# ARE213 Problem Set #2B

Peter Alstone & Frank Proulx

November 21, 2013

## Part A:   Preliminaries

### (i)   Comparison between TU and control States

**Starting with simple comparisons**   We begin with simple comparisons between the dependent outcome of interest, the natural logarithm of traffic fatalities per capita (log(fatalities per capita)), between a predefined composite treatment state "TU" (or, state #99), and all of the potential control states. The mean over the period before primary seatbelt laws were adopted in the treatment state is -1.4 and the mean for the control states is -1.7, indicating approximately a 30% lower typical fatalities rate in the treatment state than the average control state (even before the primary seatbelt law "treatment"). The trends for both shown in Figure 1 show that overall the fatalities were on the decline in both places before the treatment period.

**Roadmap**   Extracting meaningful conclusions from these data is the goal of our analysis, which will require identifying the variation in traffic fatalities that can be attributed to seat belt laws. Confounding our analysis is the fact that these data are not in the context of an RCT but are from the "real world" with messy trends and linked systems that determine outcomes. We will be applying the synthetic controls method to identify a fleet of control states as a meaningful counterfactual to measure against for our composite treatment state.
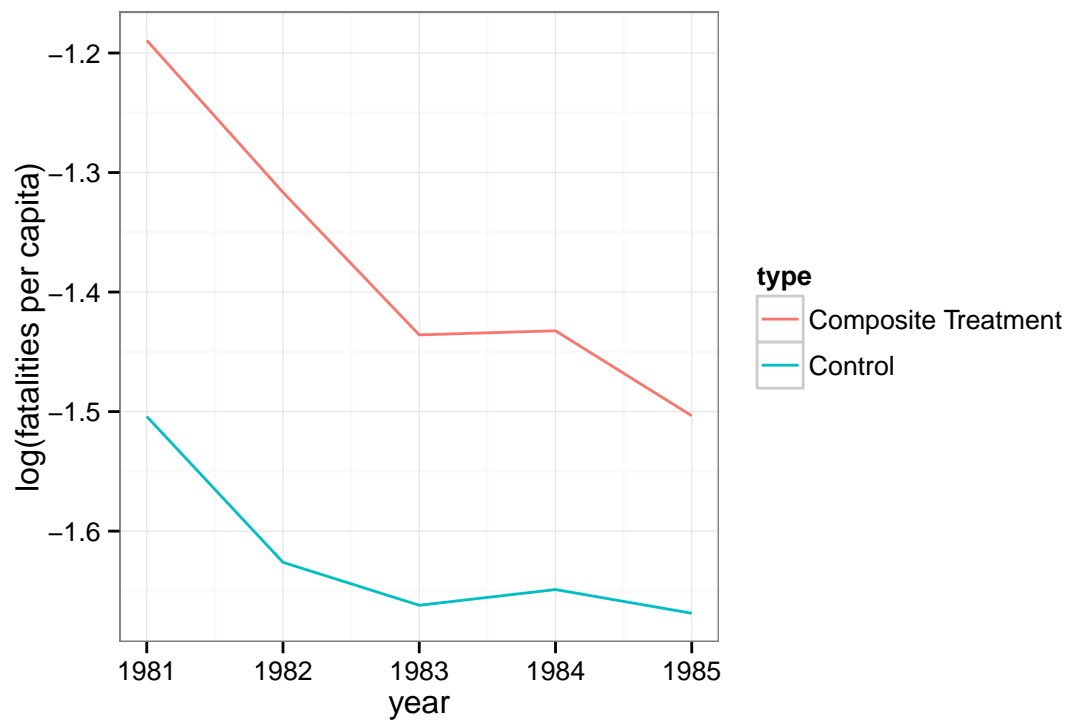
Figure 1: Trend in the dependent variable (log(fatalities per capita)) for the composite treatment state and the average of the control states.

## (ii)  "Best" control state comparison

**Sweet Home Alabama**  We observe that Alabama is the best match for the composite treatment state based on a simple comparison of log(fatalities per capita) in the year before treatment in the composite state (1985). Figure 2 below shows the distribution in the dependent variable

**Fried green covariates and other stereotypes confirmed**  Tables 1 and 2 compare the covariates for the composite treatment state and Alabama. There are broad differences between the states. Alabama has higher precipitation, lower college achievement, lower alcohol consumption, higher unemployment, etc. Additionally, the mean value for the depdentent variable of interest, log(fatalities per capita), is quite different for the two states. Examining the trends in the covariates (and dependent variable) for the two states (see Figure 3) shows that it could be construed as a coincidence that Alabama is the best match for the value of the dependent variable, since the trajectory in fatalities for both states are following opposite trends in that time and 1985 happens to be the time when they intersect. There are also important and long-term differences in precipitation and alcohol consumption.

Overall Alabama does not appear to be a particularly good match for the composite treatment state, motivating an application of synthetic controls methods to produce a better match.

# Part B:  Synthetic Controls

## (i)  Why synthetic controls?

**Unsweet Home Alabama:**  We saw earlier the difficulties in selecting an exact counterfactual match for implementing differences in differences type selection on unobservables techniques. While Alabama would appear on face value to be a good match (based on having similar outcomes in the year prior to treatment) we saw that this was coincidental and that the covariates are not a good match to the composite treatment state. Synthetic control methods are motivated by producing a "better" match by combining (synthesizing) multiple control states in a weighting scheme to create a composite control state with better match of the important covariates and dependent variable than any particular control state.
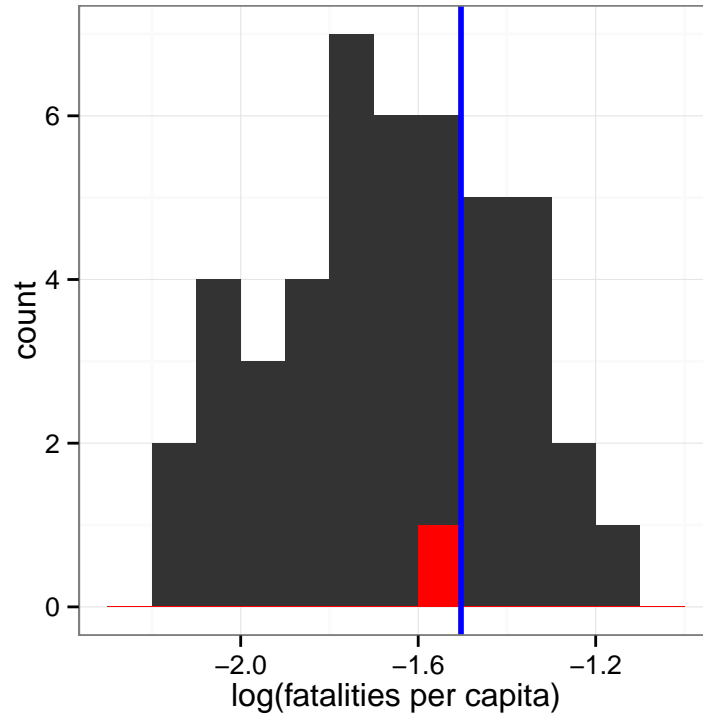
Figure 2: Distribution in traffic fatalities metric from 1985 for all control states with a vertical blue line indicating the value of the metric for the composite treatment state. The red block highlights the position of Alabama in the distribution. Alabama is the closest match to the composite treatment state for 1985, but as is shown here is one of about 11 states that is within 10% of the target value.

Table 1: Composite Treatment Group Summary

| Statistic | N | Mean | St. Dev. | Min | Max |
|---|---|---|---|---|---|
| state | 23 | 99.000 | 0.000 | 99 | 99 |
| year | 23 | 1,992.000 | 6.782 | 1,981 | 2,003 |
| college | 23 | 0.234 | 0.014 | 0.209 | 0.259 |
| beer | 23 | 1.507 | 0.074 | 1.394 | 1.670 |
| primary | 23 | 0.783 | 0.422 | 0 | 1 |
| secondary | 23 | 0.000 | 0.000 | 0 | 0 |
| population | 23 | 13,597.660 | 1,813.520 | 10,737.810 | 16,862.220 |
| unemploy | 23 | 6.085 | 1.124 | 3.855 | 8.014 |
| fatalities | 23 | 2,619.014 | 258.667 | 2,246.977 | 3,268.613 |
| totalvmt | 23 | 128,099.600 | 26,447.260 | 86,013.140 | 170,407.300 |
| precip | 23 | 2.502 | 0.289 | 1.990 | 3.104 |
| snow32 | 23 | 0.143 | 0.058 | 0.013 | 0.270 |
| rural_speed | 23 | 63.443 | 6.568 | 55.000 | 72.886 |
| urban_speed | 23 | 59.184 | 5.858 | 55.000 | 67.138 |
| logfatalpc | 23 | $-1.643$ | 0.168 | $-1.805$ | $-1.189$ |
| sqyears | 23 | 3,968,108.000 | 27,020.830 | 3,924,361 | 4,012,009 |

