

ARE213 Problem Set #1B

Peter Alstone & Frank Proulx

October 14, 2013

1 Problem #1

1.1 Part A

Under the assumption of random assignment conditional on the observables, what are the sources of misspecification bias in the estimates generated by the linear model estimated in Problem Set 1a?

Wrong functional form. In Problem Set 1A we used linear (i.e., $y = \beta x + \epsilon$) estimators to make “corrections” while the true functional form of the relationships between the covariates we included in the model were certainly not linear. By imposing a linear function on a non-linear data generating process (described by the CEF), we introduce misspecification bias in the model.

Omitted Variables Bias. We were able to use variables included in the dataset in our linear model, but not the unobserved variables that may be important for control. If omitted variables exist that both determine outcomes related to birth weight and are correlated with smoking status we will over- or under-estimate the effect (depending on the characteristics of the omission).

1.2 Part B

Now, consider a series estimator. Estimate the smoking effects using a flexible functional form for the control variables (e.g., higher order terms and interactions). What are the benefits and drawbacks to this approach?

Using a

2 Problem #2

2.1 Part A

To calculate the propensity score, we estimated a logit model of mother's tobacco use (0=non-smoker, 1=smoker) as determined by the predetermined covariates shown in Table ?? . Model #1 shows the full model using all of the covariates suspected to be predetermined. Model #2 is a reduced form of the same model, preserving just those covariates that were significant at the 1% level in Model #1.

To test whether the propensity scores predicted by these two models are significantly different we perform a likelihood ratio test between the full and reduced model. This test yields the following output:

NOTE from class: including insignificant terms can be beneficial in terms of getting better fit (by reducing

2.2 Part B

This estimation assumes unconfoundedness. and conditional independence

2.3 Part C

We used R to complete this assignment. The code is below:

```
1 # PROBLEM SET 1B
2 # ARE 213 Fall 2013
3
4
5 # Frank's Directory
6 #setwd("/media/frank/Data/documents/school/berkeley/fall13/are213/are213/ps1
7   ")
8 # Peter's Directory
9 #setwd("~/Google Drive/ERG/Classes/ARE213/are213/ps1")
10
11
12 # Packages
13 library(foreign) #this is to read in Stata data
14 library(Hmisc)
15 library(psych)
16 library(stargazer)
17 library(ggplot2) # for neato plotting tools
```

Table 1: Propensity scores calculated for mother's smoking status

| | Mother Tobacco-Use Status | |
|---------------------------------------|---------------------------|----------------------|
| | (1) | (2) |
| Mother's Race not White or Black | -1.956*** (0.134) | -1.954*** (0.133) |
| Mother's Years of Education | -0.817*** (0.028) | -0.818*** (0.028) |
| Marital status | -0.205*** (0.005) | -0.204*** (0.005) |
| Father's age | -1.256*** (0.022) | -1.251*** (0.021) |
| Father's Years of Education | 0.029*** (0.002) | 0.030*** (0.001) |
| Father Mexican | -0.131*** (0.005) | -0.131*** (0.005) |
| Father Puerto Rican | -1.961*** (0.173) | -1.957*** (0.173) |
| Father Cuban | -1.267*** (0.058) | -1.268*** (0.058) |
| Father Central or South American | -0.567 (0.364) | -0.567 (0.364) |
| Father Race Other or Unknown Hispanic | -1.933*** (0.205) | -1.932*** (0.205) |
| Plurality of Infant | -0.890*** (0.120) | -0.889*** (0.120) |
| Sex of Infant | -0.148*** (0.054) | |
| Mother's age | -0.019 (0.017) | |
| dmage | 0.003 (0.002) | |
| Constant | 2.873*** (0.088) | 2.707*** (0.064) |
| <i>N</i> | 114,610 | 114,610 |
| Log Likelihood | -44,310.690 | -44,315.790 |
| Akaike Inf. Crit. | 88,651.370 | 88,655.580 |

Notes:

