

```
1 ***** Preliminaries
2 capture log close _all // Closes any log if open //
3
4 cd "C:/Users/AN.4271/Dropbox/HHS 651/Assignments/Assignment 1/" /* Sets the Stata Working Directory (note the forward slashes.
5 Stata will try to work with either forward or backwards slashes, but Windows-style back
6 slashes sometimes interfere with functionality, so forward slashes are preferred. */
7
8 log using "assignment1log", text replace /* Starts a text-type log file called
9 "assignment1log" */
10
11 *****
12 *****          HHS 651: Assignment 1          *****
13 *****          Stata Solutions - Andrew Proctor *****
14
15
16 ***** Data Manipulation
17
18 ***** Import Dataset CSV File
19 import delimited using "prgswepl.csv", clear
20
21 ***** Question 1: Describe Dataset
22 describe, short
23
24 /* Discussion: There are 4,469 observations (individuals) and 1,328
25 variables in the dataset. */
26
27
28 ***** Question 2: Explanatory Variables
29
30 ***** 2a. Gender (gender_r)
31 *** Explore Gender Variable
32 codebook gender_r // View storage format of variable 'gender_r' //
33
34 *** Create a "Female" Indicator Variable
35 gen female = (gender_r == 2) if !missing(gender_r)
36 /* For individuals whose gender is listed in gender_r, assigns a
37 value of 1 for female if gender is equal to 2, 1 if not. Missing
38 values in gender_r would also appear as missing in the female
39 variable.*/
40
41 tabulate female // Displays the freq/percent of each value of "female."
42
43 /*
44 Discussion: The variable "gender_r" represents the gender listed
45 for each variable. When the CSV file was read into Stata, the variable
46 was interpreted as a 'numeric' type variable. 50.41% of observations
47 are male, 49.59% female, and there are no missing observations.
48 */
49
50 *** Note: Another way to create the female indicator variable would be:
51 // gen female_alt = 0 if !missing(gender_r)
```

```

52 // replace female_alt = 1 if ( gender_r == 2 & !missing(gender_r))
53 // tabulate female_alt
54
55 **** 2b. Years of Schooling (yrsqual)
56 *** Explore 'Years of Schooling' Variable
57 codebook yrsqual // View storage format of variable 'j_q04a' //
58 tabulate yrsqual
59 /* Since 'yrsqual' is a string-variable, only the first
60 9 values are shown using the codebook command. Using tabulate, we
61 see some of the observations have a missing value "D" - which means
62 "Don't Know" according to the downloaded codebook. */
63
64 ***** Format- Years of Schooling Variable
65 replace yrsqual = ".d" if yrsqual == "D"
66 /* Since we need to format the variable as a numeric (quantitative)
67 variable, we need to Stata to interpret the missing values
68 correctly. Missing values in Stata are denoted by ".", where
69 letters can follow the "." to indicate what type of missing data we
70 have. So we change "D" to ".d". */
71
72 destring(yrsqual), gen(yearsch)
73 /* Now, we need to Stata to convert the variable to numeric,
74 by parsing the text (string) values as numbers. */
75
76 tabulate yearsch // Check to make sure no more missing values.
77
78 tabulate yearsch, missing /* Note: You can see missing values again in
79 tabulate by using option, ", missing" */
80
81 summarize yearsch // Produces basic descriptive statistics for 'age'
82
83 /*
84 Discussion: The variable "yrsqual" is a derived measure of years
85 of schooling. The variable was stored in Stata as a "string" type of
86 variable (Why? Because some observations take on the non-numeric "D"
87 value). After converting the variable to numeric, we see the mean is
88 12.33, with std. dev. of 2.57, min of 6 and max of 20. There are 2
89 missing observations.
90 */
91
92 **** 2c. Age (age_r)
93 *** Explore Gender Variable
94 codebook age_r // View storage format of variable 'gender_r' //
95
96 rename age_r age /* Rename 'age_r' to 'age' (not necessary,
97 but makes regression more understandable later */
98
99 *** Generate 'Potential Experience' Variable
100 gen potent_exper = max(0, age - 19) /* Generates a 'Potential Experience'
101 variable, equal to age - 19 for
102 individuals who are at least 19,

