

```
Untitled1* x HW1.Rmd x main0.R x main.R* x
Insert Run

1 ---
2 title: "Econ 613 hw 1"
3 author: "Litong Zhang"
4 date: "Feb 28, 2021"
5 output:
6   pdf_document:
7     toc: yes
8   html_document:
9     toc: yes
10 ---
11 ## Ex 1
12
13 ```{r 1, tidy=TRUE}
14 library(dplyr)
15 library(tidyr)
16 library(doBy)
17 library(sqldf)
18 |
19 setwd("~/Desktop/613")
20 stu = read.csv('datstu.csv')
21 jss = read.csv('datjss.csv')
22 sss = read.csv('datsss.csv')
23 nrow(stu)

22 sss = read.csv('datsss.csv')
23 nrow(stu)
24 # There are 340,823 students.
25
26 jss = unique(subset(jss, select = -c(X)))
27 nrow(jss)
28 # There are 139 junior high schools.
29
30 sss = unique(subset(sss, select = c(schoolcode, sssdistrict, ssslong, ssslats)))
31 sss = sqldf('SELECT * FROM sss GROUP BY schoolcode')
32 nrow(sss)
33 # Total of 898 senior high schools.
34
35 prog = as.matrix(stu[,11:16])
36 prog = matrix(prog, ncol=1)
37 prog = unique(prog)
38 prog = prog[prog!='']
39 length(prog) # Total of 32 programs.
40
```

```
Untitled1* x HW1.Rmd x main0.R x main.R* x
Knit Insert Run
40
41 prognum = matrix(data=c(prog, 1:32),ncol=2)
42 for (i in 1:6){
43   stu[,i+10] = mapvalues(stu[,i+10], from=prognum[,1], to=prognum[,2], warn_missing = FALSE)
44   stu[,i+10][stu[,i+10] == ''] = NA
45 }
46 for (i in 1:6){
47   stu[,i+18] = stu[,i+4]*100+as.numeric(levels(stu[,i+10]))[stu[,i+10]]
48 }
49 choice = as.matrix(stu[,19:24])
50 choice = matrix(choice, ncol=1)
51 choice = unique(choice)
52 choice = choice[!is.na(choice)]
53 length(choice) # Total of 2,773 choices.
54
55 summary(stu$score)["NA's"] # Total of 179,887 missing test scores.
56
57 cc = apply(stu[,5:10], 1, function(x) length(x[!is.na(x)]) - length(unique(x[!is.na(x)])))
58 length(cc[cc != 0])
59 rm(cc) # Total of 120,071 students apply to the same school.
60
61 sum(apply(stu[,11:16],1,function(x) sign(sum(is.na(x)))))
62 # Total of 20,988 students apply to less than 6 choices.
```

```
Untitled1* x HW1.Rmd x main0.R x main.R* x
Knit Insert Run
61 sum(apply(stu[,11:16],1,function(x) sign(sum(is.na(x)))))
62 # Total of 20,988 students apply to less than 6 choices.
63 ```
[1] 340823
[1] 139
[1] 898
[1] 32
[1] 0
NA's
179887
[1] 120071
[1] 20988
64
65 ## Ex 2
66
67 ```{r 1, tidy=TRUE}
68 stu$schoolcode = NA
69 stu$prog = NA
70 for (i in 1:6){
71   stu$schoolcode[which(stu$rankplace == i)] = stu[,i+4][which(stu$rankplace == i)]
72   stu$prog[which(stu$rankplace == i)] = stu[,i+18][which(stu$rankplace == i)]%%100
```

Untitled1* x HW1.Rmd x main0.R x main.R* x

Knit Insert Run

```

71 stu$schoolcode[which(stu$rankplace == i)] = stu[,i+4][which(stu$rankplace == i)]
72 stu$prog[which(stu$rankplace == i)] = stu[,i+18][which(stu$rankplace == i)]%100
73 }
74 sch = summaryBy(score ~ schoolcode + prog, FUN = c(min, mean, length), data = stu)
75 sch = sqldf("SELECT * FROM sch LEFT JOIN sss USING(schoolcode)")
76 sch$sc = sch$schoolcode*100 + sch$prog
77
78 head(sch, n = 20)
79
80 ^