Table 2: Closest match for pre-policy fatalities: Alabama

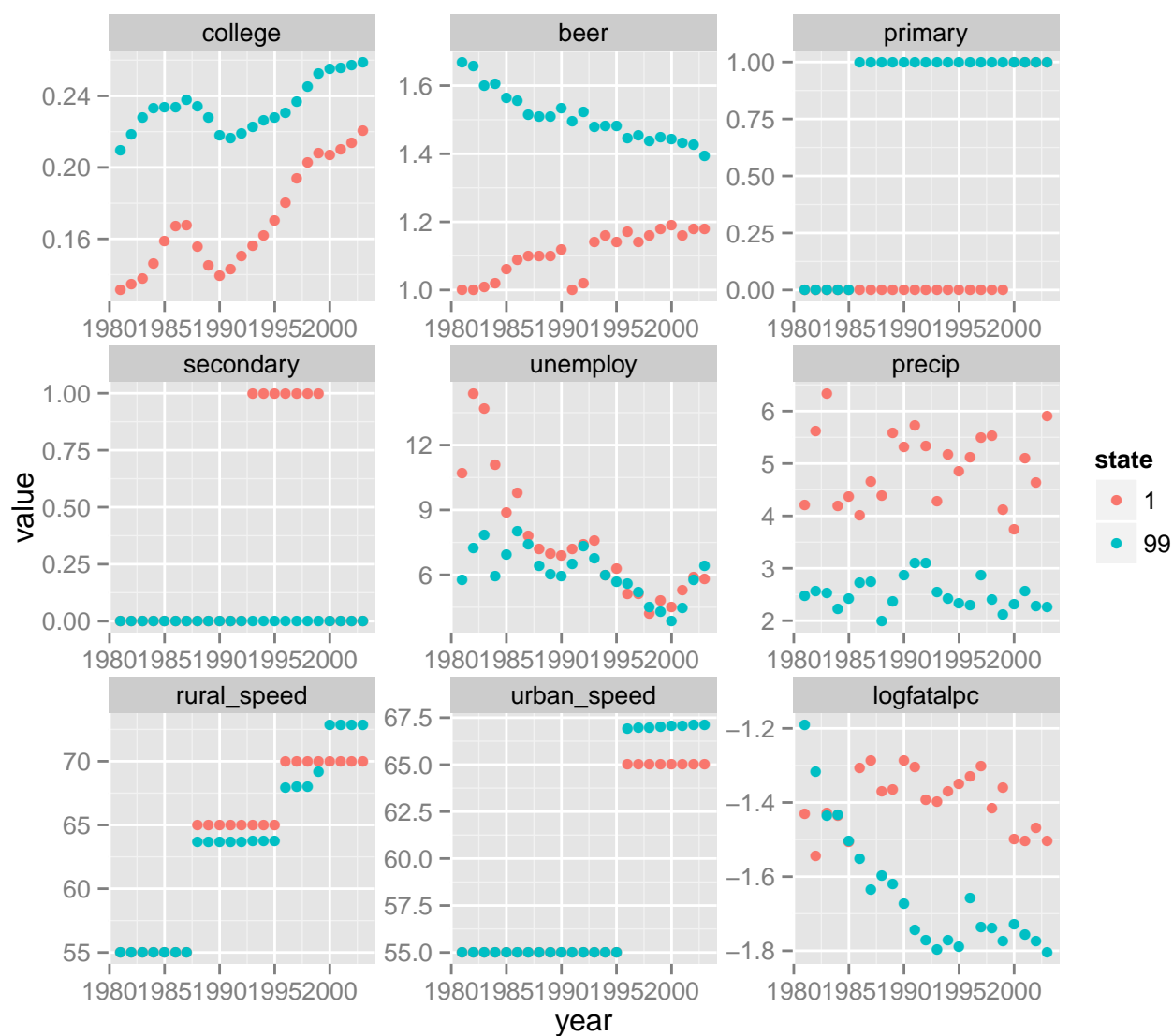| Statistic | N | Mean | St. Dev. | Min | Max |
|---|---|---|---|---|---|
| state | 23 | 1.000 | 0.000 | 1 | 1 |
| year | 23 | 1,992.000 | 6.782 | 1,981 | 2,003 |
| college | 23 | 0.170 | 0.029 | 0.131 | 0.220 |
| beer | 23 | 1.105 | 0.067 | 1.000 | 1.190 |
| primary | 23 | 0.174 | 0.388 | 0 | 1 |
| secondary | 23 | 0.304 | 0.470 | 0 | 1 |
| population | 23 | 4,185.794 | 209.389 | 3,918.533 | 4,501.862 |
| unemploy | 23 | 7.509 | 2.780 | 4.200 | 14.400 |
| fatalities | 23 | 1,036.957 | 88.042 | 839 | 1,189 |
| totalvmt | 23 | 44,826.090 | 10,109.350 | 27,852 | 58,637 |
| precip | 23 | 4.944 | 0.701 | 3.737 | 6.342 |
| snow32 | 23 | 0.000 | 0.000 | 0 | 0 |
| rural_speed | 23 | 63.696 | 6.255 | 55 | 70 |
| urban_speed | 23 | 58.478 | 4.870 | 55 | 65 |
| logfatalpc | 23 | $-1.398$ | 0.079 | $-1.543$ | $-1.286$ |
| sqyears | 23 | 3,968,108.000 | 27,020.830 | 3,924,361 | 4,012,009 |

Figure 3: Trends in the covariate (and dependent) variables for the composite treatment state (99) and Alabama (1)

**Dr. Synth-love, or how I learned to stop worrying and love econometrics:** Synthetic controls have a multi-step, iterative process for developing weighting factors to apply to control states for construction of a composite control state. The goal is to identify a weighting matrix $W$ that minimizes the distance between the treatment covariates (e.g. alcohol consumption, total VMT) and pre-intervention outcomes (log fatalities per capita) for the weighted control unit and the treatment unit. In particular, the following more formally defined criteria are sought (from Synth R package documentation):

- $\sum_{j=2}^{J+1} w_j^* \bar{Y}_j^{K_1} = \bar{Y}_1^{K_M}$ , where $j$ refers to the state, $w_j$ is state $j$'s weight, and $\bar{Y}_j^{K_m}$ denotes the pre-treatment outcome in state $j$ in year $m$

- $\sum_{j=2}^{J+1} w_j^* U_j = U_1$, where $U_i$ is a vector of covariates for state $i$

Pursuant to these criteria, the synthetic control method estimates the treatment effect as $\hat{\alpha}_{1t} = Y_{1t} - \sum_{j=2}^{J+1} w_j^* Y_{jt}$

The steps taken in this estimation by the Synth package are as follows:

1. Define a (k x 1) matrix (dubbed $X_1$) of the characteristics (covariates $U_1$ and pre-treatment outcomes $\bar{Y}_1^{K_m}$) of the treatment unit and a similar (k x J) matrix (dubbed $X_J$) for the control units.

2. Weight the control characteristics matrix with weight vector W.

3. Minimize the distance between the treatment unit characteristic matrix and the weighted control characteristics matrix with respect to the weighting matrix. Formally, that's $\min_W \sqrt{(X_1 - X_0 W)' V (X_1 - X_0 W)}$ where $V$ is chosen by default to minimize the mean square error of the estimator.

**Pros and Cons:** The upsides to Synthetic control are that one can create a better match for the treatment unit than exists in reality and that it is a method that prevents issues of selection bias (i.e., it is possible to say, "I am using synthetic control" instead of needing to justify ex post the selection of particular units to match in classic differences in differences methods). Another nice feature of the method is the use of graphical placebo testing analysis for determination of the statistical power of results. It is an elegant and compelling way to approach error analysis. A potential methodological

downside is that the method approaches a black box estimate that does not provide much intuition compared to other methods. This may manifest as a lack of trust in results from this method compared to those that are more straightforward to understand.

## (ii)    Synthesizing control

The process of creating a synthetic control unit involves 2 steps in the Synth package on [R]. First is specifying the form of the model in a "data prep" step. This is then passed to the synthetic control function to attempt implementing the algorithm described above. In practice we found that errors arise when predictors are included that do not have variation in the mean values among the control units. We used an additive process (adding more and more predictor covariates in the specification) to test whether there is variation. A sub-finding is that the computational intensity increases as covariates are added. This is a relatively small dataset but it is possible that this method could become computationally difficult with large datasets and many covariates. After the process of adding we found that there is variation in all the potentially meaningful covariates except rural and urban speed limits. Since speed limits were constant throughout the sample before 1986 they cannot be included in the synthetic controls specification. Additionally, the presence of secondary seatbelt laws does not vary in the pre-treatment period so is also left out of the potential covariates.

**Preferred specification:**   We identified that the following specifications were best for synthetic control analysis of this data:

- Covariates to include in pre-treatment "training" period: Full set of tractable and reasonable covariates. This includes alcohol consumption, VMT per capita, college educational attainment, precipitation, snowfall, and unemployment rate. We tried other combinations of covariates and found little influence on the result. We hoped a VMT-only specification would provide clarity but the gap in the pre-treatment period was biased compared to using the full set of covariates.

- Pre-treatment period: We use the full set of years available in the data, from 1981 to 1985, for the pre-treatment period. We considered dropping 1985 to avoid anticipation effects but this is not done for two

9

reasons: first, it did not have noticeable impacts on the results (i.e., the divergence between treatment and synthetic control appears between 1985-1986 regardless of whether 1985 is included), and second, because there is a very short pre-treatment period available and we wished to maximize the support of the data.