***Significant at the 1 percent level.

**Significant at the 5 percent level.

```

18 library(plyr) # for nice data tools like ddply
19 library(epicalc) # For likelihood ratio test
20 library(car) # "companion for applied regression" - recode fxn, etc.
21 library(gmodels) #for Crosstabs
22 library(splines) # for series regression
23
24 # Homebrewed functions
25 source("../util/are213-func.R")
26 source("../util/watercolor.R") # for watercolor plots
27
28 # Data
29 ps1.data <- read.dta(file="ps1.dta")
30
31 # Data Cleaning Step
32 full.record.flag <- which(ps1.data$cardiac != 9 &
33                           ps1.data$cardiac != 8 &
34                           ps1.data$lung != 9 &
35                           ps1.data$lung != 8 &
36                           ps1.data$diabetes !=9 &
37                           ps1.data$diabetes !=8 &
38                           ps1.data$herpes != 9 &
39                           ps1.data$herpes != 8 &
40                           ps1.data$chyper != 9 &
41                           ps1.data$chyper != 8 &
42                           ps1.data$phyper != 9 &
43                           ps1.data$phyper != 8 &
44                           ps1.data$pre4000 !=9 &
45                           ps1.data$pre4000 !=8 &
46                           ps1.data$preterm != 9 &
47                           ps1.data$preterm != 8 &
48                           ps1.data$tobacco != 9 &
49                           ps1.data$cigar != 99 &
50                           ps1.data$cigar6 !=6 &
51                           ps1.data$alcohol != 9 &
52                           ps1.data$drink != 99 &
53                           ps1.data$drink5 !=5 &
54                           ps1.data$wgain !=99
55                           )
56
57 ps1.data$full.record <- FALSE # initialize column as F
58 ps1.data$full.record[full.record.flag] <- TRUE #reassign level to T for full
   records
59
60 ps1.data.clean <- subset (ps1.data, full.record == TRUE)
61 ps1.data.missingvalues <- subset(ps1.data, full.record == FALSE)
62
63 # Problem 1a is only in ps1b.tex
64 # Problem 1b
65 # This is using a series estimator. I think smooth.spline() is the right
   function to use, but let me know if you think we should be doing kernel
   regression instead. I'm also not sure how to go about adding interaction
   terms
66
67
68
69 sm.flex <- with(ps1.data.clean, smooth.spline(cigar, y=dbrwt, nknots=10,
   spar = 0.7, tol = 0.0001)) # Fits a smooth line to the data

```

```

70 sm.flex.df <- data.frame(sm.flex$x, sm.flex$y) #converts the fitted values
    into a data frame for ggplot
71
72 splineplot <- ggplot(sm.flex.df, aes(x = sm.flex.x, y=sm.flex.y))
73 splineplot <- splineplot +
74   geom_point(data=ps1.data.clean, aes(x = cigar, y = dbrwt), pch = 1) +
75   geom_line(color='red') +
76   labs(x = 'Cigarettes smoked per day by mother', y= 'Birthweight')
77 splineplot
78
79 ggsave(filename = 'img/splineplot.pdf')
80
81
82 # Problem 2a
83 ps1.data.clean$tobacco.rescale <- with(ps1.data.clean, recode(tobacco, "
    2='0'", as.numeric.result=TRUE)) #rescales the tobacco use variable to
    be 0/1, where 0=no and 1 = yes
84 ps1.data.clean$dmr.rescale <- with(ps1.data.clean, recode(dmr, "2='0'"))
85
86 smoke.propensity.all <- glm(tobacco.rescale ~ as.factor(mrace3) + dmeduc +
    dmr.rescale + dfage + dfeduc + as.factor(orfath) + dplural + csex +
    dmage, data=ps1.data.clean, family = binomial()) ## Did I miss any
    predetermined covariates here?
87
88 smoke.propensity.reduced <- glm(tobacco.rescale ~ as.factor(mrace3) + dmeduc
    + dmr.rescale + dfage + dfeduc + as.factor(orfath), data=ps1.data.
    clean, family = binomial())
89
90
91 stargazer(smoke.propensity.all, smoke.propensity.reduced,
92   type = "latex",
93   covariate.labels = c("Mother's Race not White or Black", "Mother's
    Years of Education", "Marital status", "Father's age", "
    Father's Years of Education", "Father Mexican", "Father
    Puerto Rican", "Father Cuban", "Father Central or South
    American", "Father Race Other or Unknown Hispanic", "
    Plurality of Infant", "Sex of Infant", "Mother's age"),
94   style = "qje",
95   align = TRUE,
96   font.size = "scriptsize",
97   label = "tab:propensities",
98   title = "Propensity scores calculated for mother's smoking status
    \n
    logit model specification",
99   dep.var.labels = "Mother Tobacco-Use Status",
100   out = "propensitiescores.tex"
101 )
102
103
104 ps1.data.clean$propensityfull <- predict(smoke.propensity.all, type = "
    response")
105 ps1.data.clean$propensityreduced <- predict(smoke.propensity.reduced, type =
    "response")
106
107 sink(file = "lrtest.tex", append = FALSE)
108 lrtest(smoke.propensity.all, smoke.propensity.reduced) #Test whether the two
    scores are statistically different
109 sink()

