# STATS 506 HW3

Minxuan Chen

2023-10-01

## Table of contents

Github repo: https://github.com/PKUniiiiice/STATS_506

# 1 Problem 1 Vision

```
1
2  . do "K:\STATS_506\STATA\stata_hw3.Do"
3
4  . //load data
5  . cd "K:\STATS_506\STATA\"
6  K:\STATS_506\STATA
7
8  .
9  . //Part a ----------------------------------------------------------------
10 > ----------------------
11 .
12 . //import first
13 . import sasxport5 VIX_D.XPT, clear
14
15 . //save as dta
16 . save "K:\STATS_506\STATA\VIX_D.dta"
17 file K:\STATS_506\STATA\VIX_D.dta saved
18
19 . //import second
20 . import sasxport5 DEMO_D.XPT, clear
21
```

```
22  . // Merge the second dataset using the SEQN variable
23  . merge 1:1 seqn using "K:\STATS_506\STATA\VIX_D.dta"
24
25     Result                      Number of obs
26     -----------------------------------------
27     Not matched                         3,368
28         from master                     3,368  (_merge==1)
29         from using                          0  (_merge==2)
30
31     Matched                             6,980  (_merge==3)
32     -----------------------------------------
33
34  .
35  . // Keep only the matched records
36  . keep if _merge == 3
37  (3,368 observations deleted)
38
39  .
40  . //total sample size
41  . di _N
42  6980
43
44  .
45  . //End of Part a -----------------------------------------------------------
46  > ----------------------------
47  .
48  . //Part b -----------------------------------------------------------------
49  > ----------------------------
50  . //We use the variable 'VIX220 - Glasses/contact lenses worn for distance' (viq
51  > 220)
52  . //and  'RIDAGEYR - Age at Screening Adjudicated - Recode' (ridageyr)
53  .
54  . egen age_interval = cut(ridageyr), at(0(10)90) label
55
56  . table age_interval viq220, missing statistic(percent, across(viq220)) statisti
57  > c(frequency)
58
59  ----------------------------------------------------------------------
60              |       Glasses/contact lenses worn for distance
61              |       1         2         9         .        Total
62  ------------+---------------------------------------------------------
```

2

```
age_interval |
  10-        |
    Percent  |         30.36       64.25                      5.39      100.00
    Frequency |           670       1,418                       119       2,207
  20-        |
    Percent  |         29.97       61.80       0.20            8.03      100.00
    Frequency |           306         631          2              82       1,021
  30-        |
    Percent  |         32.89       58.80                       8.31      100.00
    Frequency |           269         481                        68         818
  40-        |
    Percent  |         35.09       59.75                       5.15      100.00
    Frequency |           286         487                        42         815
  50-        |
    Percent  |         53.09       43.42                       3.49      100.00
    Frequency |           335         274                        22         631
  60-        |
    Percent  |         59.30       36.01                       4.69      100.00
    Frequency |           392         238                        31         661
  70-        |
    Percent  |         63.75       31.56                       4.69      100.00
    Frequency |           299         148                        22         469
  80-        |
    Percent  |         58.10       28.77                      13.13      100.00
    Frequency |           208         103                        47         358
  Total      |
    Percent  |         39.61       54.15       0.03            6.20      100.00
    Frequency |         2,765       3,780          2             433       6,980
-----------------------------------------------------------------

. quietly collect layout (age_interval) (viq220#result[percent] viq220#result[fr
> equency])

. //If only want to see first column
. //https://grodri.github.io/stata/tables https://www.stata.com/manuals/tables.p
> df
.
. //percent
. collect layout (age_interval) (viq220[1]#result[percent] viq220[.m]#result[fre
> quency])
```

```
Collection: Table
      Rows: age_interval
   Columns: viq220[1]#result[percent] viq220[.m]#result[frequency]
   Table 1: 10 x 2


------------------------------------------------------------
             |   Glasses/contact lenses worn for distance
             |                        1              Total
             |                  Percent          Frequency
-------------+----------------------------------------------
age_interval |
  10-        |                    30.36              2,207
  20-        |                    29.97              1,021
  30-        |                    32.89                818
  40-        |                    35.09                815
  50-        |                    53.09                631
  60-        |                    59.30                661
  70-        |                    63.75                469
  80-        |                    58.10                358
  Total      |                    39.61              6,980
------------------------------------------------------------


.
. //End of Part b ---------------------------------------------------------------
> ----------------------------
.
.
. //Part c ---------------------------------------------------------------------
> ----------------------------
. //For age, we use ridageyr
. //For race, we use ridreth1
. //For gender, we use riagendr
. //For Poverty Income ratio, we use indfmpir
. //We first check how many missing values are in these variables
. misstable summarize ridageyr ridreth1 riagendr indfmpir viq220
                                                      Obs<.
                                          +----------------------------
          |                               | Unique
 Variable |     Obs=.     Obs>.     Obs<. | values       Min        Max
----------+-------------------------------+----------------------------
 indfmpir |       342               6,638 |    435         0          5
```

