

```

1 ##### Assignment 2 Data #####
2 # 1/28/2020
3 # Bahar Zafer
4
5 ### Exercise 1 Basic characteristics ###
6
7 library(readstata13)
8 library(xtable)
9 library(plyr)
10 library(dplyr)
11 library(ggplot2)
12 library(data.table)
13
14 dat = read.dta13("guide_master_file_anonymized_full.dta")
15
16 # finding number of rows and columns
17 > nrow(dat)
18 [1] 42974
19 > ncol(dat)
20 [1] 3332
21
22 # sinking descriptive characteristics of variables into "summary.txt"
23 > sink("summary.txt")
24 > sapply(dat, FUN = summary)
25 > sink()
26
27 # finding an individual identifier
28 > varnames = names(dat)
29 > varnames[grepl("id", varnames)]
30 [1] "anon_studentid" "schoolid"
31 [4] "pg_bl_middle_income" "schl_bl_confident"
32 [7] "SHS_TotalCandidates" "TotalCandidates_1"
33 [10] "TotalCandidates_2" "TotalCandidates_3"
34 [13] "a_when_decide" "a_idea"
35 [16] "b_when_decide" "b_idea"
36 [19] "alladmin_SHS_TotalCandidates" "alladmin_TotalCandidates_1"
37 [22] "alladmin_TotalCandidates_2" "alladmin_TotalCandidates_3"
38 [25] "alladmin_TotalCandidates_4"
39
40 # "anon_student_id" is likely to be the individual identifier
41
42 # checking if there are duplicates
43 > duplicates = which(duplicated(dat$anon_studentid) == 1)
44 > duplicates #returns index values
45 [1] 1001 1188 1397 2780 5127 5329 5333 5364 5965 7355 8258 8558 10226 10266
46 [18] 10493 10500 10759
47 [18] 10774 11516 12606 13169
48
49 # removing duplicates from data
50 dat = dat[-duplicates, ]
51 > which(duplicated(dat$anon_studentid) == 1)
52 integer(0)
53
54 ### Exercise 2 Small Bases ###
55
56 > varnames[1]
57 [1] "anon_studentid" # individual identifier, "anon_studentid" is the first
58 variable.
59
60 # dataset with individual identifier and all variables starting with "stud_base"
61 > dat.surv = dat[, c(1, grepl("stud_base", varnames))]
62
63 # dataset with individual identifier and all variables starting with "alladmin"

```

```

61 > dat.admin = dat[, c(1, grep("alladmin", varnames))]
62
63 # dataset of other variables
64 > dat.rest = dat[, -c(grep("stud_base", varnames), grep("alladmin", varnames))]
65
66 ### Exercise 3 Consistency ###
67
68 > gender_table = table(dat.surv$stud_base_GENDER, dat.admin$alladmin_GENDER)
69 > age_table = table(dat.surv$stud_base_age, dat.admin$alladmin_age)
70
71 # off-diagonal values in gender_table and age_table are inconsistent
72 > inconsistent_age_num = sum(age_table[upper.tri(age_table, diag = FALSE)]) +
73   sum(age_table[lower.tri(age_table, diag = FALSE)])
74 > inconsistent_gender_num = sum(gender_table[upper.tri(gender_table, diag = FALSE)]) +
75   sum(gender_table[lower.tri(gender_table, diag = FALSE)])
76 > inconsistent_gender_num
77 [1] 27091
78 > inconsistent_age_num
79 [1] 12780
80
81 > varnames[grep("bece", varnames)]
82 [1] "stud_base_taken_mock_bece" "stud_base_bece_likely"
83     "stud_base_bece_best"
84 [4] "stud_base_bece_worst" "stud_base_choice_1_bece"
85     "stud_base_choice_2_bece"
86 # only first 6 results are given.
87 # We will use "stud_base_bece_best", "stud_base_bece_likely" and "stud_base_bece_worst"
88 # checking if there is any inconsistency in the answers
89
90 inconsistent_cases = which(dat$stud_base_bece_worst > dat$stud_base_bece_best)
91 over_confidence_cases = which(dat$stud_base_bece_likely > dat$stud_base_bece_best)
92 under_confident_cases = which(dat$stud_base_bece_worst > dat$stud_base_bece_likely)
93
94 ### Exercise 4 Balance ###
95
96 # finding the treatment variables
97 > varnames[grep("treat", varnames)]
98 [1] "treatgroup" "stud_base_treatgroup" "alladmin_treatgroup"
99
100 # checking if gender is balanced across "treatgroups"
101 > test1_dat = table(dat$GENDER, dat$treatgroup)
102 > test1_dat
103
104      1      2      3      4      5      6
105 F      5 593      0 451      0 385
106 M     9 2135      3 1979      2 2184
107 M    15 2135      1 2219      0 2295
108
109 # testing whether female-male distribution is balanced across treatment 2 and 4
110 > test1_1 = prop.test(test1_dat[c(2, 3), c(2, 4)], p=c(0.5, 0.5))
111 > test1_1
112 2-sample test for given proportions with continuity correction
113
114 data:  test1_dat[c(2, 3), c(2, 4)], null probabilities c(0.5, 0.5)
115 X-squared = 7.422, df = 2, p-value = 0.02445
116 alternative hypothesis: two.sided
117 null values:
118   prop 1 prop 2
119 0.5    0.5
120 sample estimates:
121   prop 1   prop 2
122 0.5189596 0.4903537
123
124 # proportions are close to 0.5 but we can reject the null hypothesis at 0.05 percent
125 level.