```

```

103             0 otherwise. */
104
105     summarize potent_exper, detail
106
107     /*
108     Discussion: The variable "age_r" is a derived measure of age (in years)
109     of the individual. The variable is stored in Stata as numeric and there
110     are no missing observations. Using "summarize, detail" we see that the
111     mean is 22.03 years and median (50th percentile) is 23 years.
112     */
113
114     **** 2d. Cognitive Ability (using pvpsl1)
115     *** Explore 'Problem-solving scale score' Variable
116     codebook pvpsl1 // View storage format of variable 'pvpsl1' //
117
118     *** Generate Quantile of Cognitive Ability
119     egen cogn_rank = rank(pvpsl1) if !missing(pvpsl1) /* Rank of individuals'
120                    pvpsl1 if known. */
121
122     egen count_cogn = count(pvpsl1) if !missing(pvpsl1) /* Total number of
123                    nomissing observations
124                    for pvpsl1. */
125
126     *** Percentile Rank
127     gen cogn_samp_pctile = ((cogn_rank -1) / (count_cogn - 1)) * 100
128
129     /*
130     Discussion: The variable "pvpsl1" is a derived measure of an
131     individuals' problem solving ability. The variable is stored as a
132     numeric variable in Stata and there are 506 missing observations.
133     */
134
135     **** Question 3: Dependent Variable (Monthly Earnings Quintile)
136     codebook monthlyincpr // View storage format of variable 'earnhrbonus' //
137
138     *** Explore 'Employment Status'
139     codebook monthlyincpr
140
141     recode monthlyincpr (1 = 5) (2 = 17.5) (3 = 37.5) (4 = 62.5) (5 = 82.5) ///
142                    (6 = 95), gen(income_pctile)
143
144     * Alternate recode
145     // gen income_pctile = .
146     // replace income_pctile = 5 if (monthlyincpr == 1 & !missing(monthlyincpr))
147     // replace income_pctile = 17.5 if (monthlyincpr == 2 & !missing(monthlyincpr))
148     // replace income_pctile = 37.5 if (monthlyincpr == 3 & !missing(monthlyincpr))
149     // replace income_pctile = 62.5 if (monthlyincpr == 4 & !missing(monthlyincpr))
150     // replace income_pctile = 82.5 if (monthlyincpr == 5 & !missing(monthlyincpr))
151     // replace income_pctile = 95 if (monthlyincpr == 6 & !missing(monthlyincpr))
152
153     replace income_pctile = 0 if c_d05 ==2 // Assign value of 0 for unemployed.

```

```

154
155     drop if c_d05 ==3 | c_d05 == 4 // Drop if not in labor market or unknown.
156
157     codebook income_pctile // Check number of missing values of new var.
158
159     /*
160     Discussion:  The number missing observations for "monthlyincpr" is 1,236.
161     The number of missing observations for the revised measure is 122.
162     */
163
164 *** Question 4:  Regression Analysis
165 *** 4a:  Regress Income Rank on Cognitive Ability, Potential Experience, and Female Gender
166 reg income_pctile cogn_samp_pctile potent_exper i.female if ///
167     ((age >= 30) & (age <= 65))
168     /*
169     Note:  A more concise way to write the condition for age in this
170     interval is to use the command inrange as follows (I will use
171     inrange in the remainder of the solution).
172
173     Additionally, an alternative to use any 'if' condition in the
174     regression whatsoever would be the command:
175     "keep if inrange(age, 30,65)" but deleting observations outside this
176     range is both unnecessary and would make things more difficult if you
177     want to do further analysis on the full sample. */
178
179 reg income_pctile cogn_samp_pctile potent_exper i.female if ///
180     inrange(age, 30, 65)
181
182
183     /*
184     Discussion:
185
186     The coefficient on cogn_samp_pctile implies that a one percentile
187     increase in cognitive ability is estimated to shift an individual's percentile
188     of earnings up by .3391429 (that is, .3391429 percentage points if
189     percentile is expressed on a 0-1 scale).
190
191     The coefficient on potent_exper implies that a one year increase in
192     potential experience is estimated to increase ones' percentile
193     of earnings by .4132204 percentage points.
194
195     The coefficient on female suggests that being female is estimated to
196     increase the percentile of income by 12.38118 percentage points,
197     compared to being a male.
198
199     The constant estimate suggests that that the predicted percentile
200     of income for a male (female = 0) with 0 years of potential experience
201     and in the 0th percentile of cognitive ability is the 37th percentile.
202     */
203
204

