

# ARE 213 Problem Set 1B

Becky Cardinali, Yuen Ho, Sara Johns, and Jacob Leffler

Due 10/12/2020

## Question 1

In Problem Set 1a, you used linear regression to relate infant health outcomes and maternal smoking during pregnancy. Please answer the following questions.

- (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?

## Question 2

Describe the propensity score approach to the problem of estimating the average causal effect of smoking when the treatment is randomly assigned conditional on the observables. How does it reduce the dimensionality problem of multivariate matching?

We know that if we condition on observables, we will get a consistent estimate of the ATE under the assumption. However, if the observables are high dimensional, it might be difficult to find a comparison unit with the same values of the observables. From lecture, we know that it is sufficient instead to condition on the propensity score. Using the propensity score allows us to compare treated and control units with the same probability of being treated. The propensity score does not require that all values of the observables be the same and so therefore avoids problems of multidimensionality.

Try a few ways to estimate the effects of maternal smoking on birthweight:

### 2a

- a) First create the propensity score. For our purposes let's use a logit specification. First specify the logit using all of the “predetermined” covariates (don't include interactions). Next, include only those “predetermined” covariates that enter significantly in the first logit specification. How comparable are the propensity scores? If they are similar does this imply that we have the “correct” set of covariates in the logit specification used for our propensity score?

```
# get prop score using all predetermined variables
prop_all <- glm(tobacco ~ factor(stresfip) + dimage + factor(mrace3) + dmeduc +
               dtotord + disllb + dfage + factor(birmon) + factor(orfath) +
               factor(dmar) + dfeduc + dplural + factor(pre4000) + factor(preterm),
               family=binomial(link='logit'), data = mom_dt)
mom_dt[, prop_score_all := fitted(prop_all)]

# take a look at the output to see which are significant - omitting because takes up a lot of space
# summary(prop_all)
# only need to take out state of residence. birth month has a few months that are significant so will keep

# get prop score using significant variables from previous logit
```

```
prop_sig <- glm(tobacco ~ dimage + factor(mrace3) + dmeduc +
               dtotord + disllb + dfage + factor(birmon) + factor(orfath) +
               factor(dmar) + dfeduc + dplural + factor(pre4000) + factor(preterm),
               family=binomial(link='logit'), data = mom_dt)
mom_dt[, prop_score_sig := fitted(prop_sig)]

# how different are the prop scores?
prop_score_diff <- mom_dt$prop_score_all - mom_dt$prop_score_sig
summary(prop_score_diff)
```

```
##           Min.      1st Qu.      Median      Mean      3rd Qu.      Max.
## -0.2889930  0.0003105  0.0004940  0.0000000  0.0006952  0.2873649
```

*# a few outliers but not very different*

## 2b

Control directly for the estimated propensity scores using a regression analysis, and estimate an average treatment effect. State clearly the assumptions under which your estimate is correct.

```
# run regression with prop_score
prop_reg <- lm(dbrwt ~ tobacco + prop_score_sig, data = mom_dt)
```

Controlling directly for the propensity score, we get the ATE is -223.23 and is statistically significant at the 99.9% level.

## 2c

```
# create propensity weights
mom_dt[,prop_weights := ifelse(tobacco == 1,
                               1/prop_score_sig,
                               1/(1 - prop_score_sig))]

# normalize the weights
mom_dt[,norm_prop_weights := ifelse(tobacco == 1,
                                     prop_weights/sum(mom_dt[tobacco == 1, prop_weights]),
                                     prop_weights/sum(mom_dt[tobacco == 0, prop_weights]))]

# estimate ATE
tau_ipw <- sum((mom_dt$tobacco*mom_dt$dbrwt)*(mom_dt$norm_prop_weights) - ((1 - mom_dt$tobacco)*mom_dt$dbrwt))

# estimate TOT - use formula from section
tot_y1 <- sum(mom_dt$tobacco*mom_dt$dbrwt) / sum(mom_dt$tobacco)
tot_y0 <- sum(mom_dt$prop_score_sig * (1 - mom_dt$tobacco) * mom_dt$dbrwt /
              (1 - mom_dt$prop_score_sig)) /
         sum(mom_dt$prop_score_sig * (1 - mom_dt$tobacco) /
              (1 - mom_dt$prop_score_sig))

tau_tot <- tot_y1 - tot_y0
```

Using inverse propensity score weighting, we get that the ATE is now -225.29 and the TOT is -224.4.