# Part C:   ...but does it work?

## (i)   Gap between TU and synth control

We show a series of figures (4 - 9) for various specifications (including the preferred specification) below. Both the gaps (Synthetic unit - Treated unit) and the actual values have been plotted here. These plots appear to show a greater pre-treatment MSPE when only the VMT per capita covariate is used as compared with the full model. Ideally, the difference between the synthetic control unit and the treatment unit should be as close to zero as possible in the pre-treatment period, implying an MSPE approaching zero. For this reason, we select the full model using the aforementioned covariates.

The mean gap is about 0.15 on the log scale, which corresponds to approximately a 15% reduction in traffic fatalities per capita.

## (ii)   Gap between TU & preferred synth spec. and gap between each control state and its "placebo" treatment

**Graphical significance?**

**Placebo test plots:**   We developed a series of placebo test plots to investigate the statistical significance of the gap between the treatment and its synthetic control. The series of plots below (Figures 10, 11, 12, 13, and 14) show the implications of different thresholds for keeping control states in the placebo test. The criterion for inclusion is the degree of difference between the pre-treatment (1981-85) MSPE for the placebo control gap vs. the gap for the treatment state. We try a range from basically allowing all the states (50x) to a strict filter (2x). The best compromise between allowing states to participate in the model and keeping out those with bad synthetic control is
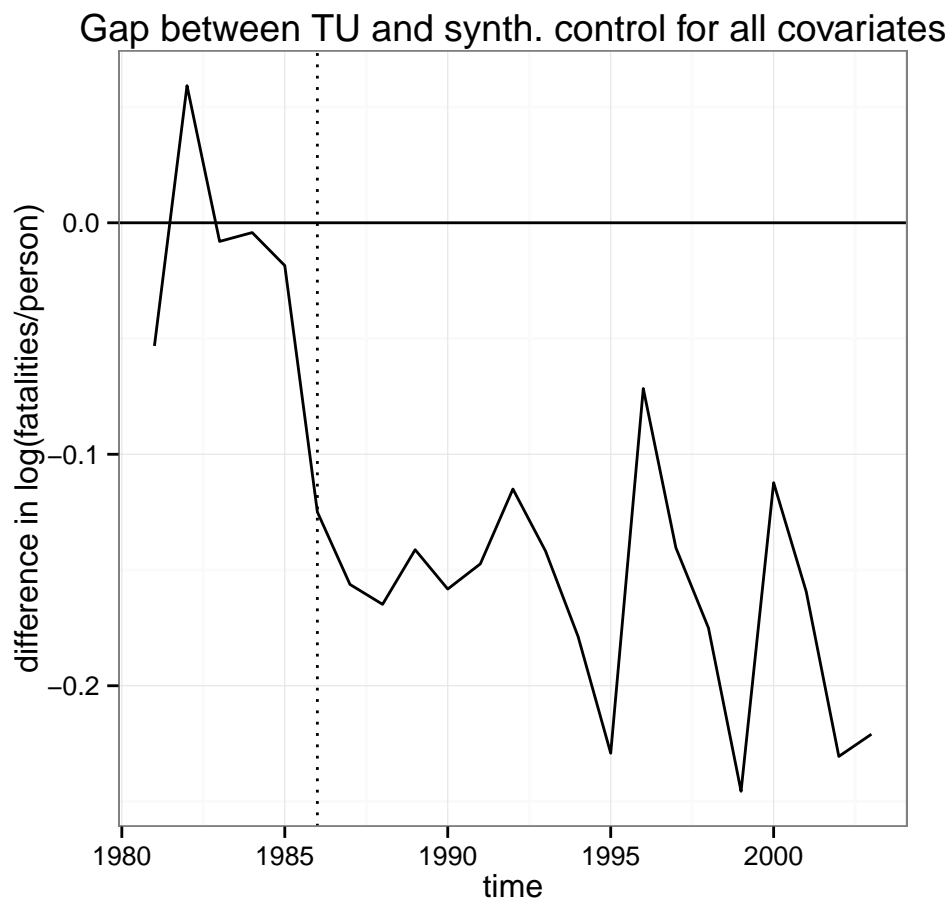
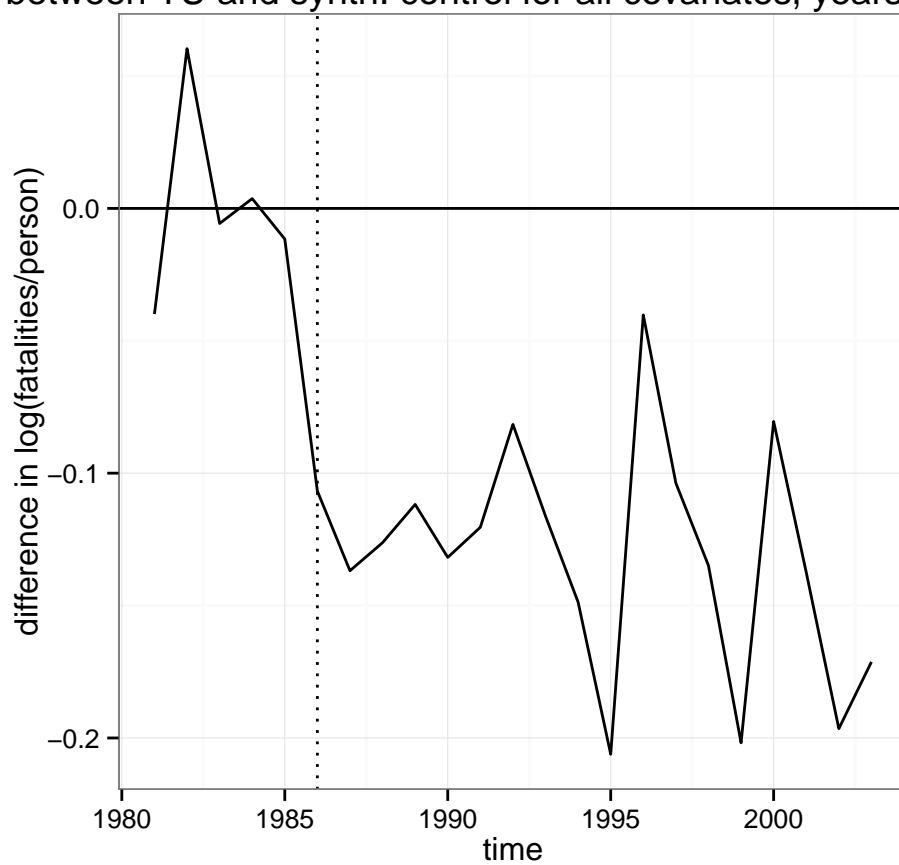Figure 4: Gap between Treatment Unit and Synthetic Control developed using all previous time periods and all covariates.