4

```
145        viq220 |        433                6,547  |      3          1          9
146    ----------------------------------------------------------------------------
147
148    .
149    . //It seems that the proportion of missing value is not large, about 10%, so we
150    >  choose to directly delete them
151    . drop if missing(indfmpir) | missing(viq220)
152    (731 observations deleted)
153
154    . misstable summarize ridageyr ridreth1 riagendr indfmpir viq220
155    (variables nonmissing or string)
156
157    .
158    . ///We treat viq220==1 as "Yes, wear", and all other values as "No, don't wear"
159    > //recode viq220
160    . recode viq220 (1=1) (else=0), generate(viq220_bin)
161    (3,594 differences between viq220 and viq220_bin)
162
163    .
164    . //ref https://www.stata.com/manuals/rlogistic.pdf
165    . //Note that race and gender shoule be categorical variables and age and PIR ar
166    > e continuous variables
167    .
168    .
169    . // Fit the first logistic regression model (age only)
170    . logistic viq220_bin ridageyr
171
172    Logistic regression                              Number of obs =  6,249
173                                                     LR chi2(1)    = 403.63
174                                                     Prob > chi2   = 0.0000
175    Log likelihood = -4058.8462                      Pseudo R2     = 0.0474
176
177    ----------------------------------------------------------------------------
178      viq220_bin | Odds ratio   Std. err.      z    P>|z|     [95% conf. interval]
179    -------------+--------------------------------------------------------------
180        ridageyr |   1.024531    .0012702    19.55   0.000     1.022044    1.027023
181           _cons |    .2923673     .015974   -22.51   0.000     .2626769     .3254136
182    ----------------------------------------------------------------------------
183    Note: _cons estimates baseline odds.
184
185    .
```

