

ARE 213 Problem Set 1A

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

Due 09/25/2020

Section 1

1. *Before getting started with the data work, first consider the table from Snow (1855) reproduced in the lecture notes (“Snow’s Table IX”). The table reports only means.
 - (a) Develop an approximate 95% confidence interval for “Deaths per 10,000 Houses” for Southwark and Vauxhall customers. Develop another 95% CI for the same quantity for Lambeth. Do the confidence intervals overlap?

Note that that we’re estimating p for a binomial distribution since deaths per 10,000 houses is the same as deaths per person (of course, scaled by persons per 10,000 households). Are we really dealing with a binomial distribution? Probably not, but it might not be a bad approximation if we think contaminated water is distributed randomly across space-time (so one person’s probability exposure and subsequent death is the same and independent of another person’s). Also, not everyone is equally susceptible to the virus (some have a higher p than others), but our estimate of p can be interpreted as an average p .

There are various ways to construct a confidence interval for an estimated binomial distribution. We use three different methods, all of which provide very similar estimates. The confidence intervals do not overlap.

```
# Southwark and Vauxhall
```

```
binom.confint(1263, 40046, method=c("asymptotic", "wilson", "agresti-coull"), type="central")
```

```
##           method    x      n      mean      lower      upper
## 1 agresti-coull 1263 40046 0.03153873 0.02987085 0.03329648
## 2   asymptotic 1263 40046 0.03153873 0.02982701 0.03325045
## 3      wilson 1263 40046 0.03153873 0.02987144 0.03329589
```

```
# Lambeth
```

```
binom.confint(98, 26107, method=c("asymptotic", "wilson", "agresti-coull"), type="central")
```

```
##           method    x      n      mean      lower      upper
## 1 agresti-coull  98 26107 0.003753783 0.003077893 0.004575688
## 2   asymptotic  98 26107 0.003753783 0.003011981 0.004495584
## 3      wilson  98 26107 0.003753783 0.003081460 0.004572122
```

- (b) Discuss either formally or intuitively the critical assumption that underlies your confidence intervals. Give a 2 or 3 sentence quote from Snow’s description (reproduced in Freedman (1991)) that supports this assumption.

To be confident that it is the choice of water company that is causing the difference in p and not some other factor, we need to be sure that there are not systematic differences between those who get their water from Southwark and Vauxhall and those who get it from Lambeth. John Snow argues that the two groups of people are comparable: “Each company supplies both rich and poor, both large houses and small; there is no difference either in the condition or occupation of the persons receiving the water of the different Companies... As there is no difference whatever in the houses or the people receiving the supply of the two Water Companies, or in any of the physical conditions with which they are surrounded, it is obvious that no experiment could have been devised which would more thoroughly test the effect of water supply on the progress of cholera than this, which circumstances placed ready made before the observer.” In that case, we are reasonably certain that the difference in water company is what causes the difference in mortality risk.

Section 2

We now move to some analysis of real data. The data portions of Problem Sets 1a and 1b are based heavily on the paper Almond, Chay, and Lee (2005), and problem sets from Ken Chay and John DiNardo based on some of the data used in the paper. The goal of this assignment is to examine the research question: what is the causal effect of maternal smoking during pregnancy on infant birthweight and other infant health outcomes. The data for the problem set is an extract of all births from the 1993 National Natality Detail Files for Pennsylvania. Each observation represents an infant-mother match. The data in Stata format can be downloaded from the bCourses website. There should be 48 variables in the data and, after you are finished with the cleaning steps described below, 114,610 observations.

The data here are “real” and quite imperfect, which will help simulate the unpleasantness of real world data work. Unlike the real world where you will confront this bleak situation largely alone, I will provide you with some hints for working your way through the raw data. You can download part of the codebook for the data to help you figure out the relevant variables.