Figure 5: Gap between Treatment Unit and Synthetic Control developed using time periods 1981-1984 and all covariates.
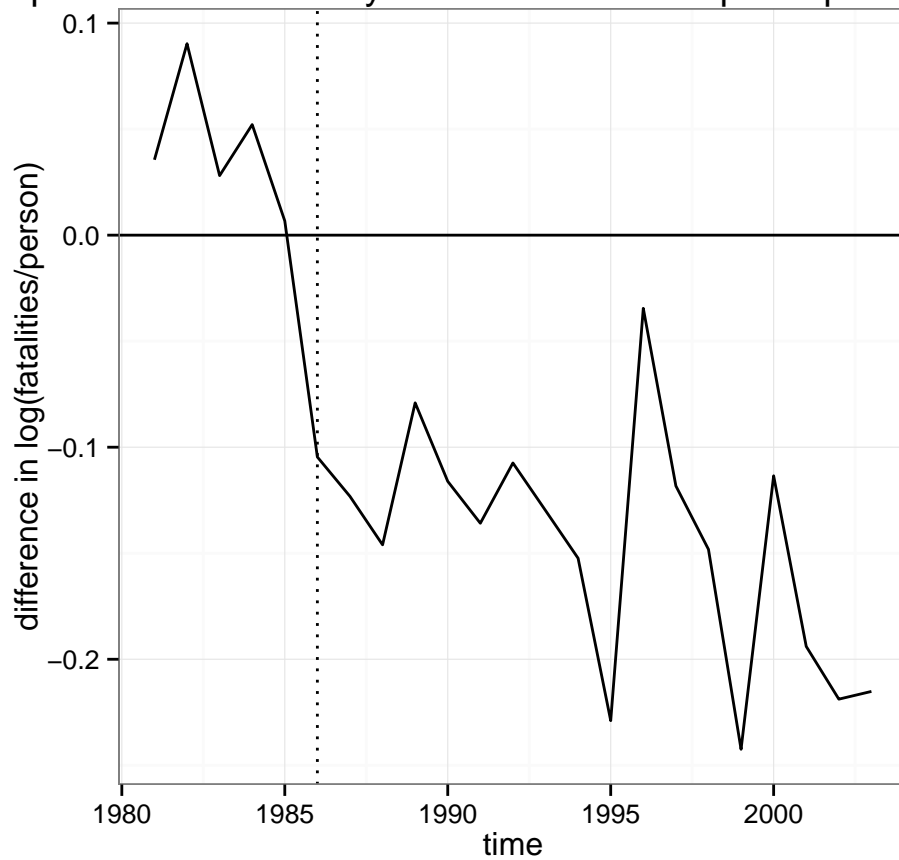
Figure 6: Gap between Treatment Unit and Synthetic Control developed using all previous time periods and VMT per capita as a covariate.
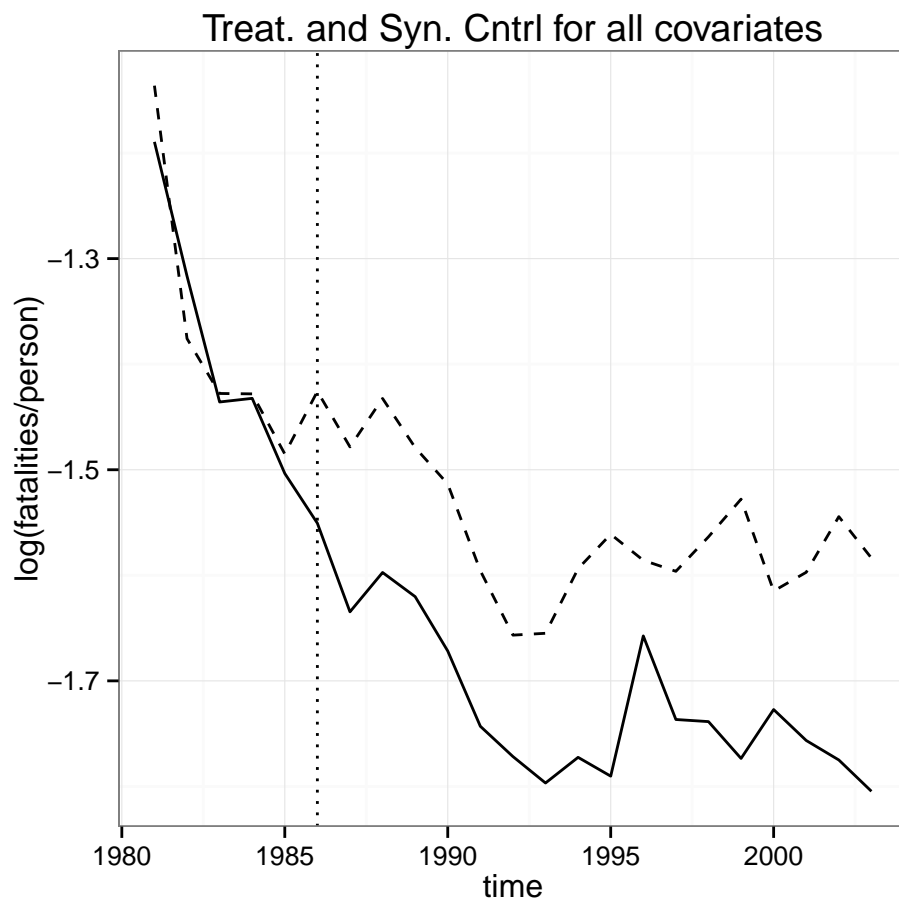
Figure 7: Treatment Unit and Synthetic Control log fatalities per capita developed using all previous time periods and all covariates.

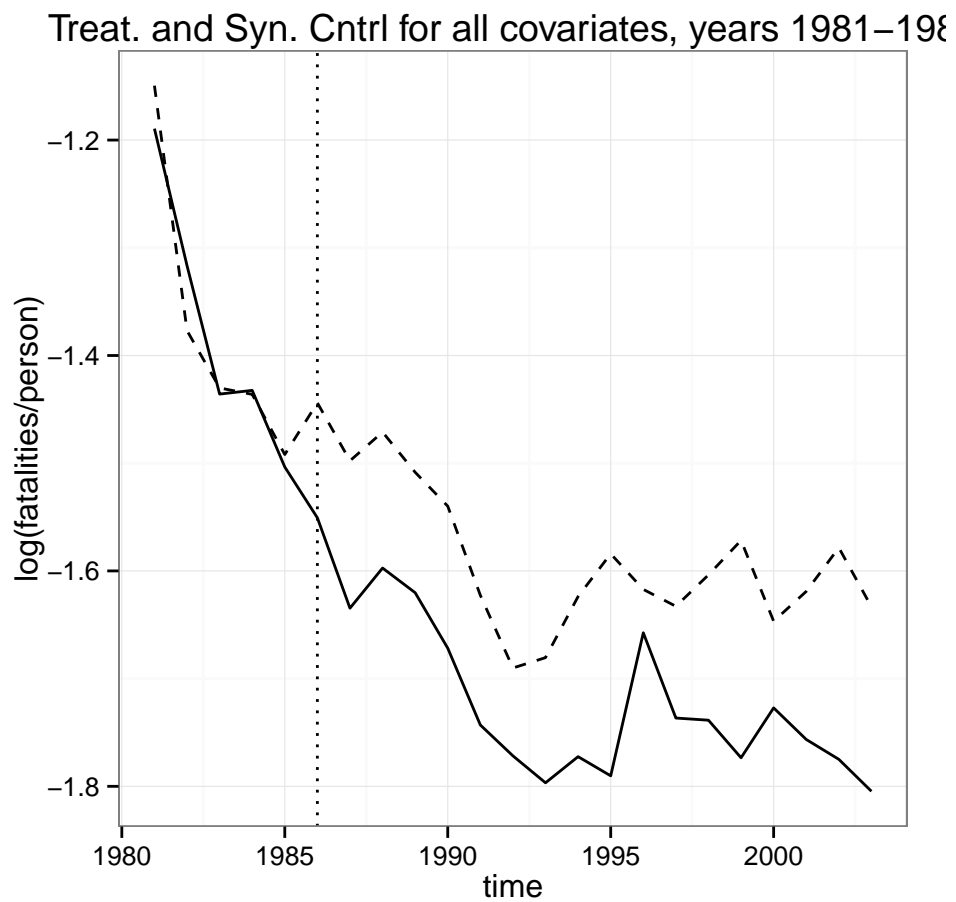Treat. and Syn. Cntrl for all covariates, years 1981–198

Figure 8: Treatment Unit and Synthetic Control log fatalities per capita developed using time periods 1981-1984 and all covariates.
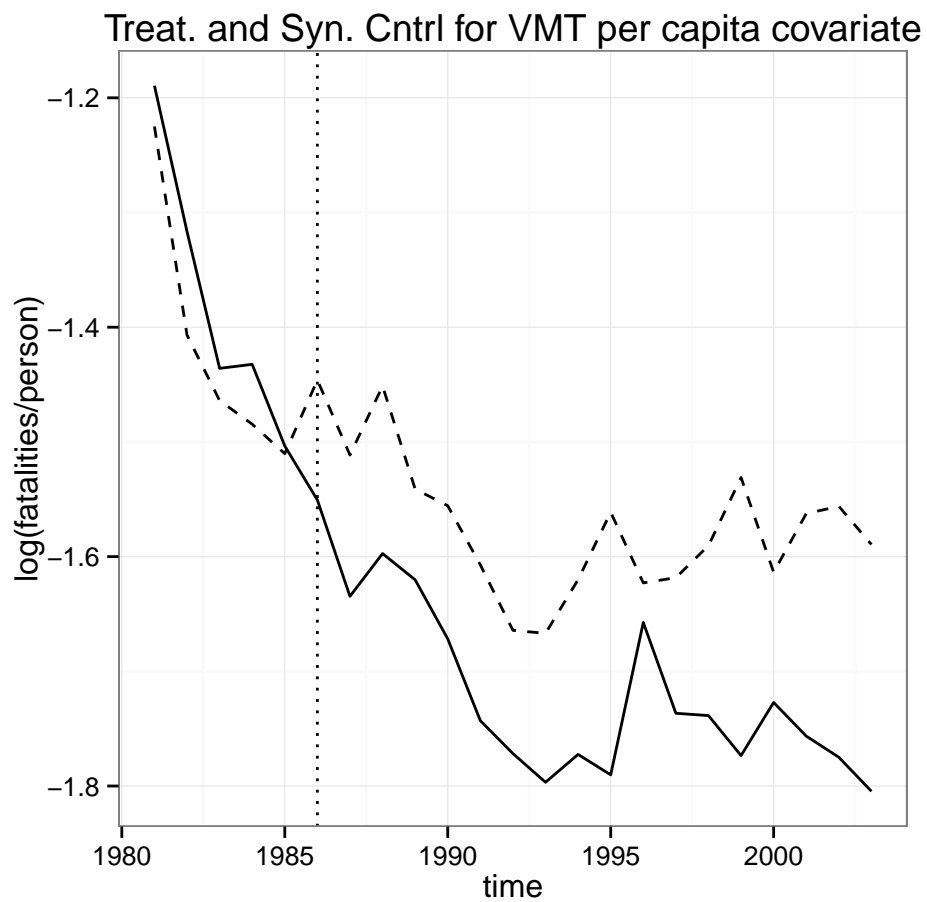
Figure 9: Treatment Unit and Synthetic Control log fatalities per capita developed using all previous time periods and VMT per capita as a covariate.