```

	schoolcode <int>	prog <dbl>	score.min <dbl>	score.mean <dbl>	score.length <dbl>	sssdistrict <chr>
1	10101	NA	284	320.2312	398	Accra Metropolitan
2	10102	NA	343	394.1492	248	Accra Metropolitan
3	10103	NA	316	353.8330	443	Accra Metropolitan
4	10104	NA	245	296.9182	220	Accra Metropolitan
5	10105	NA	260	351.2139	346	Accra Metropolitan
6	10106	NA	293	340.1013	395	Accra Metropolitan
7	10107	NA	281	311.9542	306	Accra Metropolitan
8	10108	NA	248	303.9025	318	Accra Metropolitan

Untitled1* x HW1.Rmd x main0.R x main.R* x

Knit Insert Run

8	10108	NA	248	303.9025	318	Accra Metropolitan
9	10109	NA	257	281.8233	300	Accra Metropolitan
10	10110	NA	343	408.0785	535	Accra Metropolitan

1-10 of 20 rows | 1-7 of 9 columns Previous 1 2 Next

```

81
82 ## Ex 3
83
84 {r 1, tidy=TRUE}
85 stu = sqldf("SELECT * FROM stu LEFT JOIN jss USING(jssdistrict)")
86 stu = sqldf("SELECT * FROM stu LEFT JOIN sss USING(schoolcode)")
87 stu$dist = sqrt((69.172*(stu$point_x - stu$ssslong) * cos(stu$point_y/57.3))^2
88 + (69.172*(stu$point_y - stu$ssslat))^2)
89 head(stu, n = 20)
90
91 ^

```

	X <int>	score <int>	agey <int>	male <int>	schoolcode1 <int>	schoolcode2 <int>	schoolcode3 <int>	schoolcode4 <int>
1	1	NA	16	0	50112	50107	50202	50202

	X	score	agey	male	schoolcode1	schoolcode2	schoolcode3	schoolcode4
	<int>	<int>	<int>	<int>	<int>	<int>	<int>	<int>
1	1	NA	16	0	50112	50107	50202	50202
2	2	NA	17	0	70102	70602	70107	70105
3	3	NA	19	0	50702	50705	50115	50706
4	4	NA	23	1	90501	90403	90101	9090401
5	5	NA	15	0	51802	51701	50205	50207
6	6	NA	15	0	10102	50103	51701	50202
7	7	NA	22	1	80301	80401	80302	80402
8	8	NA	19	1	40301	40401	40402	40302
9	9	NA	19	1	21303	21303	21201	21201
10	10	NA	16	0	80101	90401	50503	50901

1-10 of 20 rows | 1-9 of 32 columns

Previous 1 2 Next

```

93 ## Ex4
94
95 ```{r 1, tidy=TRUE}
96 # cutoff (V33:V38):
97 for (i in 1:6){
98   stu[,i+32] = mapvalues(stu[,i+18], from=sch[,9], to=sch[,3], warn_missing = FALSE)
99   stu[,i+32][stu[,i+32] > 500] = NA
100 }
101 # quality (V39:V44):
102 for (i in 1:6){
103   stu[,i+38] = mapvalues(stu[,i+18], from=sch[,9], to=sch[,4], warn_missing = FALSE)
104   stu[,i+38][stu[,i+38] > 500] = NA
105 }
106 # longitude (V45:V50):
107 for (i in 1:6){
108   stu[,i+44] = mapvalues(stu[,i+18], from=sch[,9], to=sch[,7], warn_missing = FALSE)
109   stu[,i+44][stu[,i+44] > 100] = NA
110 }
111 # latitude (V51:V55):
112 for (i in 1:6){
113   stu[,i+50] = mapvalues(stu[,i+18], from=sch[,9], to=sch[,8], warn_missing = FALSE)
114   stu[,i+50][stu[,i+50] > 100] = NA
115 }