```

```

205  *** 4b: Add Exper^2 and Age
206  reg income_pctile cogn_samp_pctile c.potent_exper##c.potent_exper ///
207  i.female age if inrange(age, 30, 65)
208
209  /*
210  Discussion:
211
212  Age: The age variable is omitted. If you look at the top of the
213  regression output, it notes that age is omitted because of
214  collinearity (Stata automatically detects perfect collinearity and drops
215  one of the collinear variables. Age here is a linear function of potential
216  experience and the constant, since age = potentexper + 19. This is a violation
217  of the MLR Assumption 3, which is simply "no perfect collinearity."
218
219  Square of Potential Experience: The quadratic of experience is
220  negative and significant. This indicates that the benefit of an
221  additional year of experience is diminishing as the years of
222  experience one already has increases. Omission of a relevant quadratic
223  term like this is a common example of the misspecification of functional
224  form that is a violation of MLR Assumption 4 (zero conditional mean) for
225  estimating the true model.
226
227  R^2: The R^2 in the second model is higher than the first (0.1668 as
228  opposed to 0.1551), indicating adding the square of experience increases
229  the total amount of explained variation in income percentile. R^2 will
230  never decrease with the addition of subsequent variables. To see this,
231  note that  $R^2 = 1 - (\text{Sum of Squared Residuals} / \text{Total Sum of Squares})$ .
232  Everything except the Sum of Squared Residuals are the same across
233  the two models, and since the second model contains all predictors from
234  the first model, the sum of squared residuals will be no greater than in
235  the first model.
236
237  */
238
239  *** 4c: Compare School Years vs Cognitive Ability
240  reg income_pctile cogn_samp_pctile potent_exper i.female if inrange(age, 30, 65)
241  scalar R2model4a = e(r2_a) // Save R^2 as a scalar. (Also in reg output)
242
243  reg income_pctile yearsch potent_exper i.female if inrange(age, 30, 65)
244  scalar R2model4c = e(r2_a) // Save R^2 as a scalar. (Also in reg output)
245
246  display R2model4a - R2model4c /* Displays difference in R^2 output. Note: For
247  the assignment, you could just compare them
248  from the regression output of each model. */
249
250  /*
251  Discussion:
252
253  The two models perform nearly identically, with the
254  regression model from 4(a) explaining .066432% more of the variation in
255  income quintiles.

```

```

256
257     (Not graded) Potential Problems with Either Model:
258     The two models preview common challenges in applied econometrics we will
259     discuss in subsequent lectures. As you can see from the covariance matrix
260     below, Cov(cogn_samp_pctile, yearsch) is not equal to zero, and both appear
261     likely to affect incomes, implying omitted variable bias (i.e. a violation
262     of MLR Assumption 4). One response would be to control for both cognitive
263     ability and schooling. But this brings up an issue from Ch.3: endogeneity.
264     The basic idea is that OLS is biased if you include explanatory variables
265     that are caused by other variables in the model. If cognitive ability
266     increases years of schooling, then years of schooling is endogenous when you
267     both are in the model. Equally, one might imagine that, as individual gains
268     more years of schooling, their cognitive ability increases. If this is
269     true, cognitive ability is also endogenous to schooling (when two variables
270     causally influence each other, this is a particular type of endogeneity called
271     simultaneity).
272
273     */
274
275     correlate cogn_samp_pctile yearsch, covariance
276
277
278 **** Extra Question for three person groups
279
280 **** Question 5(a) Explore Structure of the variable "g_q03h" - which is
281 ** 'Skill use work - Numeracy - How often - Use advanced math or statistics'
282 codebook g_q03h
283
284     /*
285     From looking at 'math use at work' with the codebook command, we
286     see that this variable takes on only 9 unique values, meaning that
287     all values are displayed by Codebook. From this, we can see right
288     away that we have the following 'Missing value' indicators that need
289     to be relabelled: 'D', 'N', 'R', and 'V'.
290     */
291
292 **** Question 5(b) Suitably reformat g_q03h and provide the mean and
293 **** standard deviation using the original vaue scheme.
294
295     *** Recode Missing Values for g_q03h
296     replace g_q03h = ".d" if g_q03h=="D"
297     replace g_q03h = ".n" if g_q03h=="N"
298     replace g_q03h = ".r" if g_q03h=="R"
299     replace g_q03h = ".v" if g_q03h=="V"
300
301     *** Convert g_q03h to a numeric variable by destringing
302     destring g_q03h, replace
303
304     *** Produce summary statistics for g_q03h using original coding of
305     *** use frequencies
306     summarize g_q03h