126 # distribution is not balanced.
127
128 # same test across treatment 2, 4 and 6
129 > test1_2 = prop.test(test1_dat[c(2, 3), c(4, 6)], p=c(0.5, 0.5)) # p-value = 0.00362
130 > test1_3 = prop.test(test1_dat[c(2, 3), c(2, 6)], p=c(0.5, 0.5)) # p-value = 0.04415

```

```

125
126
127 # checking if age is balanced across "treatgroups"
128 > test2_dat = table(dat$stud_base_age, dat$treatgroup)
129
130 # testing if each group of age (from 14 to 18) is balanced across treatment 2, 4 and 6
131
132 > test2_1 = prop.test(test2_dat[14:18, c(2, 6)], p=c(0.5, 0.5, 0.5, 0.5, 0.5))
133 > test2_1
134 5-sample test for given proportions without continuity correction
135
136 data: test2_dat[14:18, c(2, 6)], null probabilities c(0.5, 0.5, 0.5, 0.5, 0.5)
137 X-squared = 11.866, df = 5, p-value = 0.03667
138 alternative hypothesis: two.sided
139 null values:
140   prop 1 prop 2 prop 3 prop 4 prop 5
141   0.5    0.5    0.5    0.5    0.5
142 sample estimates:
143   prop 1    prop 2    prop 3    prop 4    prop 5
144 0.5027408 0.4855282 0.5485294 0.5607477 0.5666667
145
146 > test2_2 = prop.test(test2_dat[14:18, c(2, 4)], p=c(0.5, 0.5, 0.5, 0.5, 0.5)) #
147 p-value = 0.001667
148 > test2_3 = prop.test(test2_dat[14:18, c(4, 6)], p=c(0.5, 0.5, 0.5, 0.5, 0.5)) #
149 p-value = 0.5455
150
151 # checking if SHS region is balanced across "treatgroups"
152 > test3_dat = table(dat$SHSregioncode, dat$treatgroup)
153
154 # Balance tests across treatment 2, 4, and 6
155 > test3_1 = prop.test(test3_dat[,c(2, 4)])
156
157 > test3_1
158 12-sample test for equality of proportions without continuity correction
159
160 data: test3_dat[, c(2, 4)]
161 X-squared = 87.801, df = 11, p-value = 4.49e-14
162 alternative hypothesis: two.sided
163 sample estimates:
164   prop 1    prop 2    prop 3    prop 4    prop 5    prop 6    prop 7    prop 8    prop
165   9    prop 10
166 0.5680077 0.4842105 0.6190476 0.4767442 0.5581395 0.4761905 0.4943332 0.7151163
167 0.9000000 0.7364341
168 prop 11    prop 12
169 0.4928910 0.6071429
170
171 > test3_2 = prop.test(test3_dat[,c(2, 6)]) # p-value < 2.2e-16
172 > test3_3 = prop.test(test3_dat[,c(4, 6)]) # p-value = 4.178e-08
173
174 ### Exercise 5 Recoding and Histogram ###
175
176 # searching for the variable of the highest education level desired
177 > varnames[grep("educ",varnames)]
178
179 # Recoding the variable, "stud_educ_want", using its levels.