10x (Figure 12). In that formulation the treatment state is relatively consistently below nearly all the placebo control states until the end of the period. We discount the results towards the end of the treatment period because it is far beyond the support of the pre-treatment period. These results indicate that there is very likely a significant decrease in fatalities from primary seatbelt laws.
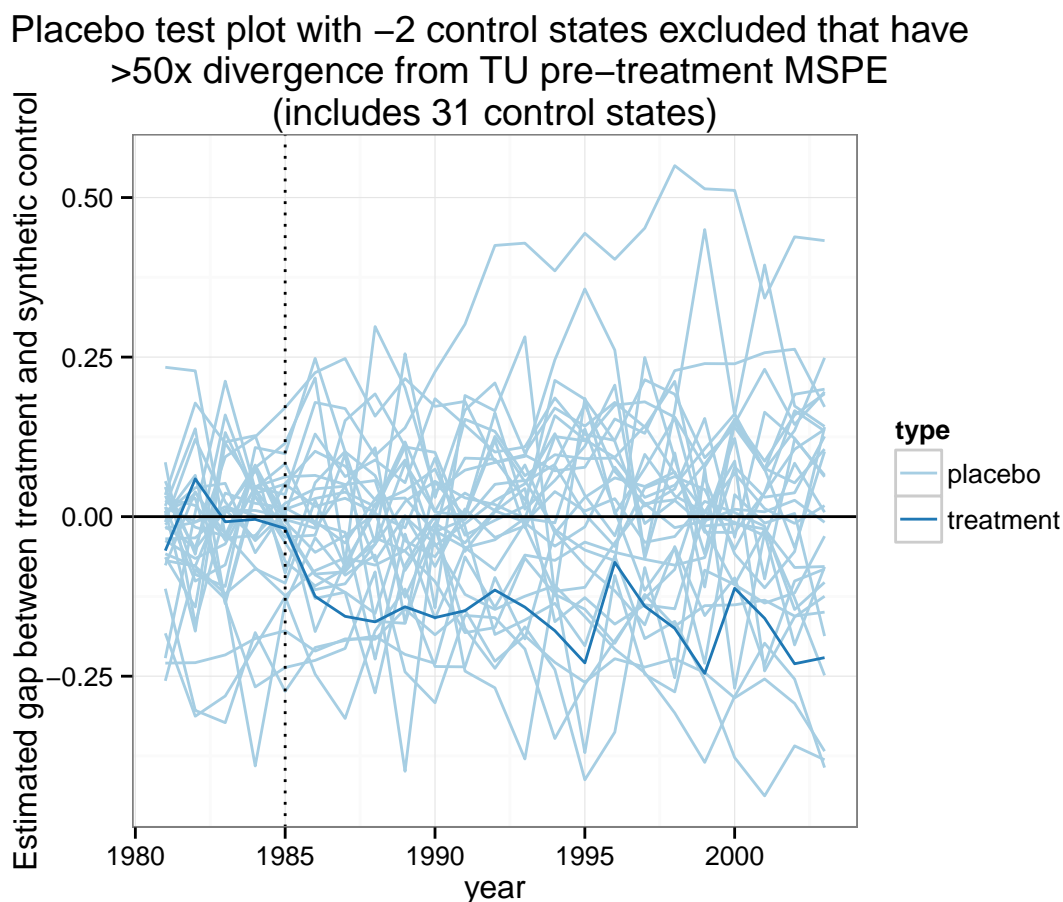


Figure 10: Placebo test results comparing treatment with placebo analysis on all control states. This plot does not exclude any states from the analysis.

Figure 11: Placebo test results comparing treatment with placebo analysis on all control states. This plot excludes control states with MSPE in the pre-treatment period greater than 20x the MSPE for the treatment state.

Placebo test plot with 5 control states excluded that have >10x divergence from TU pre−treatment MSPE (includes 24 control states)

Figure 12: Placebo test results comparing treatment with placebo analysis on all control states. This plot excludes control states with MSPE in the pre-treatment period greater than 10x the MSPE for the treatment state

Figure 13: Placebo test results comparing treatment with placebo analysis on all control states. This plot excludes control states with MSPE in the pre-treatment period greater than 5x the MSPE for the treatment state
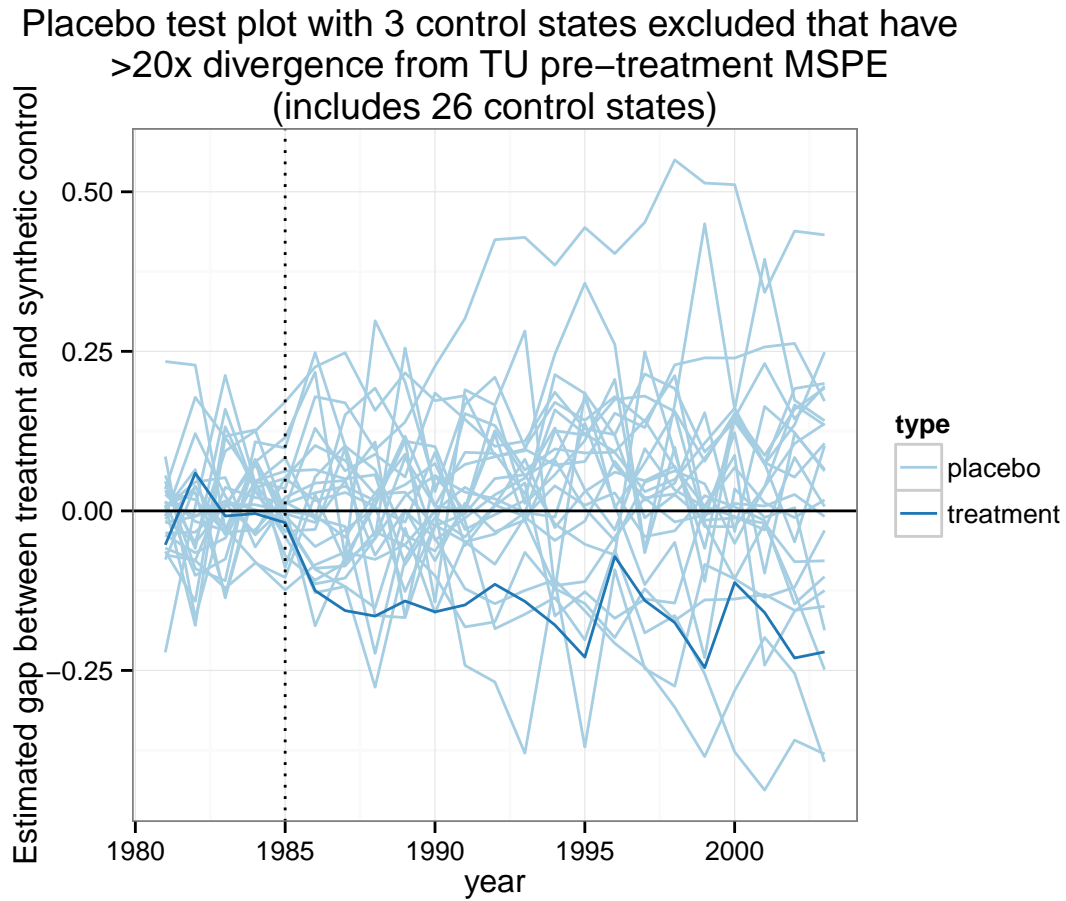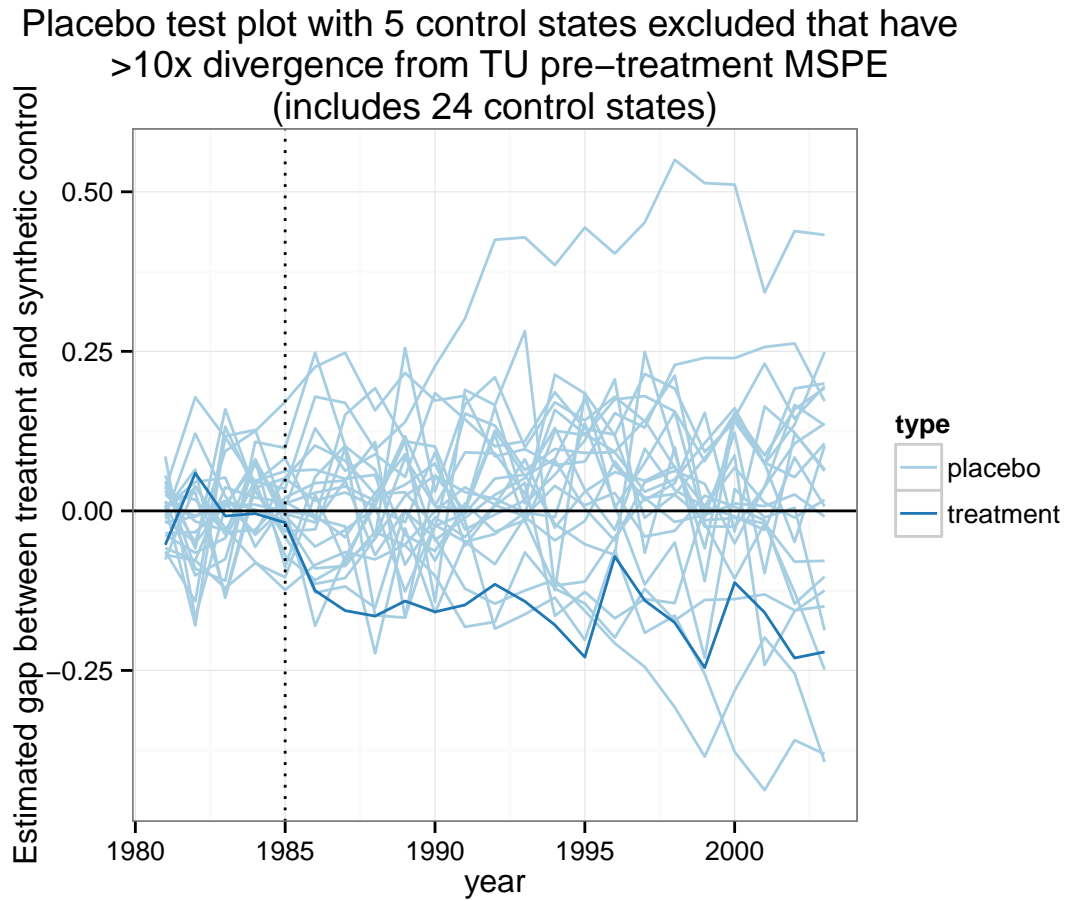
Figure 14: Placebo test results comparing treatment with placebo analysis on all control states. This plot excludes control states with MSPE in the pre-treatment period greater than 2x the MSPE for the treatment state