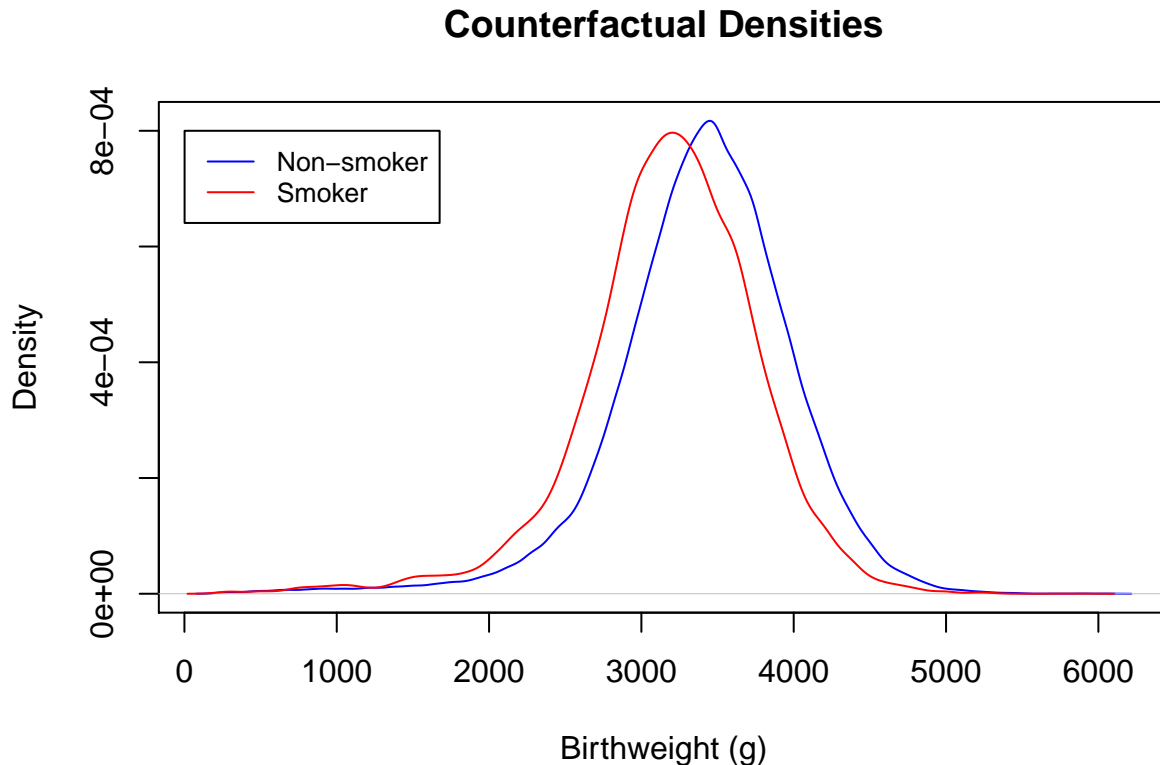
## 2d

Estimate the counterfactual densities relevant for the above part with a kernel density estimator. That is, estimate the density of birthweight (or log birthweight) if everyone smoked and again if no one smoked. Hint: Consider directly