```

```
307
308      /*
309      The mean (pre-transformation) of this variable is 1.287818 and the
310      standard deviation is 0.7269503.
311      */
312
313 **** Question 5(c) - Recode g_q03h so that the values represents number of
314 *** times each month an individual uses advanced math or statistics at work
315 recode g_q03h (1 = 0) (1 = 0.5) (3 = 2.5) (4 = 12) (5 = 20) ///
316      , gen(mathuseatwork)
317
318      /*
319      This question highlights a common problem in applied work, which is
320      that survey data often uses an ordinal or interval approach to
321      asking retrospective information. You as the researcher must then
322      decide how to make that interpretable numerically and justify it.
323
324      In assigning values here myself, I assume that individuals work
325      4 5-day work weeks per month, for a total of 20 work days. So if an
326      individual reports they use math at work "everyday," (5 in the old
327      schema) that equates to 20 days per month.
328
329      "Never" (1 in original coding) is straightforwardly represented as
330      0 times per month.
331
332      For less than once a month (1), I code this as
333      as the midpoint between 0 and 1, i.e. 0.5 days per month.
334
335      For less than once a week but at least once a month (3), this
336      should be less than four (i.e. at most 3) according to my
337      assumptions about a 4 week work month, but greater than 1. I again
338      use the midpoint of (1,3), that is is 2.5 days per month.
339
340      For at least once a week but not every day (4), this again should be
341      less than 20 but less than 4. So once again taking the midpoint of
342      (4,20), I code this as 12 days per month.
343      */
344
345 **** Summarize recoded math use at work variable
346 summarize mathuseatwork
347
348      /*
349      The mean of the variable after transforming it to be more directly
350      interpretable is .7665993 and the standard deviation is 2.663051.
351      */
352
353 **** Question 5(d) - Regressions relating to a math use at work -> cognitive
354 **** ability -> income pathway.
355
356 *** Question 5(d) (i) Regression of Cognitive Ability on math use at work
357 reg cogn_samp_pctile mathuseatwork if (inrange(age, 30, 65) & (c_d05==1))
```

```

358
359 *** Question 5(d) (ii) Regression of Earnings Pctile on Cognitive Ability
360 reg income_pctile cogn_samp_pctile if (inrange(age, 30, 65) & (c_d05==1))
361
362 *** Question 5(d) (iii) Regression of Earnings Pctile on math use at work
363 reg income_pctile mathuseatwork if (inrange(age, 30, 65) & (c_d05==1))
364
365 /*
366 Discussion:
367
368 Regression 5(d) (i) suggests that for each additional day per month
369 that an individual uses advanced math at work, their percentile of
370 cognitive ability increases by 1.723182, which is statistically
371 significant (p-value < 0.01). It's not immediately required for
372 this question, but you may note that these estimates seem almost
373 implausibly high - as we will discuss further in 5(f).
374
375 Regression 5(d) (ii), like analysis in question 4, suggests that
376 cognitive ability has a positive impact on earning, with a
377 1 percentile increase in positive ability estimated to increase
378 earnings percentile by 0.2753122, which is statistically
379 significant (p-value < 0.01). If both this relationship and the
380 relationship from 5(d) (i) are indeed correct, then math use at
381 work should have a direct effect on earnings percentile via this
382 pathway.
383
384 Regression 5(d) (iii) estimates that cognitive ability does indeed
385 have an effect earnings percentile - in fact even larger than the
386 estimated effect through the cognitive ability - earnings pathway.
387 An increase in math use of work by once a month is estimated to
388 increase earnings percentile by 1.811131, which is statistically
389 significant (p-value < 0.01). Again, these results are implausibly
390 high - raising the spector of reverse cauality / endogeneity and
391 foreshadowing 5(f).
392 */
393
394
395 **** Question 5(e) - Regressions relating to an erroneous math use at work
396 **** -> years of schooling -> income pathway.
397
398 *** Question 5(e) (i) Regression of years of schooling on math use at work
399 reg yearsch mathuseatwork if (inrange(age, 30, 65) & (c_d05==1))
400
401 *** Question 5(e) (i) Regression of income percentile on years of schooling
402 reg income_pctile yearsch if (inrange(age, 30, 65) & (c_d05==1))
403
404 /*
405 Discussion:
406
407 Regression 5(e) (i) estimates that math use at work
408 has a positive, statistically significant effect on years of