# (iii)   MSPE Ratios

**... Was it significant??**   As shown in Figure 15, we find that the MSPE ratio for the TU is about 20, which is the second highest among the combined set of the TU and the control states. Florida (as a placebo) had a higher MSPE ratio, about 120. This sows seeds of doubt but we still find that the TU has a higher apparent effect then 93% of units.

Figure 15: Ratio of post-treatment mean square percentage error (mspe) to pre-treatment mspe by state

## Part D:   Compare with FE Model

Our central estimate for the impact of seatbelt laws using synthetic control methods is a 15% reduction in fatality rate. This is roughly double the estimates we found using a fixed effects model (8% with significance at a 0.05 level) in the previous assignment. It is notable that many of the alternative estimates for the coefficient on primary seatbelt laws had higher (but not significant) results closer to those we find with the synthetic controls method.

The differences may stem from a reduced sample of treatment states in the synthetic control method. Since only four states with primary seatbelt laws were included (out of nearly 20 that eventually have the laws) it could be the case that the analysis is biased.
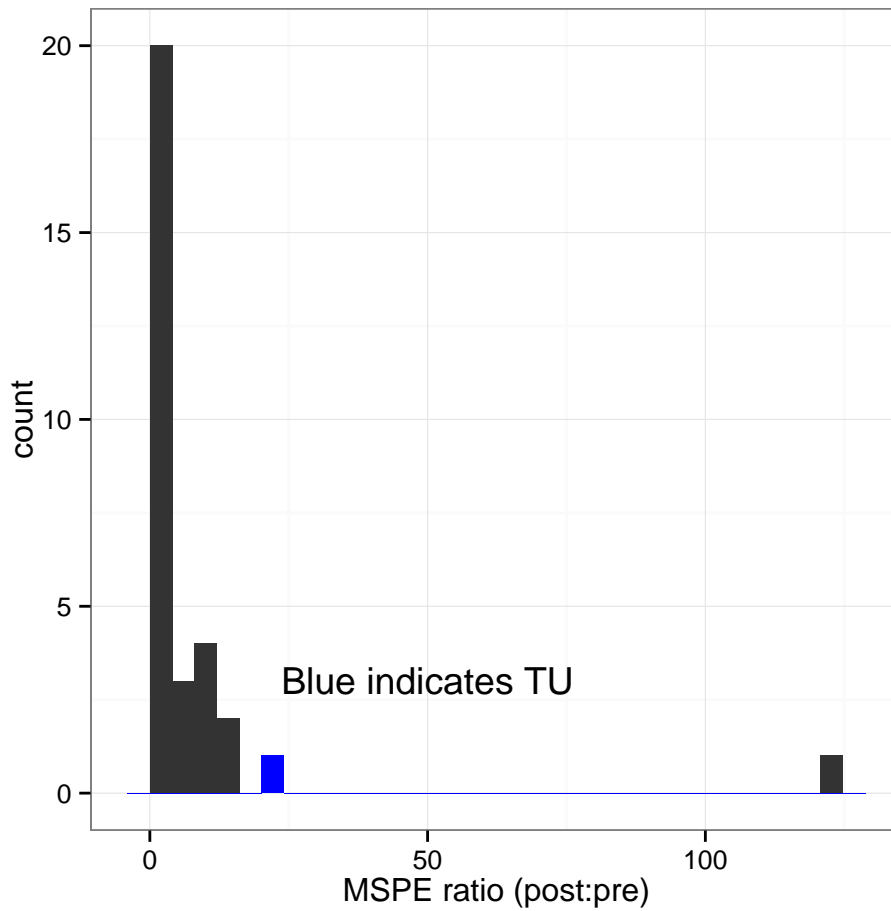
Another source of difference is in the structural way the estimates are made. Where the fixed effects model compares pre and post treatment, the synthetic controls method compares only in the post-treatment period. The fixed effects model includes first and second order corrections for year which should address some of these issues, but to the extent that the trend in national-level drivers for traffic fatalities does not obey the quadratic functional form there are potential errors that are not present in the synthetic control method, since synthetic controls does not impose a linear function in the same way on the covariates.