applying the Hirano, Imbens, and Ridder propensity score reweighting scheme in the context of estimating the densities of the treated and control groups (rather than the means of the treated and control groups). Stata has very useful preprogrammed commands. In addition to using the preprogrammed Stata command to compute/graph the kernel density over the entire range of birthweight, please also calculate by hand the kernel estimator at birthweight equals 3,000 grams (and provide the code you wrote that shows the calculation of the kernel estimator at this single point). Play around with a bandwidth starting with half the default Stata bandwidth. Choose the same bandwidth for all the pictures, and produce a (beautiful, production quality) figure depicting both densities.

```
d0 <- density(mom_dt[tobacco==0,dbrwt], kernel="gaussian", bw="sj",
              adjust=1, weights=mom_dt[tobacco==0,norm_prop_weights])
d1 <- density(mom_dt[tobacco==1,dbrwt], kernel="gaussian", bw="sj",
              adjust=1, weights=mom_dt[tobacco==1,norm_prop_weights])

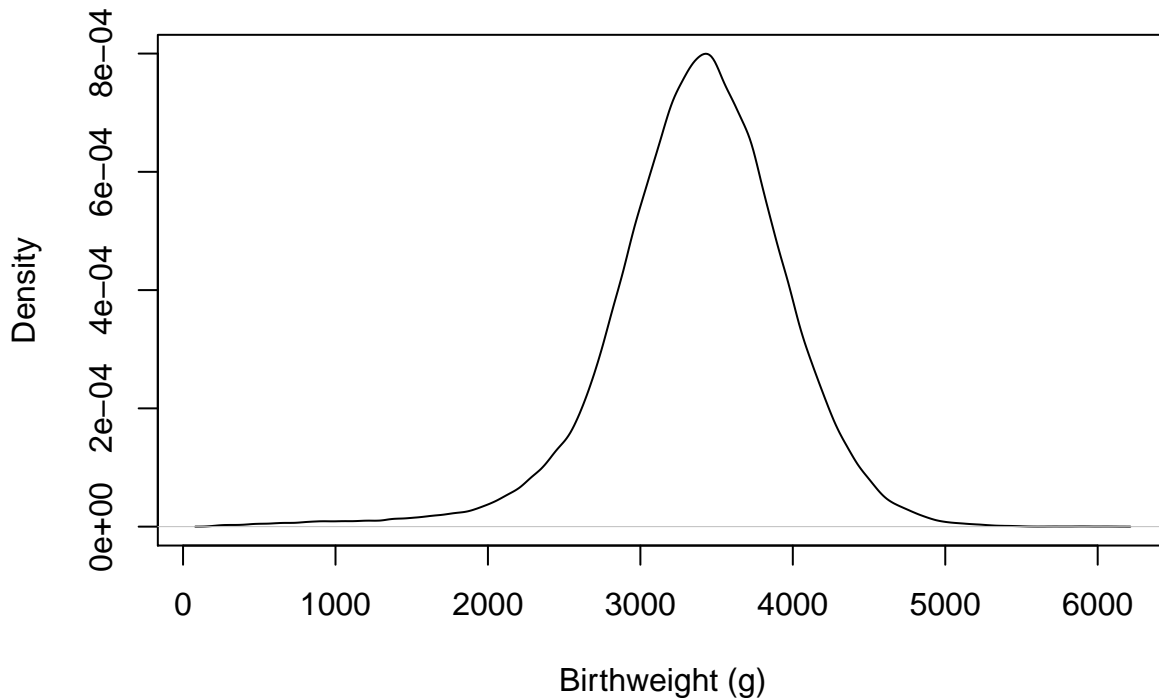
plot(d0, col="blue", main="Counterfactual Densities",
     xlab = "Birthweight (g)")
lines(d1, col="red")
legend(1, 0.0008, legend=c("Non-smoker", "Smoker"),
      col=c("blue","red"), lty=1, cex=0.8)
```



As we expect from our estimates of the ATE, the density of the birthweights for the treated group is shifted to the left of the control group (lower birthweights).

```
d_all <- density(mom_dt$dbrwt, kernel="gaussian", bw="sj", adjust=1)
plot(d_all, col="black", main="Kernel Density, Full range of birthweights",
     xlab = "Birthweight (g)")
```

## Kernel Density, Full range of birthweights



```
# kernel estimator at 3000 g
# use bandwidth selected above: h = 49
h <- 49
ke_3000 <- 1 / (nrow(mom_dt) * h * sqrt(2*pi)) * sum(exp(-0.5 * (((3000-mom_dt$dbrwt)/h)^2)))
```

We plot the kernel density over the entire range of birthweights, using the Gaussian kernel and selecting the bandwidth using the methods of Sheather & Jones (1991) which uses pilot estimation of derivatives. For the entire range of birthweights, this gives us a bandwidth of 49. We then calculate the kernel estimator by hand at a weight of 3000 grams using the Gaussian kernel and this bandwidth. Our formula is

$$\hat{f}(x) = \frac{1}{nh} \frac{1}{\sqrt{2\pi}} \sum_i^n e^{-1/2(\frac{x-x_i}{h})^2}$$