```

```

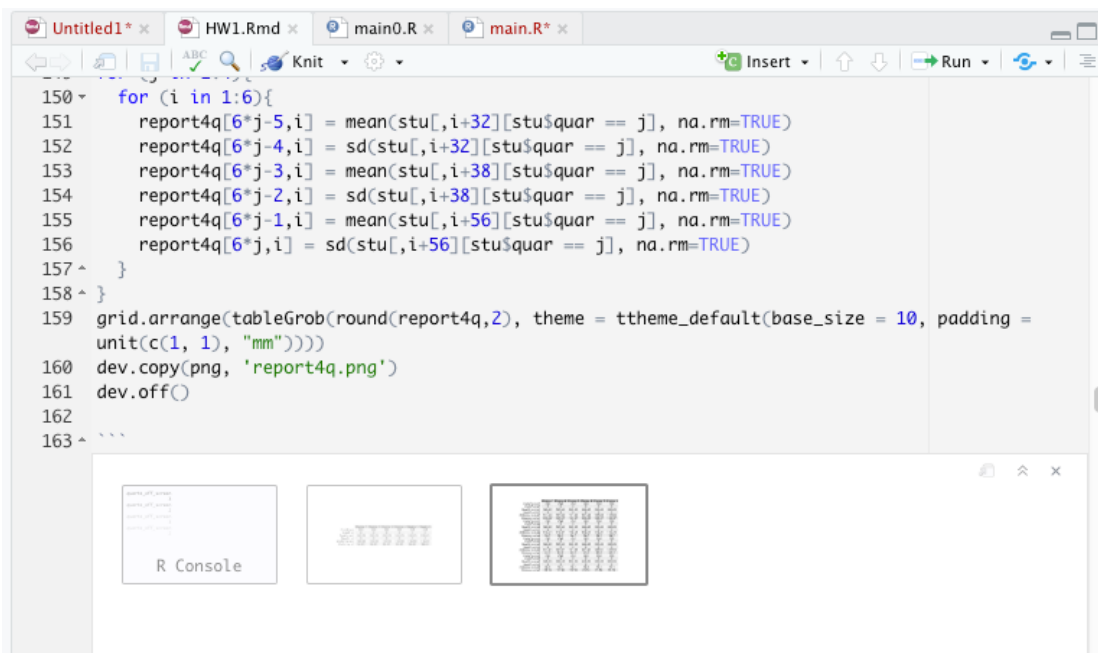
114 stu[,1+50][stu[,1+50] > 100] = NA
115 }
116 # distance (V56:V61):
117 for (i in 1:6){
118   stu[,i+56] = sqrt((69.172*(stu$point_x - stu[,i+44]) * cos(stu[,i+50]/57.3))^2
119                 + (69.172*(stu$point_y - stu[,i+50]))^2)
120 }
121
122 # create tables to report Exercise 4:
123
124 report4 = matrix(rep(NA, 36), ncol=6, byrow=TRUE)
125 colnames(report4) = c('Choice 1', 'Choice 2', 'Choice 3', 'Choice 4', 'Choice 5', 'Choice 6')
126 rownames(report4) = c('Cutoff avg', 'Cutoff sd', 'Quality avg', 'Quality sd', 'Distance avg',
127                       'Distance sd')
127 for (i in 1:6){
128   report4[1,i] = mean(stu[,i+32], na.rm=TRUE)
129   report4[2,i] = sd(stu[,i+32], na.rm=TRUE)
130   report4[3,i] = mean(stu[,i+38], na.rm=TRUE)
131   report4[4,i] = sd(stu[,i+38], na.rm=TRUE)
132   report4[5,i] = mean(stu[,i+56], na.rm=TRUE)
133   report4[6,i] = sd(stu[,i+56], na.rm=TRUE)
134 }

```

```

135 grid.arrange(tableGrob(round(report4,2), theme = ttheme_default(base_size = 10, padding =
136   unit(c(1, 1), "mm"))))
137 setwd("~/Desktop/613")
138 dev.copy(png, 'report4.png')
139 dev.off()
140 # students' quartile of score:
141 stu$quar = cut(stu$score, breaks = quantile(stu$score, probs = seq(0, 1, .25), na.rm = TRUE),
142               include.lowest = TRUE, labels = 1:4)
143 report4q = matrix(rep(NA, 36*4), ncol=6, byrow=TRUE)
144 colnames(report4q) = colnames(report4)
145 rownames(report4q) = c('Cutoff avg q1', 'Cutoff sd q1', 'Quality avg q1', 'Quality sd q1',
146                       'Distance avg q1', 'Distance sd q1',
147                       'Cutoff avg q2', 'Cutoff sd q2', 'Quality avg q2', 'Quality sd q2',
148                       'Distance avg q2', 'Distance sd q2',
149                       'Cutoff avg q3', 'Cutoff sd q3', 'Quality avg q3', 'Quality sd q3',
150                       'Distance avg q3', 'Distance sd q3',
151                       'Cutoff avg q4', 'Cutoff sd q4', 'Quality avg q4', 'Quality sd q4',
152                       'Distance avg q4', 'Distance sd q4')
149 for (j in 1:4){
150   for (i in 1:6){
151     report4q[6*j-5,i] = mean(stu[,i+32][stu$quar == j], na.rm=TRUE)