180 > Educ_want = dat$stud_base_educ_want
181 > Educ_want = as.factor(Educ_want)
182 > levels(Educ_want) = c("Junior high school", "Technical or vocational training",
183 "Senior high school", "nursing or teacher training", "Polytechnic", "University")
184
185 # the histogram of desired education by gender
186 > par(mfrow = c(1, 2))
187 > hist(dat$stud_base_educ_want[dat$stud_base_GENDER == "F"], xlab = "Female", main =
188 "Desired Education")
189 > hist(dat$stud_base_educ_want[dat$stud_base_GENDER == "M"], xlab = "Male", main =
190 "Desired Education")
191
192 ### Exercise 6 Manipulating the Data ###

```

```

187 library(reshape2)
188
189 # finding variables related with school/program choice
190 varnames[grep("mychoice",varnames)]
191 # recovered variables of choice:
192 # stud_base_mychoice_1_pgm, stud_base_mychoice_2_pgm, stud_base_mychoice_3_pgm,
stud_base_mychoice_4_pgm
193
194 # creating a dataset of program choices
195 > my_dat1 = cbind(dat$anon_studentid, dat$stud_base_mychoice_1_pgm,
dat$stud_base_mychoice_2_pgm, dat$stud_base_mychoice_3_pgm, dat$stud_base_mychoice_4_pgm)
196 > colnames(my_dat1) = c("id", "choice1", "choice2", "choice3", "choice4")
197 > my_dat1 = as.data.frame(my_dat1)
198 > my_dat1 = melt(my_dat1, id = c("id"))
199
200 # frequency table of programs by ranked choice number
201 > choiceFreq_programs = table(my_dat1$variable, my_dat1$value)
202 > choiceFreq_programs
203
AGRICULTURAL SCIENCE BUSINESS
204 choice1 34465 105 696
205 choice2 35018 125 739
206 choice3 35407 120 714
207 choice4 35834 131 602
208
GENERAL ARTS HOME ECONOMICS OTHER SCIENCE
209 choice1 3018 1565 510 1640
210 choice2 2802 1392 454 1404
211 choice3 2731 1282 458 1252
212 choice4 2614 1239 426 1146
213
TECHNICAL SKILLS VISUAL ARTS
214 choice1 272 682
215 choice2 251 768
216 choice3 248 741
217 choice4 263 698
218
219
220
221 # creating a dataset of region choices
222 > varnames[grep("region",varnames)]
223
224 > my_dat2 = cbind(dat$anon_studentid, dat$stud_ful_mychoice1_region,
dat$stud_ful_mychoice2_region, dat$stud_ful_mychoice3_region,
dat$stud_ful_mychoice4_region)
225 > colnames(my_dat2) = c("id", "regchoice1", "regchoice2", "regchoice3", "regchoice4")
226 > my_dat2 = as.data.frame(my_dat2)
227 > my_dat2 = melt(my_dat2, id = "id")
228
229 # frequency table of regions ranked by choice number
230 > choiceFreq_regions = table(my_dat2$variable, my_dat2$value)
231 > choiceFreq_regions
232
233 Ashanti Brong Ahafo Central Eastern Greater Accra
234 regchoice1 31466 10210 132 244 117 103
235 regchoice2 31544 10316 142 160 116 70
236 regchoice3 31661 10235 144 158 115 56
237 regchoice4 31799 10180 155 128 103 53
238
Northern Upper East Upper West Volta Western
239 regchoice1 339 165 117 15 45
240 regchoice2 288 163 116 16 22
241 regchoice3 261 158 114 23 28
242 regchoice4 259 146 95 14 21
243
244
245 # searching for the variables of expected score needed for admission
246 varnames[grep("score",varnames)]
247 # "alladmin_score1", "alladmin_score2", "alladmin_score3" and "alladmin_score4" are
retrieved.
248
249 my_dat3 = cbind(dat$anon_studentid, dat$alladmin_score1, dat$alladmin_score2,
dat$alladmin_score3, dat$alladmin_score4)

```

```

250 colnames(my_dat3) = c("id", "score1", "score2", "score3", "score4")
251 my_dat3 = as.data.frame(my_dat3)
252 my_dat3 = melt(my_dat3, id = c("id"))
253
254 my_dat4 = cbind(my_dat1, my_dat3)
255 colnames(my_dat4) = c("id", "choice_rank", "program", "id", "score_order",
  "score_values")
256
257 # plotting the average, sd, the 25th and 75th quantiles by ranked choice
258 avg = tapply(my_dat4$score_values, my_dat4$choice_rank, mean, na.rm = TRUE)
259 sd = tapply(my_dat4$score_values, my_dat4$choice_rank, sd, na.rm = TRUE)
260 avg=as.matrix(avg)
261 sd=as.matrix(sd)
262 quant = tapply(my_dat4$score_values, my_dat4$choice_rank, quantile, na.rm = TRUE)
263 quant25th = cbind(c(1,2,3,4), as.matrix(c(quant$`choice1`[2], quant$`choice2`[2],
  quant$`choice3`[2], quant$`choice4`[2])))
264 quant75th = cbind(c(1,2,3,4), as.matrix(c(quant$`choice1`[4], quant$`choice2`[4],
  quant$`choice3`[4], quant$`choice4`[4])))
265
266 par(mfrow = c(2,2))
267 plot(avg, xlab = "choice rank", ylab = "mean score")
268 plot(sd, xlab = "choice rank", ylab = "std. dev. of score")
269 plot(quant25th[,1], quant25th[,2], xlab = "choice rank", ylab = "25th Quantile")
270 plot(quant75th[,1], quant75th[,2], xlab = "choice rank", ylab = "75th Quantile")
271
272 # selecting 5 ids randomly