```


409 schooling. Regression 5(e)(ii) then suggests that years
 410 of schooling has a positive, statistically significant effect on
 411 earnings percentile.
 412

413 This would point to a second causal pathway
 414 for math use at work to effect earnings, but thinking about
 415 regression 5(e)(i) - it doesn't make any sense under our assumptions.
 416 If schooling strictly predates math use at work, then math use at
 417 work cannot effect schooling. Instead, what we very likely have is
 418 reverse causality - an individual's schooling instead affects their
 419 math use at work. To see that a coefficient will be different from
 420 zero when the true relationship runs in reverse of what is estimated,
 421 consider the expression for Beta in terms of the sample correlation
 422 and standard deviations:

423 - For regression of y on x, the coefficient on x is:

$$424 \quad \beta_x = \text{Corr}(x,y) * (\text{StdDev } x / \text{StdDev } y)$$

425 - And for the regression of x on y, the coefficient on y is:

$$426 \quad \beta_y = \text{Corr}(x,y) * (\text{StdDev}_y / \text{StdDev}_x)$$

427
 428 Since the fraction $(\text{StdDev}_y / \text{StdDev}_x)$ and it's inverse are always
 429 strictly positive, then for nonzero $\text{Corr}(x,y)$, running regression
 430 in the 'wrong' direction (from y to x) will always yield a nonzero
 431 coefficient with the same sign as the effect in the right direction
 432 (from x to y).
 433

434 To demonstrate this argument, we run a regression interchanging
 435 our dependent and independent variables in 5(e)(i).
 436

437 */

438
 439 *** Demonstrating that regression can't tell us the direction of causality
 440 `reg mathuseatwork yearsch`

441
 442
 443 **** Question 5(f) - Inference from 5(d) in light of 5(e)

444 /*

445 Discussion:

446 In 5(e), we see a rather stark case where causality cannot run in
 447 the direction estimated by OLS, where math use at work is estimated
 448 to increase years of schooling that predates work.
 449

450
 451 This same concern is likely to extend to the relationship in 5(d).
 452 Individuals with higher cognitive ability are probably more likely
 453 to work in jobs with greater use of advanced math. In general,
 454 there is likely to be the same issue of simultaneity in the
 455 relationship between math use at work and cognitive ability.
 456

457 Generally, this question highlights the difficulty in finding good
 458 variables where there is no concern about OVB or reverse causality.
 459

```
460      Specifically, extending the logic from 5(d), it seems reasonable to
461      believe that higher paying jobs may often require greater use of
462      mathematics - irrespective of someone's aptitude or qualifications.
463      Hence, rather than higher math use 'causing' higher earnings, higher
464      earnings in these situations would be 'causing' more math use. But
465      since more math use might actually have the effect we originally
466      hypothesized - increasing cognitive ability and thereby leading to
467      greater earnings - it's hard to disentangle these two effects.
468
469      The potentially problematic nature of the relationship between math
470      use at work and cognitive ability highlights another possible
471      challenge to the regression we have specified in 4(a): while
472      cognitive ability is likely to influence earnings, earnings may also
473      be affecting the measurement of cognitive ability through higher
474      math use at better paid jobs.
475
476      Note: Questions 5 is meant to get at the questions of
477      reverse causality and simultaneity more in-depth. The timing of
478      effects problem in 5(e) is meant especially to highlight that
479      causality can't run in the direction specified. But it is also
480      possible to make a critique centered entirely around more typical
481      omitted variable bias (OVB). Students who don't address reverse
482      causality but instead make a clear and well-reasoned analysis to
483      this question using OVB will still earn full credit.
484      */
485      *****
486      log close _all
487
```