5

```
186   . // Store the results
187   . eststo model1
188
189   .
190   . // Fit the second logistic regression model (age, race, gender)
191   . logistic viq220_bin ridageyr i.ridreth1 i.riagendr
192
193   Logistic regression                              Number of obs =  6,249
194                                                    LR chi2(6)    = 584.06
195                                                    Prob > chi2   = 0.0000
196   Log likelihood = -3968.6291                      Pseudo R2     = 0.0685
197
198   ------------------------------------------------------------------------------
199      viq220_bin | Odds ratio   Std. err.      z    P>|z|     [95% conf. interval]
200   -------------+----------------------------------------------------------------
201        ridageyr |     1.0226   .0013241    17.26   0.000     1.020009    1.025199
202                 |
203        ridreth1 |
204               2 |   1.169508   .1959093     0.93   0.350     .8421995    1.624021
205               3 |   1.895064   .1363291     8.89   0.000     1.645846    2.182019
206               4 |   1.293781   .1015763     3.28   0.001     1.109257    1.509002
207               5 |   1.885095   .2612655     4.57   0.000     1.436681    2.473465
208                 |
209      2.riagendr |   1.650228   .0891912     9.27   0.000     1.484357    1.834634
210           _cons |   .1650721   .0132324   -22.47   0.000     .1410718    .1931555
211   ------------------------------------------------------------------------------
212   Note: _cons estimates baseline odds.
213
214   .
215   . eststo model2
216
217   .
218   . // Fit the third logistic regression model (age, race, gender, Poverty Income
219   > ratio)
220   . logistic viq220_bin ridageyr i.ridreth1 i.riagendr indfmpir
221
222   Logistic regression                              Number of obs =  6,249
223                                                    LR chi2(7)    = 625.24
224                                                    Prob > chi2   = 0.0000
225   Log likelihood = -3948.0387                      Pseudo R2     = 0.0734
226
```

```
227    ------------------------------------------------------------------------------
228      viq220_bin |  Odds ratio   Std. err.      z    P>|z|     [95% conf. interval]
229    -------------+----------------------------------------------------------------
230        ridageyr |    1.02245     .001324    17.15   0.000     1.019858    1.025048
231                 |
232        ridreth1 |
233               2 |    1.124663   .1892328     0.70   0.485     .8087261    1.564023
234               3 |    1.652417    .124123     6.69   0.000     1.426201    1.914514
235               4 |     1.23222   .0975979     2.64   0.008      1.05504    1.439155
236               5 |     1.70633   .2391229     3.81   0.000     1.296513    2.245688
237                 |
238      2.riagendr |    1.673821   .0908852     9.49   0.000     1.504841    1.861777
239         indfmpir |     1.12011   .0198248     6.41   0.000      1.08192    1.159647
240           _cons |    .1330474   .0116811   -22.97   0.000     .1120144    .1580298
241    ------------------------------------------------------------------------------
242    Note: _cons estimates baseline odds.
243
244    .
245    . eststo model3
246
247    .
248    . // Create a table to display results using esttab
249    . // https://repec.org/bocode/e/estout/hlp_esttab.html
250    . esttab model1 model2 model3, ///
251    >         con ///
252    >         not ///
253    >     stats(N r2_p aic) ///
254    >     eform ///
255    >     varwidth(15) ///
256    >     title("Logistic Regression Results") ///
257    >     label
258
259    Logistic Regression Results
260    ----------------------------------------------------------------
261                          (1)            (2)            (3)
262                RECODE of ~c   RECODE of ~c   RECODE of ~c
263    ----------------------------------------------------------------
264    RECODE of viq~c
265    Age at Screen~R      1.025***       1.023***       1.022***
266    Race/Ethnicit~1                            1              1
267    Race/Ethnicit~2                        1.170          1.125
268    Race/Ethnicit~3                        1.895***       1.652***
```

7

```
269  Race/Ethnicit~4                                          1.294**          1.232**
270  Race/Ethnicit~5                                          1.885***         1.706***
271  Gender=1                                                 1                1
272  Gender=2                                                 1.650***         1.674***
273  Family PIR                                                                1.120***
274  Constant                       0.292***         0.165***         0.133***
275  -------------------------------------------------------------------
276  N                                 6249             6249             6249
277  r2_p                            0.0474           0.0685           0.0734
278  aic                             8121.7           7951.3           7912.1
279  -------------------------------------------------------------------
280  Exponentiated coefficients
281  * p<0.05, ** p<0.01, *** p<0.001
282
283  .
284  . //End of Part c -----------------------------------------------------------
285  > ----------------------------
286  .
287  . //Part d ------------------------------------------------------------------
288  > ----------------------------
289  .
290  . //Note that in the output table, the odds ratio of Gender=2 is significant, th
291  > erefore,
292  . //the odds of men and women being wears of glasses/contact lenses for distance
293  >  vision differs.
294  .
295  . //We use chi-square test (Pearson's Chi-Squared Test of Independence)
296  . tabulate riagendr viq220_bin, chi2
297
298              |   RECODE of viq220
299              |   (Glasses/contact
300              |     lenses worn for
301              |         distance)
302     Gender  |         0            1 |      Total
303  -----------+--------------------+----------
304          1 |     1,919        1,134 |      3,053
305          2 |     1,675        1,521 |      3,196
306  -----------+--------------------+----------
307      Total |     3,594        2,655 |      6,249
308
309          Pearson chi2(1) =   69.7397   Pr = 0.000
```

```
310
311   .
312   . //From the result, p-value is 0.000, therefore we conclude that gender and wea
313   > ring or not
314   . //are not independent, in other words, the proportion of wearers of glasses/co
315   > ntact lenses for distance vision differs between men and women
316   .
317   . //End of Part d -----------------------------------------------------------
318   > ----------------------------
319   .
320   .
321   .
322   .
323   .
324   .
325   end of do-file
326
327   .
328
329
330
331
```