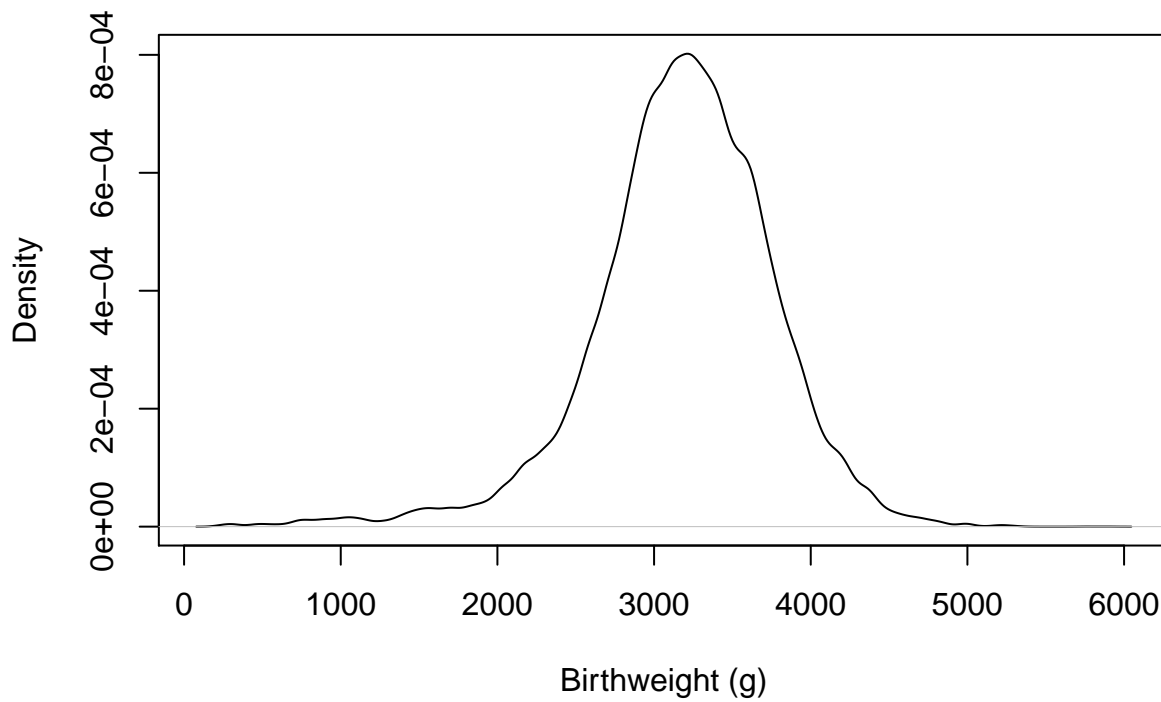
where  $x = 3000$  and  $h = 49$ . We get that the density at 3000g is  $5.4 \times 10^{-4}$ , which visually corresponds to our plot of the density.

## 2e

Take one of your densities and display an estimate of the density using different bandwidths as well as the one you settled on. What happens with bigger (smaller) bandwidths?

```
# bw of treated group is 68.9, what if we used 49.8 like control?
d1_low <- density(mom_dt[tobacco==1,dbrwt], kernel="gaussian", bw=49.8, adjust=1, weights=mom_dt[tobacco==1,n
plot(d1_low, col="black", main="Kernel Density, Smoker: Lower Bandwidth",
      xlab = "Birthweight (g)")
```

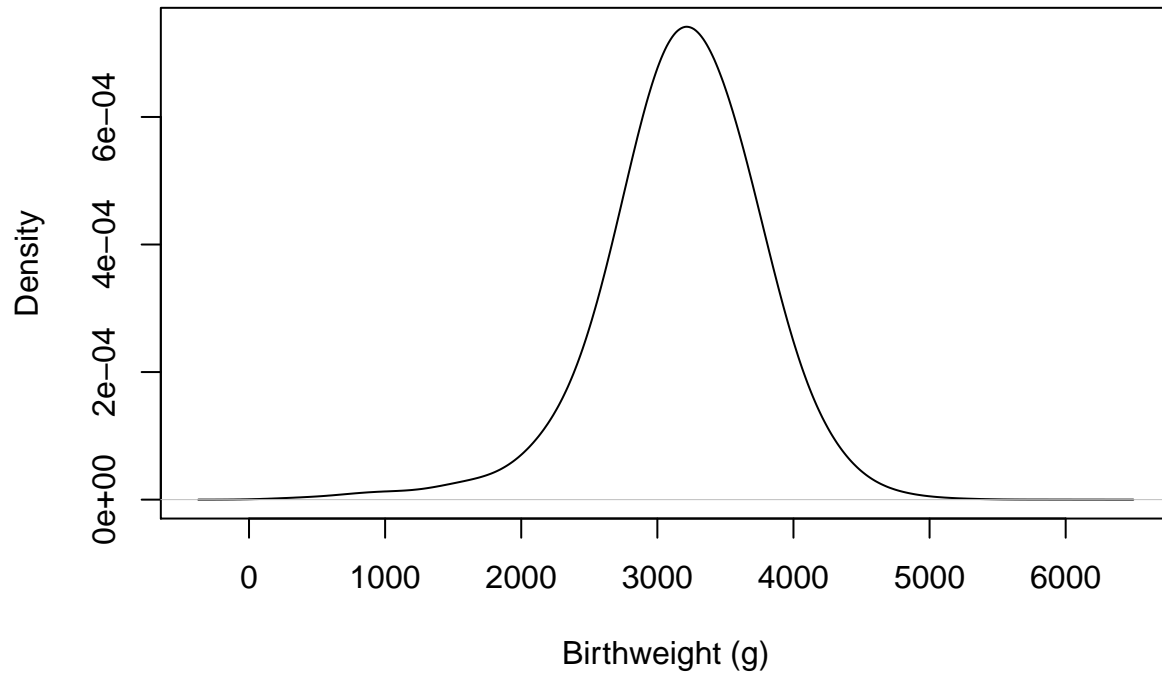
## Kernel Density, Smoker: Lower Bandwidth



```
# what if we raised it?
```

```
d1_high <- density(mom_dt[tobacco==1,dbrwt], kernel="gaussian", bw=200, adjust=1, weights=mom_dt[tobacco==1,n  
plot(d1_high, col="black", main="Kernel Density, Smoker: Higher Bandwidth",  
      xlab = "Birthweight (g)")
```