2. The first order of business is to go through the code book, decide on the relevant variables, and process the data. This involves several steps:
 - (a) Fix missing values. In the data set several variables take on a value of, say, 9999 if missing. We have already checked for missing observations for about 2/3 of the variables. The remaining variables need to be checked and are the last 15 in the variables list (i.e. from ‘cardiac’ to ‘wgain’). Refer to the codebook for missing value codes. Produce an analysis data set that drops any observations with missing values.

```
# According to the codebook, for the following medical risk factor variables, 8 corresponds to
# "Factor not on certificate" and 9 corresponds to "Factor not classifiable": cardiac, lung, diabetes,
#herpes, chyper, phyper, pre4000, preterm

med_risk_factors <- c('cardiac', 'lung', 'diabetes', 'herpes', 'chyper', 'phyper', 'pre4000', 'preterm')

for (var in med_risk_factors){
  mom_dt[var] <- replace(mom_dt[var], which(mom_dt[var] == 8, arr.ind = TRUE), NA)
  mom_dt[var] <- replace(mom_dt[var], which(mom_dt[var] == 9, arr.ind = TRUE), NA)
}

# Below, arr.ind = TRUE returns the indices at which the row equals a certain value

# According to the codebook, for tobacco, 9 corresponds to "Unknown or not stated"
mom_dt$tobacco <- replace(mom_dt$tobacco, which(mom_dt$tobacco == 9, arr.ind = TRUE), NA)

# According to the codebook, for cigar, 99 corresponds to "Unknown or not stated"
mom_dt$cigar <- replace(mom_dt$cigar, which(mom_dt$cigar == 99, arr.ind = TRUE), NA)

# According to the codebook, for cigar6, 6 corresponds to "Unknown or not stated"
mom_dt$cigar6 <- replace(mom_dt$cigar6, which(mom_dt$cigar6 == 6, arr.ind = TRUE), NA)

# According to the codebook, for alcohol, 9 corresponds to "Unknown or not stated"
```

```

mom_dt$alcohol <- replace(mom_dt$alcohol, which(mom_dt$alcohol == 9, arr.ind = TRUE), NA)

# According to the codebook, for drink, 99 corresponds to "Unknown or not stated"
mom_dt$drink <- replace(mom_dt$drink, which(mom_dt$drink == 99, arr.ind = TRUE), NA)

# According to the codebook, for drink5, 5 corresponds to "Unknown or not stated"
mom_dt$drink5 <- replace(mom_dt$drink5, which(mom_dt$drink5 == 5, arr.ind = TRUE), NA)

# According to the codebook, for wgain (assuming that's wtgain in codebook),
# 99 corresponds to "Unknown or not stated"
mom_dt$wgain <- replace(mom_dt$wgain, which(mom_dt$wgain == 99, arr.ind = TRUE), NA)

# Make indicator for missing, will drop after comparison
setDT(mom_dt)
mom_dt[, miss := ifelse(complete.cases(mom_dt), 0, 1)]

```

- (b) If this were a real research project you would want to consider other approaches to missing data besides termination with extreme prejudice. What observations do you have to drop because of missing data? Might this affect your results? Do the data appear to be missing completely at random? How might you assess whether the data appear to be missing at random?

We know from the last problem set that if the data are missing at random then dropping them should not affect our results of the effect of smoking on birth weight. However, if the missing data is correlated with the treatment (smoking) or the outcome (birth weight) then it could bias our results. In tables 1 and 2 at the end of this document, we present the means and standard deviations for a number of the variables between the missing and nonmissing group, the difference in the means, t-statistic, and p-value (under the null that the difference in means is 0). We also present the proportion of missing data by each categorical group, though we do not calculate a statistical test for these variables. From the table, it does appear that there are differences in the missing and nonmissing data. For example, the mothers in the missing data are younger, less educated, less likely to be married, have more previous children, received less prenatal care, have a shorter time since the last birth, and have a lower gestation period. As discussed in the last problem set, we could formally assess whether the data is missing at random by regressing an indicator for the missing variable on the treatment.

```

# Compare missing to non missing
# categorical
birth_attendant <- mom_dt[,.(prop_miss = mean(miss)), by = birattnnd]
county_pop <- mom_dt[,.(prop_miss = mean(miss)), by = cntocpop]
state <- mom_dt[,.(prop_miss = mean(miss)), by = stresfip]
state <- state[order(stresfip)]
race <- mom_dt[,.(prop_miss = mean(miss)), by = mrace3]

category_dt <- rbind(birth_attendant, county_pop, state, race, use.names = FALSE)
category_dt[, Variable := c("Attendant: M.D.", "Attendant: D.O.", "Attendant: C.N.M.",
                           "Attendant: Other Midwife",
                           "Attendant: Other", "County: 100k-250k", "County: 250k-500k", "County: 500k-1m",
                           "County: 1 million +", "State: Foreign", "State: AZ", "State: CA",
                           "State: CO", "State: CT", "State: DE", "State: DC", "State: FL", "State: GA",
                           "State: IL", "State: IA", "State: KY", "State: ME", "State: MD", "State: MA",
                           "State: MI", "State: MN", "State: MO", "State: NE", "State: NV", "State: NJ",
                           "State: NY", "State: NC", "State: ND", "State: OH", "State: OK", "State: PA",
                           "State: RI", "State: SC", "State: SD", "State: TN", "State: TX", "State: VA",
                           "State: WA", "State: WV", "State: WY", "Race: White", "Race: Black", "Race: Other")]

category_dt <- category_dt[, 2:3]
colnames(category_dt) <- c("Proportion Missing", "Variable")
setcolorder(category_dt, c("Variable", "Proportion Missing"))