## Part E:   Appendix: Code Listings

```
 1 ## Frank's wd
 2 ## setwd("/media/frank/Data/documents/school/berkeley/fall13/are213/are213/
     ps2")
 3 ## Peter's wd
 4 # setwd("~/Google Drive/ERG/Classes/ARE213/are213/ps2")
 5
 6 library(foreign) #this is to read in Stata data
 7 library(Hmisc)
 8 library(psych)
 9 library(stargazer)
10 library(ggplot2) # for neato plotting tools
11 library(plyr) # for nice data tools like ddply
12 library(car) # "companion for applied regression" - recode fxn, etc.
13 library(gmodels) #for Crosstabs
14 library(plm) # for panel data
15 library(Synth) #for synthetic control
16 library(reshape)
17
18 source("../util/are213-func.R")
19
```

```
20  ps2a.data <- read.dta('traffic_safety2.dta')
21  ps2a.datakey <- data.frame(var.name=names(ps2a.data), var.labels = attr(ps2a
        .data, "var.labels"))
22  ps2a.data$logfatalpc <- with(ps2a.data, log(fatalities/population))
23  ps2a.data$sqyears <- with(ps2a.data, year^2)
24
25  # make state labels dataframe
26  state.labels <- attr(ps2a.data, "label.table")
27  state.labels$state.name <- attr(state.labels$state_number,"names")
28  state.labels <- as.data.frame(state.labels)
29  state.labels$state <- state.labels$state_number
30  state.labels$state.name <- as.character(state.labels$state.name)
31  state.labels$state_number <- NULL
32  state.labels <- rbind(state.labels, c("TU", 99))
33
34
35
36  tu.label <- data.frame(state_number = 99, state.name = "TU")
37
38
39  ## Problem 1 --------------------
40
41  ## Part a -------------
42  ## Subpart i ------------
43                                          # Average pre-period log fatalities
                                              per cap.
44  pre.treatment.avg <- with(subset(ps2a.data, (state == 99 & primary == 0)),
        mean(logfatalpc))
45  pre.control.avg <- with(subset(ps2a.data, (state != 99 & primary == 0)),
        mean(logfatalpc))
46
47                                          # Define state types
48  ps2a.data$type <- "Control"
49
50  for(i in 1:length(ps2a.data$year)){
51    if(ps2a.data$state[i] == 99){
52        ps2a.data$type[i] <- "Composite Treatment"}
53
54    if(ps2a.data$state[i] %in% c(4,10,30,41)){
55        ps2a.data$type[i] <- "Treatment State"}
56  }
57
58  use.rows <- which(ps2a.data$type != "Treatment State")
59
60  pre.period.data <- ddply(ps2a.data[use.rows,], .(year, type), summarize,
61                            logfatalpc = mean(logfatalpc))
62
63                                          # Plot of average fatality rates pre
                                              -1986.
64  preperiod.plot <- ggplot(data = subset(pre.period.data, year < 1986), aes(x
        = year, y = logfatalpc, color = type))
65  preperiod.plot <- preperiod.plot +
66      geom_line() +
67      theme_bw() +
68      ylab("log(fatalities per capita)")
69
```

```
70 ggsave(filename = "img-p2b.logfatTrend.pdf", plot = preperiod.plot, width =
       6, height = 4)
71
72 ## Subpart ii
73                                             # Compare logfatalpc in year before
                                                   treatment
74
75 year.before <- ddply(subset(ps2a.data, year == 1985 & type != "Treatment
       State"), .(state), summarize,
76                       type = type,
77                       logfatalpc = logfatalpc)
78
79 target.fatal <- year.before$logfatalpc[which(year.before$type== "Composite
       Treatment")]
80 year.before$distance <- abs(year.before$logfatalpc - target.fatal)
81 year.before$distance[which(year.before$type== "Composite Treatment")] <- NA
82
83
84                                             # plot all states and TU
85 pdf("img-p2b-compareStates.pdf", width = 4, height = 4)
86
87 ggplot(year.before, aes(logfatalpc)) +
88     geom_histogram(binwidth=0.1) +
89     geom_histogram(data=subset(year.before, state==1), aes(logfatalpc),
           binwidth=0.1, fill="red") +
90     theme_bw() +
91     xlab("log(fatalities per capita)") +
92     geom_vline(aes(xintercept=target.fatal), size = 1, color = "blue")
93
94 dev.off()
95
96
97 best.yb.match <- which(year.before$distance == min(year.before$distance, na.
       rm=T))
98
99 print(paste("The state number for the closest year before match is", year.
       before$state[best.yb.match]))
100
101                                             # The best match is Alabama.
102
103                                             # Tables comparing Alabama to the
                                                   composite treatment group
104 stargazer(subset(ps2a.data, state==99),
105           out = "tab-ps2b-1a.tex",
106           title = "Composite Treatment Group Summary",
107           label = "tab:a21")
108
109 stargazer(subset(ps2a.data, state==1),
110           out = "tab-ps2b-1b.tex",
111           title = "Closest match for pre-policy fatalities: Alabama",
112           label = "tab:a22")
113
114
115 # graphical comparison between states
116
117 prep.gg.compare <- melt.data.frame(subset(ps2a.data, state==99 | state ==1),
       id.vars = c("state","year"))
```

```
118
119 prep.gg.compare$value <- as.numeric(prep.gg.compare$value)
120 prep.gg.compare$state <- as.factor(prep.gg.compare$state)
121
122 vars.we.care.about <- c("beer", "college", "primary", "secondary", "unemploy
        ", "logfatalpc", "rural_speed", "urban_speed", "precip")
123 rows.we.care.about <- which(prep.gg.compare$variable %in% vars.we.care.about
        )
124
125 gg.compare.states <- ggplot(prep.gg.compare[rows.we.care.about,], aes(x=year
        , y=value, color=state)) +
126     geom_point() +
127     facet_wrap("variable", scales = "free")
128
129 pdf("img-ps2b-compareStatesFacets.pdf", width=7, height = 6)
130 gg.compare.states
131 dev.off()
132
133 ## Part B - Synthetic control method
134
135
136                                             # Identify which states can be "
                                                control" (i.e., those that never
                                                have a primary seatbelt law)
137
138 state.policy <- ddply(ps2a.data, .(state), summarize, primary.max = max(
        primary))
139 good.states <- which(state.policy$primary.max < 1)
140 control.states <- state.policy$state[good.states]
141 all.states <- c(control.states, 99)
142
143 ## subpart ii
144
145
146 # add to dataset
147 ps2a.data$vmt_percapita <- with(ps2a.data, totalvmt / population)
148 ps2a.data <- join(ps2a.data, state.labels)
149
150 #helper function: extract useful data from synth results
151 get.plot.data <- function(synth.res, dataprep.res, label="unlabeled", type =
        "unknown"){
152     out <- data.frame(time=dataprep.res$tag$time.plot)
153
154     synthetic.trend <- dataprep.res$Y0plot %*% synth.res$solution.w
155     treatment.trend <- dataprep.res$Y1plot
156     gap <- treatment.trend - synthetic.trend
157
158     out$synth <- as.numeric(synthetic.trend)
159     out$treat <- as.numeric(treatment.trend)
160     out$gap <- as.numeric(gap)
161     out$spe <- out$gap^2
162     out$label <- label
163     out$type <- type
164
165
166     return(out)
167 }
```

```
168
169 #helper function: run dataprep, synth, and get.plot.data for a set of inputs
       .
170 run.syn <- function(predictor.set, time.prior, treatment.state = std.
       treatment.state, control.states, label = "unlabeled", type = "unknown"){
171
172 # synthetic controls dataprep with specified set of tractable predictors
173     dataprep.results <- dataprep(foo = ps2a.data,
174                                  predictors = predictor.set,
175                                  predictors.op = c("mean"),
176                                  dependent = "logfatalpc",
177                                  unit.variable = "state",
178                                  time.variable = "year",
179                                  treatment.identifier = treatment.state,
180                                  controls.identifier = control.states,
181                                  time.predictors.prior = time.prior,
182                                  time.optimize.ssr = time.prior,
183                                  time.plot = c(1981:2003),
184                                  unit.names.variable = "state.name"
185                                  )
186
187     synth.results <- synth(data.prep.obj = dataprep.results)
188
189     output <- get.plot.data(synth.results, dataprep.results, label, type)
190
191     return(output)
192 }
193
194 ## helper function: makes a gap plot so we can get them to look consistent
195 make.gap.plot <- function(syntheticresults, finame, desc){
196     gap.plot <- ggplot(data = syntheticresults, aes(y = gap, x = time))
197     gap.plot <- gap.plot +
198         geom_path() +
199             geom_vline(xintercept = 1986, linetype = 'dotted') +
200                 theme_bw() + geom_hline(yintercept = 0) + ylab("difference
                        in log(fatalities/person)") +
201                     labs(title = paste("Gap between TU and synth. control
                            for", desc))
202     ggsave(filename = paste0('img-gap-', finame, '.pdf'), plot = gap.plot,
           width=5, height = 5)
203
204     split.plot <- ggplot(data = syntheticresults, aes(y = treat, x = time))
205     split.plot <- split.plot +
206         geom_path() +
207             geom_path(aes(y = synth), linetype = 'dashed') +
208                 geom_vline(xintercept = 1986, linetype = 'dotted') +
209                     theme_bw() + ylab('log(fatalities/person)') +
210                         labs(title = paste("Treat. and Syn. Cntrl for",
                                desc))
211     ggsave(filename = paste0('img-split-', finame, '.pdf'), plot = split.
           plot, width=5, height = 5)
212
213     return(gap.plot)
214 }
215
216
```

```
217                                          # default entries to run.syn
                                             function
218 std.treatment.state <- 99
219 full.predictor.set <- c("college", "precip", "snow32", "beer", "vmt_
        percapita", "unemploy")
220 full.time.prior <- c(1981:1985)
221
222 ## Part C: Pretty pictures
223 ## Part (i)
224 full.synthesis <- run.syn(full.predictor.set, full.time.prior, 99, control.
        states)
225 full.gap.plot <- make.gap.plot(full.synthesis, 'full', 'all covariates')
226
227 full.synthesis1984 <- run.syn(full.predictor.set, c(1981:1984), 99, control.
        states)
228 full.1984.gap.plot <- make.gap.plot(full.synthesis1984, 'full1984', 'all
        covariates, years 1981-1984')
229
230 vmt.synthesis <- run.syn("vmt_percapita", full.time.prior, 99, control.
        states)
231 vmt.gap.plot <- make.gap.plot(vmt.synthesis, 'vmt', 'VMT per capita
        covariate')
232
233 # Script to generate a placebo test plot
234 ## Part (ii)
235
236 placebo.test <- run.syn(full.predictor.set, full.time.prior, treatment.state
        = 99, control.states, label = 99, type = "treatment")
237
238 for(state in control.states){
239     updated.control <- control.states[which(control.states!=state)]
240     placebo.additional.result <- run.syn(full.predictor.set, full.time.prior
            , state, updated.control, label = state, type = "placebo")
241     placebo.test <- rbind(placebo.test, placebo.additional.result)
242     print(paste("Finished with state number", state, "and you should be
            patient for the rest to finish :)"))
243 }
244
245 placebo.test$treated <- (placebo.test$time > 1985)
246
247 MPSEs <- ddply(placebo.test, .(label), summarize,
248                 preMPSE = mean(spe[!treated]),
249                 postMPSE = mean(spe[treated]))
250 MPSEs$ratio <- with(MPSEs, postMPSE/preMPSE)
251
252
253 # Exclusion threshold for placebo plot based on
254 preThreshold <- 50 #exclude states with pre-intervention MSPE greater than x
            times that for TU
255 tu.mspe <- MPSEs$preMPSE[which(MPSEs$label==99)]
256 include.states <- MPSEs$label[which(MPSEs$preMPSE < tu.mspe * preThreshold)]
257 include.states.rows <- which(placebo.test$label %in% include.states)
258
259 ex.number <- length(control.states) - length(include.states) -1
260 in.number <- length(include.states)
261
262 # plot all the lines
```

```
263  placeboTest <- ggplot(placebo.test[include.states.rows,], aes(time, gap,
         group=label))
264  placeboTest <- placeboTest +
265      geom_line(aes(color=type))+
266      geom_vline(aes(xintercept = 1985), linetype="dotted") +
267      theme_bw() +
268      xlab("year") +
269      ylab("Estimated gap between treatment and synthetic control") +
270      scale_color_brewer(palette = "Paired") +
271      geom_hline(yintercept = 0) +
272      ggtitle(paste0("Placebo test plot with ",ex.number," control states
             excluded that have \n >", preThreshold, "x divergence from TU pre-
             treatment MSPE\n(includes ", in.number, " control states)"))
273  ggsave(paste0("img-placeboTest", preThreshold, ".pdf"), plot = placeboTest,
         width = 6, height = 5)
274
275
276  ## Part(iii)
277
278
279  MPSE.plot <- ggplot(MPSEs, aes(ratio))
280  MPSE.plot <-  MPSE.plot +
281    geom_histogram() +
282    geom_histogram(data=subset(MPSEs, label==99), aes(ratio), fill="blue") +
283    theme_bw() +
284    annotate(geom="text", x=50, y=3, label="Blue indicates TU") +
285    xlab("MSPE ratio (post:pre)")
286
287  ggsave(filename = "img-mspeRatio.pdf", plot = MPSE.plot, width = 5, height =
         5)
288
289  # # Visual exploration of plots....only works if you have the native
         dataprep and synth objects.  We would rather make our own plots :)
290  # gaps.plot(seatbelts.synth, syn.data.full)
291  # path.plot(seatbelts.synth, syn.data.full, Ylim = c(-1.3, -2))
292
293  # # Tables for synthetic controls results
294  # syn.table <- synth.tab(seatbelts.synth, syn.data.full,3)
```