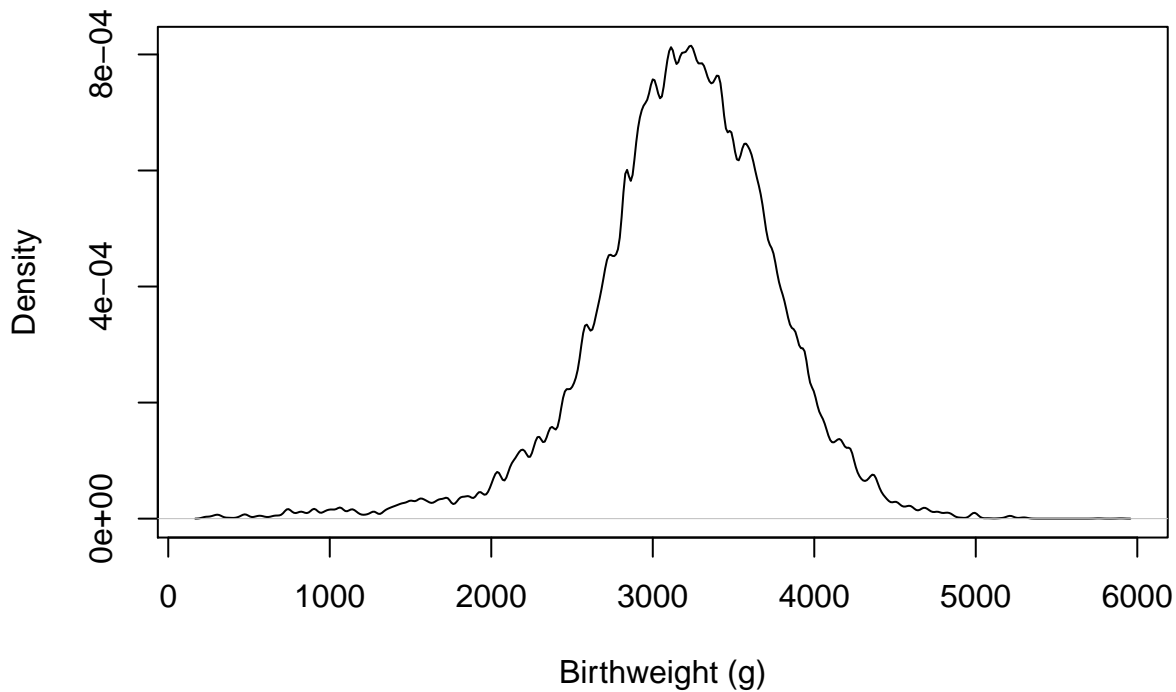
## Kernel Density, Smoker: Higher Bandwidth



```
# what if we made it very low
```

```
d1_verylow <- density(mom_dt[tobacco==1,dbrwt], kernel="gaussian", bw=20, adjust=1, weights=mom_dt[tobacco==1])  
plot(d1_verylow, col="black", main="Kernel Density, Smoker: Very Low Bandwidth",  
      xlab = "Birthweight (g)")
```

### Kernel Density, Smoker: Very Low Bandwidth



2f

What are the benefits of the weighting approach (from part c)? What are the potential drawbacks? Pay particular attention to the issue of people with extremely high and extremely low values of the propensity score.

2g

Present your findings and interpret the results on the relationship between birthweight and smoking. For the estimates in parts (b) and (c), consider which of the following conditions must hold in order for that estimate to be valid: i. The treatment effect heterogeneity is linear in the propensity score. ii. The treatment effect heterogeneity is not linear in the propensity score. iii. The decision to smoke is completely randomly assigned. iv. Conditional on the exogenous variables the decision to smoke is randomly assigned.