```

```

# numeric
compare_dt <- transpose(mom_dt[,lapply(.SD, mean, na.rm=TRUE), .SDcols = c(2, 6, 9:18, 20, 22,
                                                                           25:30, 33:43, 45:46, 48),
                        by = miss])
compare_dt <- cbind(compare_dt, transpose(mom_dt[,lapply(.SD, sd, na.rm=TRUE), .SDcols = c(2, 6, 9:18,
                                                                           20, 22, 25:30,
                                                                           33:43, 45:46, 48),
                        by = miss]))
colnames(compare_dt) <- c("Nonmiss means", "Miss means", "Nonmiss sd", "Miss sd")
compare_dt <- compare_dt[2:35,]
compare_dt[, Variable := c("Hospital", "Mother age", "Mother educ", "Marital status", "Prenatal adequacy",
                           "Number living child", "Number dead or living child",
                           "Total live birth or terminations", "Birth order", "Month prenatal began",
                           "Number prenatal visits", "Time since last birth", "Father age",
                           "Father educ", "Gestation", "Child sex", "Birth weight", "Number born",
                           "One min Apgar", "Five min Apgar", "Anemia", "Cardiac disease",
                           "Lung disease", "Diabetes", "Herpes", "Chron. hypertension",
                           "Preg. hypertension", "Previous heavy birth", "Previous preterm", "Tobacco use",
                           "Number cigarettes", "Alcohol use", "Number drinks", "Weight gain")]

formulas <- paste("mom_dt$", names(mom_dt)[c(2, 6, 9:18, 20, 22, 25:30, 33:43, 45:46, 48)], "~ mom_dt$miss")
t_test <- t(sapply(formulas, function(f) {
  res <- t.test(as.formula(f))
  c(res$statistic, p.value=res$p.value)
}))

colnames(t_test) <- c("t-stat", "p-value")

compare_dt <- cbind(compare_dt, t_test)
compare_dt[, Difference := `Nonmiss means` - `Miss means`]

setcolorder(compare_dt, c("Variable", "Nonmiss means", "Nonmiss sd", "Miss means", "Miss sd", "Difference", "p-value"))

mom_dt <- na.omit(mom_dt)

```

(c) Produce a summary table describing the final analysis data set.

We create summary tables similar to the tables in 2b, but this time we compare the means of smokers vs non-smokers. These tables are tables 3 and 4 at the end of the document. To see means/standard deviations for the entire dataset, refer to the nonmissing data columns of Table 2. We will discuss the differences between the groups in question 3b.

```

# Recode to binary 0/1 treatment
# tobacco is 1: yes, tobacco use during pregnancy and 2: no tobacco use during pregnancy
mom_dt[, tobacco := ifelse(tobacco==2, 0, 1)]

# Compare smoker to nonsmoker
# Categorical
birth_attendant <- mom_dt[,.(prop = mean(tobacco)), by = birattnd]
county_pop <- mom_dt[,.(prop = mean(tobacco)), by = cntocpop]
state <- mom_dt[,.(prop = mean(tobacco)), by = stresfip]
state <- state[order(stresfip)]
race <- mom_dt[,.(prop = mean(tobacco)), by = mrace3]

category_sum_dt <- rbind(birth_attendant, county_pop, state, race, use.names = FALSE)
category_sum_dt[, Variable := c("Attendant: M.D.", "Attendant: D.O.", "Attendant: C.N.M.",
                                "Attendant: Other Midwife",