./ps2b.r

```
 1  # Econometrics helper functions for [R]
 2  #
 3  # Peter Alstone and Frank Proulx
 4  # 2013
 5  # version 1
 6  # contact: peter.alstone AT gmail.com
 7
 8  # Category: Data Management -------------
 9
10
11  # Category: Data Analysis ----------------
12
13  # Function: Find adjusted R^2 for subset of data
14  # This requires a completed linear model...pull out the relevant y-values
         and residuals and feed them to function
```

```
15 # [TODO @Peter] Improve function so it can simply evaluate lm or glm object,
       add error handling, general clean up.
16 adjr2 <- function(y,resid){
17   r2 <- 1-sum(resid^2) / sum((y-mean(y))^2)
18   return(r2)
19 } #end adjr2
20
21
22 # Category: Plots and Graphics -------------
23
24 ## Function for arranging ggplots. use png(); arrange(p1, p2, ncol=1); dev.
       off() to save.
25 require(grid)
26 vp.layout <- function(x, y) viewport(layout.pos.row=x, layout.pos.col=y)
27 arrange_ggplot2 <- function(..., nrow=NULL, ncol=NULL, as.table=FALSE) {
28   dots <- list(...)
29   n <- length(dots)
30   if(is.null(nrow) & is.null(ncol)) { nrow = floor(n/2) ; ncol = ceiling(n/
          nrow)}
31   if(is.null(nrow)) { nrow = ceiling(n/ncol)}
32   if(is.null(ncol)) { ncol = ceiling(n/nrow)}
33   ## NOTE see n2mfrow in grDevices for possible alternative
34   grid.newpage()
35   pushViewport(viewport(layout=grid.layout(nrow,ncol) ) )
36   ii.p <- 1
37   for(ii.row in seq(1, nrow)){
38     ii.table.row <- ii.row
39     if(as.table) {ii.table.row <- nrow - ii.table.row + 1}
40     for(ii.col in seq(1, ncol)){
41       ii.table <- ii.p
42       if(ii.p > n) break
43       print(dots[[ii.table]], vp=vp.layout(ii.table.row, ii.col))
44       ii.p <- ii.p + 1
45     }
46   }
47 }
48
49 robust <- function(model){ #This calculates the Huber-White Robust standard
       errors -- code from http://thetarzan.wordpress.com/2011/05/28/
       heteroskedasticity-robust-and-clustered-standard-errors-in-r/
50     s <- summary(model)
51     X <- model.matrix(model)
52     u2 <- residuals(model)^2
53     XDX <- 0
54
55     for(i in 1:nrow(X)) {
56         XDX <- XDX +u2[i]*X[i,]%*%t(X[i,])
57     }
58
59 # inverse(X'X)
60     XX1 <- solve(t(X)%*%X)
61
62 #Compute variance/covariance matrix
63     varcovar <- XX1 %*% XDX %*% XX1
64
65 # Degrees of freedom adjustment
66     dfc <- sqrt(nrow(X))/sqrt(nrow(X)-ncol(X))
```

```
67
68      stdh <- dfc*sqrt(diag(varcovar))
69
70      t <- model$coefficients/stdh
71      p <- 2*pnorm(-abs(t))
72      results <- cbind(model$coefficients, stdh, t, p)
73      dimnames(results) <- dimnames(s$coefficients)
74      results
75  }
76
77  ## Two functions for clustered standard errors below from: http://people.su.
        se/~ma/clustering.pdf -------
78
79  clx <-
80    function(fm, dfcw, cluster){
81      # R-codes (www.r-project.org) for computing
82      # clustered-standard errors. Mahmood Arai, Jan 26, 2008.
83
84      # The arguments of the function are:
85      # fitted model, cluster1 and cluster2
86      # You need to install libraries 'sandwich' and 'lmtest'
87
88      # reweighting the var-cov matrix for the within model
89      library(sandwich);library(lmtest)
90      M <- length(unique(cluster))
91      N <- length(cluster)
92      K <- fm$rank
93      dfc <- (M/(M-1))*((N-1)/(N-K))
94      uj  <- apply(estfun(fm),2, function(x) tapply(x, cluster, sum));
95      vcovCL <- dfc*sandwich(fm, meat=crossprod(uj)/N)*dfcw
96      coeftest(fm, vcovCL) }
97
98  mclx <-
99    function(fm, dfcw, cluster1, cluster2){
100     # R-codes (www.r-project.org) for computing multi-way
101     # clustered-standard errors. Mahmood Arai, Jan 26, 2008.
102     # See: Thompson (2006), Cameron, Gelbach and Miller (2006)
103     # and Petersen (2006).
104     # reweighting the var-cov matrix for the within model
105
106     # The arguments of the function are:
107     # fitted model, cluster1 and cluster2
108     # You need to install libraries 'sandwich' and 'lmtest'
109
110     library(sandwich);library(lmtest)
111     cluster12 = paste(cluster1,cluster2, sep="")
112     M1  <- length(unique(cluster1))
113     M2  <- length(unique(cluster2))
114     M12 <- length(unique(cluster12))
115     N   <- length(cluster1)
116     K   <- fm$rank
117     dfc1  <- (M1/(M1-1))*((N-1)/(N-K))
118     dfc2  <- (M2/(M2-1))*((N-1)/(N-K))
119     dfc12 <- (M12/(M12-1))*((N-1)/(N-K))
120     u1j  <- apply(estfun(fm), 2, function(x) tapply(x, cluster1,  sum))
121     u2j  <- apply(estfun(fm), 2, function(x) tapply(x, cluster2,  sum))
122     u12j <- apply(estfun(fm), 2, function(x) tapply(x, cluster12, sum))
```

```
123    vc1    <-   dfc1*sandwich(fm, meat=crossprod(u1j)/N )
124    vc2    <-   dfc2*sandwich(fm, meat=crossprod(u2j)/N )
125    vc12   <-  dfc12*sandwich(fm, meat=crossprod(u12j)/N)
126    vcovMCL <- (vc1 + vc2 - vc12)*dfcw
127    coeftest(fm, vcovMCL)}
128
129 ## Function to compute ols standard errors , robust, clustered...
130 ## Based on http://diffuseprior.wordpress.com/2012/06/15/standard-robust-and
       -clustered-standard-errors-computed-in-r/
131 ols.hetero <- function(form, data, robust=FALSE, cluster=NULL,digits=3){
132   r1 <- lm(form, data)
133   if(length(cluster)!=0){
134     data <- na.omit(data[,c(colnames(r1$model),cluster)])
135     r1 <- lm(form, data)
136   }
137   X <- model.matrix(r1)
138   n <- dim(X)[1]
139   k <- dim(X)[2]
140   if(robust==FALSE & length(cluster)==0){
141     se <- sqrt(diag(solve(crossprod(X)) * as.numeric(crossprod(resid(r1))/(n
           -k))))
142     res <- cbind(coef(r1),se)
143   }
144   if(robust==TRUE){
145     u <- matrix(resid(r1))
146     meat1 <- t(X) %*% diag(diag(crossprod(t(u)))) %*% X
147     dfc <- n/(n-k)
148     se <- sqrt(dfc*diag(solve(crossprod(X)) %*% meat1 %*% solve(crossprod(X)
           )))
149     res <- cbind(coef(r1),se)
150   }
151   if(length(cluster)!=0){
152     clus <- cbind(X,data[,cluster],resid(r1))
153     colnames(clus)[(dim(clus)[2]-1):dim(clus)[2]] <- c(cluster,"resid")
154     m <- dim(table(clus[,cluster]))
155     dfc <- (m/(m-1))*((n-1)/(n-k))
156     uclust  <- apply(resid(r1)*X,2, function(x) tapply(x, clus[,cluster],
           sum))
157     se <- sqrt(diag(solve(crossprod(X)) %*% (t(uclust) %*% uclust) %*% solve
           (crossprod(X)))*dfc)
158     res <- cbind(coef(r1),se)
159   }
160   res <- cbind(res,res[,1]/res[,2],(1-pnorm(abs(res[,1]/res[,2])))*2)
161   res1 <- matrix(as.numeric(sprintf(paste("%.",paste(digits,"f",sep=""),sep=
         ""),res)),nrow=dim(res)[1])
162   rownames(res1) <- rownames(res)
163   colnames(res1) <- c("Estimate","Std. Error","t value","Pr(>|t|)")
164   return(res1)
165 }
```

../util/are213–func.R