273 > random5ids = trunc(runif(5, as.numeric(min(levels(my_dat4$id))),
  as.numeric(max(levels(my_dat4$id)))))
274 > random5ids
275 [1] 8707 8897 828 3165 837
276 # finding index values of selected ids in the dataset
277 > id_index_vec = c(which(my_dat4$id == random5ids[1]),
  which(my_dat4$id == random5ids[2]),
  which(my_dat4$id == random5ids[3]),
  which(my_dat4$id == random5ids[4]),
  which(my_dat4$id == random5ids[5]))
281
282 # table reporting how the expected scores change by choice ranks for the selected 5
  students
283 > table_random5 = my_dat4[id_index_vec,]
284 > table_random5
285
286      id choice_rank      program      id score_order score_values
287 8707    8707    choice1 HOME ECONOMICS 8707      score1          37
288 51660   8707    choice2 HOME ECONOMICS 8707      score2          67
289 94613   8707    choice3 HOME ECONOMICS 8707      score3          85
290 137566  8707    choice4             8707      score4          79
291 8897    8897    choice1             8897      score1          34
292 51850   8897    choice2             8897      score2          60
293 94803   8897    choice3             8897      score3          73
294 137756  8897    choice4             8897      score4          56
295 828     828    choice1    GENERAL ARTS 828     score1          48
296 43781   828    choice2             828     score2          68
297 86734   828    choice3             828     score3          79
298 129687  828    choice4             OTHER 828     score4          55
299 3165    3165    choice1             3165     score1          29
300 46118   3165    choice2             3165     score2          38
301 89071   3165    choice3             3165     score3          11
302 132024  3165    choice4             3165     score4          35
303 837     837    choice1             SCIENCE 837     score1          57
304 43790   837    choice2             SCIENCE 837     score2          76
305 86743   837    choice3             SCIENCE 837     score3          90
306 129696  837    choice4             SCIENCE 837     score4          70
307
308 # expected score values of choice1 are the lowest for the students with ids 8707, 8897,
  828 and 837.
309 # expected score values of choice4 are lower than those of choice3 for the students
  with ids 8707, 8897, 828 and 837.
310
311 # Creating a dummy variable "reverting"

```

```

312 # "reverting" = 1 if the score is decreasing in choice rank.
313
314 reg_dat = dat[complete.cases(dat$alladmin_score1), ]
315 reg_dat = reg_dat[complete.cases(reg_dat$alladmin_score2), ]
316 reg_dat = reg_dat[complete.cases(reg_dat$alladmin_score3), ]
317 reg_dat = reg_dat[complete.cases(reg_dat$alladmin_score4), ]
318
319 reverting = vector()
320 for (i in 1:nrow(reg_dat)){
321   if(reg_dat$alladmin_score1[i] > reg_dat$alladmin_score2[i]) {
322     reverting[i] = 1
323   } else {reverting[i] = 0}
324 }
325
326 reg_dat = cbind(reg_dat, reverting)
327
328 ### Exercise 7 Probit ###
329
330 # preparing dataset to be used in probit estimation
331 consistent = as.numeric(reg_dat$stud_base_bece_worst > reg_dat$stud_base_bece_best)
332 over_confident = as.numeric(reg_dat$stud_base_bece_likely > reg_dat$stud_base_bece_best)
333 under_confident = as.numeric(reg_dat$stud_base_bece_worst >
334   reg_dat$stud_base_bece_likely)
335
336 reg_dat = cbind(reg_dat, consistent, over_confident, under_confident)
337
338 # running a probit model of reverting on measure of consistency, overconfidence,
339 # underconfidence
340 # as well as gender, age, and expected education.
341
342 > myprobit = glm(reverting ~ consistent + over_confident + under_confident + GENDER +
343   alladmin_age + stud_base_educ_want,
344   family=binomial(link = "probit"), data = reg_dat)
345
346 > myprobit
347
348 Call:  glm(formula = reverting ~ consistent + over_confident + under_confident +
349   GENDER + alladmin_age + stud_base_educ_want, family = binomial(link =
350   "probit"),
351   data = reg_dat)
352
353 Coefficients:
354 (Intercept)          consistent          over_confident
355 under_confident          GENDERF          GENDERM
356 -2.49127          0.20185          -0.01555          0.07223
357 0.19891          0.01510
358 alladmin_age  stud_base_educ_want
359 0.02447          -0.01386
360
361 Degrees of Freedom: 10067 Total (i.e. Null); 10060 Residual
362 (31688 observations deleted due to missingness)
363 Null Deviance: 3145
364 Residual Deviance: 3124 AIC: 3140
365

```