```



	Choice 1	Choice 2	Choice 3	Choice 4	Choice 5	Choice 6
Cutoff avg	284	284	284	284	284	284
Cutoff sd	0	0	0	0	0	0
Quality avg	320.23	320.23	320.23	320.23	320.23	320.23
Quality sd	0	0	0	0	0	0
Distance avg	116.99	116.99	116.99	116.99	116.99	116.99
Distance sd	96.83	96.83	96.83	96.83	96.83	96.83

	Choice 1	Choice 2	Choice 3	Choice 4	Choice 5	Choice 6
Cutoff avq q1	284	284	284	284	284	284
Cutoff sd q1	0	0	0	0	0	0
Quality avq q1	320.23	320.23	320.23	320.23	320.23	320.23
Quality sd q1	0	0	0	0	0	0
Distance avq q1	131.63	131.63	131.63	131.63	131.63	131.63
Distance sd q1	105.37	105.37	105.37	105.37	105.37	105.37
Cutoff avq q2	284	284	284	284	284	284
Cutoff sd q2	0	0	0	0	0	0
Quality avq q2	320.23	320.23	320.23	320.23	320.23	320.23
Quality sd q2	0	0	0	0	0	0
Distance avq q2	116.41	116.41	116.41	116.41	116.41	116.41
Distance sd q2	92.7	92.7	92.7	92.7	92.7	92.7
Cutoff avq q3	284	284	284	284	284	284
Cutoff sd q3	0	0	0	0	0	0
Quality avq q3	320.23	320.23	320.23	320.23	320.23	320.23
Quality sd q3	0	0	0	0	0	0
Distance avq q3	101.91	101.91	101.91	101.91	101.91	101.91
Distance sd q3	84.4	84.4	84.4	84.4	84.4	84.4
Cutoff avq q4	284	284	284	284	284	284
Cutoff sd q4	0	0	0	0	0	0
Quality avq q4	320.23	320.23	320.23	320.23	320.23	320.23
Quality sd q4	0	0	0	0	0	0
Distance avq q4	76.1	76.1	76.1	76.1	76.1	76.1
Distance sd q4	74.18	74.18	74.18	74.18	74.18	74.18

```

164
165 ## Ex 5
166
167 ```{r 1, tidy=TRUE}
168 set.seed(1)
169 X1 = runif(10000, 1, 3)
170 X2 = rgamma(10000, 3, scale=2)
171 X3 = rbinom(10000, 1, .3)
172 eps = rnorm(10000, 2, 1)
173 b = c(.5, 1.2, -.9, .1)
174 X = matrix(c(rep(1, 10000), X1, X2, X3), ncol=4)
175 Y = (X %>% as.matrix(b)) + eps
176 ydum = as.numeric(Y > mean(Y))
177
178
179 ^ ```
180

```

```

181 ## Ex 6
182
183 ```{r 1, tidy=TRUE}
184 cor(X1, Y)
185 # The correlation is significantly different from 1.2.
186
187 coefYX = solve(t(X)%*%X)%*%t(X)%*%Y
188 residYX = Y - X%*%coefYX
189 varYX = (t(residYX)%*%residYX)/(length(Y)-length(b))
190 vcovYX = c(varYX) * solve(t(X)%*%X)
191 seYX = sqrt(diag(vcovYX))
192
193 # likelihood function for the probit model:
194 probitl = function(beta, y, X){
195   p = pnorm(X%*%beta)
196   l = sum(y*log(p)) + sum((1-y)*log(1-p))
197   return(l)
198 }
199 probitlb = function(par){
200   -probitl(par, ydum, X)
201 }

```

```

199 ~ probitlb = function(par){
200   -probitl(par, ydum, X)
201 ~ }
202 # gradient function for the probit model:
203 ~ dprobitl = function(beta, y, X){
204   p = pnorm(X%*%beta)
205   f = dnorm(X%*%beta)
206   n = length(y)
207   k = length(beta)
208   d = t(matrix(rep(f/p, k), nrow=n)*X)%*%y -
209       t(matrix(rep(f/(1-p), k), nrow=n)*X)%*%(1-y)
210   return(d)
211 ~ }
212 ~ dprobitlb = function(par){
213   -dprobitl(par, ydum, X)
214 ~ }
215 # steepest ascent optimization algorithm:
216 par = rep(0,4)
217 for (i in 1:10000){
218   par = par - .00001 * c(dprobitlb(par))
219 ~ }
220
221 ~ ...