```

```

110
111 #Problem 2b - Estimating a regression model using propensity scores
112
113 sm.propensityregression <- lm(dbrwt ~ tobacco.rescale + (propensityreduced *
114   tobacco.rescale) + propensityreduced, ps1.data.clean)
115
116 #calculation of average treatment effect:
117 coefficients(sm.propensityregression)[2] + coefficients(sm.
118   propensityregression)[4]*mean(ps1.data.clean$propensityreduced)
119
120 stargazer(sm.propensityregression,
121   type = "latex",
122   covariate.labels = c("Delta1", "Beta", "Delta2", "Constant"),
123   style = "qje",
124   align = TRUE,
125   font.size="footnotesize",
126   label = "tab:propensitymodel",
127   title = "Model of effects of tobacco use on birthweight using
128     propensity score as a control",
129   dep.var.labels = "Mother Tobacco-Use Status",
130   out = "propensityscoremodel.tex"
131 )
132
133 #Problem 2c - Using reweighting with propensity scores
134
135 term1 <- with(ps1.data.clean, sum((tobacco.rescale*dbrwt)/propensityreduced)/
136   sum(tobacco.rescale/propensityreduced))
137 term2 <- with(ps1.data.clean, sum(((1-tobacco.rescale)*dbrwt)/(1-
138   propensityreduced))/sum((1-tobacco.rescale)/(1-propensityreduced)))
139
140 weightingestimator <- term1-term2 #This should be the average treatment
141   effect
142
143 term1.T <- with(subset(ps1.data.clean, tobacco.rescale=1), sum((tobacco.
144   rescale*dbrwt)/propensityreduced)/sum(tobacco.rescale/propensityreduced)
145 )
146 #term2.T <- with(subset(ps1.data.clean, tobacco.rescale=1), sum(((1-tobacco.
147   rescale)*dbrwt)/(1-propensityreduced))/sum((1-tobacco.rescale)/(1-
148   propensityreduced)))
149
150 weightingestimator.T <- term1.T-term2.T #This should be the average
151   treatment on treated
152
153
154 # Problem 2d - Kernel Density Estimator
155 tot.propensity.nosm <- with(subset(ps1.data.clean, tobacco.rescale == 0),
156   sum(propensityreduced))
157 tot.propensity.sm <- with(subset(ps1.data.clean, tobacco.rescale == 1), sum(
158   propensityreduced))
159 h <- 35 # This is to play with the bandwidth
160
161 kerndensity.nosm <- with(subset(ps1.data.clean, tobacco.rescale == 0),
162   density(dbrwt, #if nobody smoked
163     kernel = "epanechnikov",
164     bw = h,

```

```

153     weights = propensityreduced/tot.propensity.nosm))
154 kerndensity.nosm.df <- data.frame(kerndensity.nosm[1], kerndensity.nosm[2])
155
156 kerndensity.sm <- with(subset(ps1.data.clean, tobacco.rescale == 1), density
157   (dbrwt, #if everybody smoked
158     kernel = "epanechnikov",
159     bw = h,
160     weights = propensityreduced/tot.propensity.sm))
161 kerndensity.sm.df <- data.frame(kerndensity.sm[1], kerndensity.sm[2])
162
163 kerndensity.plot <- ggplot(kerndensity.nosm.df, aes(x, y))
164 kerndensity.plot <- kerndensity.plot +
165   geom_line(linetype = 'dotted') +
166   geom_line(data = kerndensity.sm.df, aes(x, y)) +
167   labs(title = paste("Density of birthweights estimated using \n propensity
168     score-weighted kernel regression \n Bandwidth=", as.factor(h)), x = "
169     Birthweight (grams)", y = "Density") +
170   guides(linetype = "Legend") # Having trouble getting a legend.
171
172 kerndensity.plot
173
174 ggsave(file = 'img/kerndensity.pdf', plot = kerndensity.plot)
175
176 #2d - calculating kernel value by 'hand' at dbrwt = 3000
177 # y_pred = sum_onj(K(pi-pj/h)*yj)/sum_onj(K(pi-pj/h))
178
179 #h <- 30
180 #kernel.epa <- function(u){
181 #return(if (abs(u) < 1){0.75*(1-u*u)})}
182 #}
183 #propensity3000 <- # I'm not sure how we get the estimated propensity score
184 # at dbrwt = 3000
185
186 #Problem 2e
187 ## This is in progress
188 #kerndensity.plot.bws <- ggplot()
189 #for(h in seq(from = 15, to = 40, by = 5)) {
190 #paste0('kerndensity.sm.', h) <- with(subset(ps1.data.clean, tobacco.rescale
191 # == 1), density(dbrwt, #if everybody smoked
192 #   kernel = "epanechnikov",
193 #   bw = h,
194 #   weights = propensityreduced/tot.propensity.sm))
195 #paste0('kerndensity.sm.df.', h) <- data.frame(paste0('kerndensity.sm.', h
196 # , '[1]'), paste0('kerndensity.sm.', h, '[2]'))
197 #kerndensity.plot.bws <- kerndensity.plot.bws +
198 #   geom_line(data = paste0('kerndensity.sm.df.', h), aes(x, y)) +
199 #   labs(title = "Density of birthweights estimated using \n propensity score
200 #     -weighted kernel regression", x = "Birthweight (grams)", y = "Density")
201 #}
202
203 #kerndensity.plot.bws

```

ps1b.R

```

1 # Econometrics helper functions for [R]
2 #

```