
	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)
A. All applicants ( $N = 3542$ )								
Dependent var. mean	0.354		0.339		0.331		0.318	
Voucher winner	0.072 (0.016)	0.059 (0.015)	0.069 (0.016)	0.056 (0.014)	0.072 (0.016)	0.059 (0.014)	0.068 (0.016)	0.056 (0.014)
Male		-0.052 (0.014)		-0.053 (0.014)		-0.043 (0.014)		-0.045 (0.014)
Age		-0.160 (0.005)		-0.156 (0.005)		-0.153 (0.005)		-0.149 (0.005)
B. Female applicants ( $N = 1789$ )								
Dependent var. mean	0.387		0.372		0.361		0.348	
Voucher winner	0.067 (0.023)	0.056 (0.021)	0.069 (0.023)	0.057 (0.021)	0.071 (0.023)	0.060 (0.021)	0.073 (0.023)	0.062 (0.021)
Age		-0.168 (0.006)		-0.164 (0.006)		-0.160 (0.006)		-0.156 (0.006)
C. Male applicants ( $N = 1752$ )								
Dependent var. mean	0.320		0.304		0.302		0.288	
Voucher winner	0.079 (0.022)	0.063 (0.020)	0.071 (0.022)	0.055 (0.020)	0.074 (0.022)	0.059 (0.020)	0.065 (0.022)	0.050 (0.020)
Age		-0.153 (0.007)		-0.148 (0.007)		-0.146 (0.007)		-0.141 (0.006)

*Notes:* Robust standard errors are shown in parentheses. The sample used to construct this table includes all Bogotá applicants with valid ID numbers and valid age data (i.e., ages 9 to 25 at application). The sample is the same as in Table 1, column 5.

- Focusing on the first two columns, what’s the overall ICFES match rate?
- (Continue to focus on the first two columns) Do vouchers appear to boost ICFES match rates? Are the estimated effects on match rates significantly different from zero?
- Are the estimated treatment effects in columns 1 and 2 robust to inclusion of controls? Briefly explain why or why not.

## ADDITIONAL PRACTICE QUESTION

- (c) [Extra credit] Consider Table 3 from the 2002 paper (reproduced below). Each row of this table reports regression estimates of treatment effects on different outcomes, ranging from effects on scholarship use in the first row to effects on the number of years enrolled at the bottom. Columns 2-4 show voucher treatment effects for applicants from Bogota while columns 5-6 show estimates computed using data from Bogota and an additional town. Within samples, different columns report estimates from regressions with different sets of controls.

TABLE 3—EDUCATIONAL OUTCOMES AND VOUCHER STATUS

Dependent variable	Bogotá 1995				Combined sample	
	Loser means (1)	No controls (2)	Basic controls (3)	Basic +19 barrio controls (4)	Basic controls (5)	Basic +19 barrio controls (6)
Using any scholarship in survey year	0.057 (0.232)	0.509 (0.023)	0.504 (0.023)	0.505 (0.023)	0.526 (0.019)	0.521 (0.019)
Ever used a scholarship	0.243 (0.430)	0.672 (0.021)	0.663 (0.022)	0.662 (0.022)	0.636 (0.019)	0.635 (0.019)
Started 6th grade in private	0.877 (0.328)	0.063 (0.017)	0.057 (0.017)	0.058 (0.017)	0.066 (0.016)	0.067 (0.016)
Started 7th grade in private	0.673 (0.470)	0.174 (0.025)	0.168 (0.025)	0.171 (0.024)	0.170 (0.021)	0.173 (0.021)
Currently in private school	0.539 (0.499)	0.160 (0.028)	0.153 (0.027)	0.156 (0.027)	0.152 (0.023)	0.154 (0.023)
Highest grade completed	7.5 (0.960)	0.164 (0.053)	0.130 (0.051)	0.120 (0.051)	0.085 (0.041)	0.078 (0.041)
Currently in school	0.831 (0.375)	0.019 (0.022)	0.007 (0.020)	0.007 (0.020)	-0.002 (0.016)	-0.002 (0.016)
Finished 6th grade	0.943 (0.232)	0.026 (0.012)	0.023 (0.012)	0.021 (0.011)	0.014 (0.011)	0.012 (0.010)
Finished 7th grade (excludes Bogotá 97)	0.847 (0.360)	0.040 (0.020)	0.031 (0.019)	0.029 (0.019)	0.027 (0.018)	0.025 (0.018)
Finished 8th grade (excludes Bogotá 97)	0.632 (0.483)	0.112 (0.027)	0.100 (0.027)	0.094 (0.027)	0.077 (0.024)	0.074 (0.024)
Repetitions of 6th grade	0.194 (0.454)	-0.066 (0.024)	-0.059 (0.024)	-0.059 (0.024)	-0.049 (0.019)	-0.049 (0.019)
Ever repeated after lottery	0.224 (0.417)	-0.060 (0.023)	-0.055 (0.023)	-0.051 (0.023)	-0.055 (0.019)	-0.053 (0.019)
Total repetitions since lottery	0.254 (0.508)	-0.073 (0.028)	-0.067 (0.027)	-0.064 (0.027)	-0.058 (0.022)	-0.057 (0.022)
Years in school since lottery	3.7 (0.951)	0.058 (0.052)	0.034 (0.050)	0.031 (0.050)	0.015 (0.044)	0.012 (0.043)
Sample size	562		1,147		1,577	

*Notes:* The table reports voucher losers' means and the estimated effect of winning a voucher. Numbers in parentheses are standard deviations in the column of means and standard errors in columns of estimated voucher effects. The samples used to estimate 7th- and 8th-grade completion effects exclude Bogotá 1997. The sample size for these outcomes is 1,304 in columns (5) and (6). The regression estimates are from models that include controls for city, year of application, phone access, age, type of survey and instrument, strata of residence, and month of interview.

- i. Give a rationale for reporting results on intermediate outcomes like whether winners were more likely than losers to use a scholarship or to enroll in private school.
- ii. Briefly assess program impact: does it seem like vouchers boosted winners' human capital as Milton Friedman would have hoped?