```

```

      "Attendant: Other", "County: 100k-250k", "County: 250k-500k", "County: 500k-1m",
      "County: 1 million +", "State: Foreign", "State: AZ", "State: CA",
      "State: CO", "State: CT", "State: DE", "State: DC", "State: FL", "State: GA",
      "State: IL", "State: IA", "State: KY", "State: ME", "State: MD", "State: MA",
      "State: MI", "State: MN", "State: MO", "State: NE", "State: NV", "State: NJ",
      "State: NY", "State: NC", "State: ND", "State: OH", "State: OK", "State: PA",
      "State: RI", "State: SC", "State: SD", "State: TN", "State: VA",
      "State: WA", "State: WV", "State: WY", "Race: White", "Race: Black", "Race: Other"
category_sum_dt <- category_sum_dt[, 2:3]
colnames(category_sum_dt) <- c("Proportion Smoking", "Variable")
setcolorder(category_sum_dt, c("Variable", "Proportion Smoking"))

#numeric
summary_dt <- transpose(mom_dt[,lapply(.SD, mean, na.rm=TRUE), .SDcols = c(2, 6, 9:18, 20, 22,
                                                                    25:30, 33:41, 43, 45:46, 48),
                        by = tobacco])
summary_dt <- cbind(summary_dt, transpose(mom_dt[,lapply(.SD, sd, na.rm=TRUE), .SDcols = c(2, 6, 9:18, 20, 22,
                                                                    25:30, 33:41, 43,
                                                                    45:46, 48),
                        by = tobacco]))
colnames(summary_dt) <- c("Nonsmoker means", "Smoker means", "Nonsmoker sd", "Smoker sd")
summary_dt <- summary_dt[2:34,]
summary_dt[, Variable := c("Hospital", "Mother age", "Mother educ", "Marital status", "Prenatal adequacy",
                          "Number living child", "Number dead or living child",
                          "Total live birth or terminations", "Birth order", "Month prenatal began",
                          "Number prenatal visits", "Time since last birth", "Father age", "Father educ",
                          "Gestation", "Child sex", "Birth weight", "Number born", "One min Apgar",
                          "Five min Apgar", "Anemia", "Cardiac disease", "Lung disease", "Diabetes",
                          "Herpes", "Chron. hypertension", "Preg. hypertension", "Previous heavy birth",
                          "Previous preterm", "Number cigarettes", "Alcohol use", "Number drinks", "Weight g

formulas <- paste("mom_dt$", names(mom_dt)[c(2, 6, 9:18, 20, 22, 25:30, 33:41, 43, 45:46, 48)], "~ mom_dt$tobacco")
t_test <- t(sapply(formulas, function(f) {
  res <- t.test(as.formula(f))
  c(res$statistic, p.value=res$p.value)
}))

colnames(t_test) <- c("t-stat", "p-value")

summary_dt <- cbind(summary_dt, t_test)
summary_dt[, Difference := `Nonsmoker means` - `Smoker means`]

setcolorder(summary_dt, c("Variable", "Nonsmoker means", "Nonsmoker sd", "Smoker means", "Smoker sd", "Difference"))

```

3. The next part of the assignment is to try to estimate the “causal” effect of maternal smoking during pregnancy on infant birth weight. Let’s start out using techniques that are familiar, and think about whether they are likely to work in this context. Answer the following questions.

- (a) Compute the mean difference in APGAR scores (both five and one minute versions) as well as birthweight by smoking status.

```

# According to the codebook, omaps is the one minute APGAR score and fmaps is the five minute APGAR score
# Both are a score from 0-10
# dbrwt (assuming that corresponds to dbirwt in codebook) is birthweight in grams

smoker <- subset(mom_dt, mom_dt$tobacco == 1)

```

```
nonsmoker <- subset(mom_dt, mom_dt$tobacco == 0)

# Mean difference in one minute APGAR score by smoking status
mean_diff_1min_apgar <- mean(smoker$omaps) - mean(nonsmoker$omaps)
print(mean_diff_1min_apgar)
```

```
## [1] -0.01743508
```

```
# Mean difference in five minute APGAR score by smoking status
mean_diff_5min_apgar <- mean(smoker$fmaps) - mean(nonsmoker$fmaps)
print(mean_diff_5min_apgar)
```

```
## [1] -0.0001498085
```

```
# Mean difference in birthweight by smoking status
mean_diff_birthweight <- mean(smoker$dbrwt) - mean(nonsmoker$dbrwt)
print(mean_diff_birthweight)
```

```
## [1] -240.4778
```

- (b) Under what circumstances can one identify the average treatment effect of maternal smoking by comparing the unadjusted difference in mean birth weight of infants of smoking and non-smoking mothers? Estimate its impact under this assumption. Provide and comment on some evidence for or against the validity of the assumption.