```

```

220
221 ~ ...

```

```

      [,1]
[1,] 0.2162589

```

```

222
223 ~ ## Ex 7
224
225 ~ ```{r 1, tidy=TRUE}
226 # probit model:
227 optimprobit = optim(par=rep(0,4), fn=probitlb, hessian=T)
228 # logit model:
229 ~ logitl = function(beta, y, X){
230   p = 1/(1 + exp(-X%*%beta))
231   l = sum(y*log(p)) + sum((1-y)*log(1-p))
232   return(l)
233 ~ }
234 ~ logitlb = function(par){
235   -logitl(par, ydum, X)
236 ~ }
237 optimlogit = optim(par=rep(0,4), fn=logitlb, hessian=T)

```

```

237 optimlogit = optim(par=rep(0,4), fn=logitlb, hessian=T)
238 # linear probability model:
239 ~ lpms = function(par){
240   s = (ydum - X%*%par)^2
241   return(sum(s))
242 ~ }
243 optimlpm = optim(par=rep(0,4), fn=lpms, hessian=T)
244 # The probit model returns parameters closest to the actual values,
245 # and the other two models have approximately the same accuracy.
246 tprobit = optimprobit$par/sqrt(diag(solve(optimprobit$hessian)))
247 tlogit = optimlogit$par/sqrt(diag(solve(optimlogit$hessian)))
248 tlpm = optimlpm$par/sqrt(diag(solve(optimlpm$hessian)))
249 # All the estimates from all the three models are statistically
250 # significant, except the last parameter in the linear probability
251 # model under 95% confidence level.
252
253 ~ ...
254

```



```

255 ~ ## Ex 8
256
257 ~ ```{r 1, tidy=TRUE}
258 # probit model:  $F'(x) = (-\exp(-x^2/2))/\sqrt{2\pi}$ 
259 # logit model:  $F'(x) = \exp(-x)/(1 + \exp(-x))^2$ 
260 meoptimprobit = dnorm(X%%optimprobit$par)%%optimprobit$par
261 meoptimlogit = (exp(-X%%optimlogit$par)/(1 + exp(-X%%optimlogit$par))^2)%%optimlogit$par
262 # probit standard error:
263 vcovprobit = solve(optimprobit$hessian)
264 jacprobit = array(0,dim=c(4,4,10000))
265 ~ for (k in 1:10000){
266   jacprobit[,k] = diag(4)
267   f = dnorm(X[k,]%%optimprobit$par)
268   ~ for (i in 1:4){
269     ~ for (j in 1:4){
270       jacprobit[i,j,k] = f*(jacprobit[i,j,k]-
271         (X[k,]%%optimprobit$par)*optimprobit$par[j]*X[k,i])
272     }
273   }
274 ~ }
275 jbprobit = rowMeans(jacprobit, dims=2)
276 deltaprobit = sqrt(diag(t(jbprobit)%%vcovprobit%%jbprobit))

```

```

275 ~ }
274 ~ }
275 jbprobit = rowMeans(jacprobit, dims=2)
276 deltaprobit = sqrt(diag(t(jbprobit)%%vcovprobit%%jbprobit))
277 # logit standard error:
278 vcovlogit = solve(optimlogit$hessian)
279 jaclogit = array(0,dim=c(4,4,10000))
280 ~ for (k in 1:10000){
281   jaclogit[,k] = diag(4)
282   f = 1/(1 + exp(-X[k,]%%optimlogit$par))
283 ~ for (i in 1:4){
284   ~ for (j in 1:4){
285     jaclogit[i,j,k] = f * (1-f) * (
286       jaclogit[i,j,k] + (1-2*f)*(optimlogit$par[j]*X[k,i])
287     )
288   }
289 ~ }
290 ~ }
291 jblogit = rowMeans(jaclogit, dims=2)
292 deltalogit = sqrt(diag(t(jblogit)%%vcovlogit%%jblogit))
293
294
295 ~ ```
296
179:4 ~ ## Ex 5 ~

```