We can identify the average treatment effect by comparing the unadjusted difference in means if the treatment (smoking) is randomly assigned. Under this assumption, we would estimate that smoking results in an effect of -240.478 grams on birth weight. The mean of all birth weights is 3373.291 so this represents a decrease of -7.1 percent.

However, based on our table in 2c, random assignment is likely not a valid assumption. Smokers and non smokers are different on many dimensions, some of which may affect birth weight. For example, smokers started prenatal care later and had fewer prenatal visits.

- (c) Suppose that maternal smoking is randomly assigned conditional on the other observable “predetermined” determinants of infant birth weight. First discuss which (if any) of the variables contained in the data set can clearly be considered to be predetermined. In general, what kinds of variables can be considered predetermined and what kinds of variables cannot?

A variable can be considered “predetermined” if it affects selection into treatment, but is not in turn affected by the treatment variable - that is while there is correlation between the predetermined and the treatment variable, that correlation should be driven by the effect of the predetermined variable on the treatment variable, and not the reverse direction.

In our data set, our treatment variable is smoking during pregnancy, so any variables that may be affected by smoking during pregnancy would not be considered predetermined variables. For example, it is likely that smoking during pregnancy affects prenatal care, alcohol use during pregnancy, weight gain during pregnancy, birth month (i.e. whether the birth is premature), and the health condition of the mother (i.e. anemia, diabetes, cardiac disease, etc.). Those variables would not be considered predetermined.

However, variables that could affect the likelihood of a mother smoking during pregnancy, but would not in turn be affected by a mother’s prenatal smoking status, would be considered predetermined. In our data set, variables such as state of residence, population density of county, and mother’s race, all arguably fit this criteria and can be considered predetermined variables. Other variables, such as mother’s age, mother’s education, mother’s marital status, father’s race,

father's age, and father's education, could also be considered predetermined variables, although there is more ambiguity in these cases. Particularly if we think smoking during pregnancy reflects a higher probability of smoking before pregnancy and in early life more generally, then early smoking habits could affect the age at which a women becomes pregnant, the education level she has when she is pregnant, her marital status, and if married, the spouse she chooses and subsequently their characteristics. If that is the case, then smoking during pregnancy could affect these variables and thus they should not be considered "predetermined." However, if we define our treatment narrowly as smoking during pregnancy, without inferring anything about smoking habits prior to pregnancy, than the choice to smoke during pregnancy should not affect a mother's age, education level, father's age, or father's education (marital status could still be ambiguous). From this perspective, these additional variables could arguably also be considered predetermined.

- (d) What does "selection on observables" imply about the relationship between maternal smoking and unobservable determinants of birth weight conditional on the observables? Use a basic linear regression model, in conjunction with your answer to (c), to estimate the impact of smoking and report your estimates. Under what circumstances is the average treatment effect identified?

The key assumption underlying a "selection on observables" design is that the treatment is as good as randomly assigned after we condition on observables. In other words, we assume that we observe *all* the factors that affect treatment assignment (smoking) and are correlated with the potential outcomes (birth weight). If there is systematic selection into treatment, we assume this selection is only a function of the observables. That is, we assume that maternal smoking is uncorrelated with unobservable determinants of birth weight conditional on the observables. If these assumptions hold, then a regression of maternal smoking on birth weight, conditioning on observables, will estimate the ATE.

In a selection on observables design we estimate two models, based on whether we take a more strict or more relaxed classification of which variables are predetermined (as discussed in part c above):

With a stricter definition of predetermined variables, we estimate

$$birthweight_i = \alpha_i + \beta Smoking_i + \delta_1 state_i + \delta_2 county_{pop}_i + \delta_3 mother_race + \epsilon_i \quad (1)$$

With a more relaxed definition of predetermined variables, we estimate

$$birthweight_i = \alpha_i + \beta Smoking_i + \delta_1 state_i + \delta_2 county_{pop}_i + \delta_3 mother_race + \delta_4 mother_age_i + \delta_5 mother_educ_i + \delta_6 marital_status_i + \delta_7 father_age_i + \delta_8 father_educ + \epsilon_i \quad (2)$$

```
# Selection on observables model - strict definition of predetermined variables
lm1 <- lm(dbrwt ~ tobacco + factor(stresfip) + factor(cntocpop) + factor(mrace3), mom_dt)
summary(lm1)
```

```
##
## Call:
## lm(formula = dbrwt ~ tobacco + factor(stresfip) + factor(cntocpop) +
##     factor(mrace3), data = mom_dt)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -3234.0  -305.5    26.0   354.0  2824.5
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   3296.006    134.798   24.451 < 2e-16 ***
## tobacco      -242.453     4.634  -52.322 < 2e-16 ***
## factor(stresfip)4    85.725    426.127    0.201  0.8406
## factor(stresfip)6   210.286    213.065    0.987  0.3237
## factor(stresfip)8   -99.370    426.128   -0.233  0.8156
## factor(stresfip)9  -157.620    316.029   -0.499  0.6180
```



```
## factor(stresfip)10 103.121 139.774 0.738 0.4607
## factor(stresfip)11 -62.031 426.129 -0.146 0.8843
## factor(stresfip)12 -38.988 175.311 -0.222 0.8240
## factor(stresfip)13 84.466 233.410 0.362 0.7174
## factor(stresfip)17 255.408 289.015 0.884 0.3769
## factor(stresfip)19 -51.101 426.128 -0.120 0.9045
## factor(stresfip)21 604.048 289.020 2.090 0.0366 *
## factor(stresfip)23 275.994 587.375 0.470 0.6384
## factor(stresfip)24 144.715 141.141 1.025 0.3052
## factor(stresfip)25 230.863 289.012 0.799 0.4244
## factor(stresfip)26 205.112 356.530 0.575 0.5651
## factor(stresfip)27 509.482 587.384 0.867 0.3857
## factor(stresfip)29 607.804 587.377 1.035 0.3008
## factor(stresfip)31 445.994 587.375 0.759 0.4477
## factor(stresfip)32 160.804 587.377 0.274 0.7843
## factor(stresfip)34 94.469 135.422 0.698 0.4854
## factor(stresfip)36 19.337 144.377 0.134 0.8935
## factor(stresfip)37 88.028 213.071 0.413 0.6795
## factor(stresfip)38 -357.272 587.382 -0.608 0.5430
## factor(stresfip)39 48.005 138.263 0.347 0.7284
## factor(stresfip)40 169.643 426.129 0.398 0.6906
## factor(stresfip)42 124.439 134.774 0.923 0.3558
## factor(stresfip)44 -130.272 587.382 -0.222 0.8245
## factor(stresfip)45 -466.187 316.031 -1.475 0.1402
## factor(stresfip)46 -45.272 587.382 -0.077 0.9386
## factor(stresfip)47 172.794 269.504 0.641 0.5214
## factor(stresfip)51 -21.694 183.648 -0.118 0.9060
## factor(stresfip)53 -102.736 356.528 -0.288 0.7732
## factor(stresfip)54 67.089 149.162 0.450 0.6529
## factor(stresfip)55 81.060 316.032 0.256 0.7976
## factor(cntocpop)1 38.266 5.362 7.137 9.6e-13 ***
## factor(cntocpop)2 40.537 4.478 9.053 < 2e-16 ***
## factor(cntocpop)3 30.190 4.953 6.095 1.1e-09 ***
## factor(mrace3)2 -201.149 12.184 -16.509 < 2e-16 ***
## factor(mrace3)3 -242.488 5.496 -44.119 < 2e-16 ***
## ---
## Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 571.7 on 114569 degrees of freedom
## Multiple R-squared: 0.04586, Adjusted R-squared: 0.04553
## F-statistic: 137.7 on 40 and 114569 DF, p-value: < 2.2e-16
```

```
# Selection on observables model - relaxed definition of predetermined variables
```

```
lm2 <- lm(dbrwt ~ tobacco + factor(stresfip) + factor(cntocpop) + factor(mrace3) + dimage + dmeduc + factor(dmar) + dfeduc, data = mom_dt)
summary(lm2)
```

```
##
## Call:
## lm(formula = dbrwt ~ tobacco + factor(stresfip) + factor(cntocpop) +
##     factor(mrace3) + dimage + dmeduc + factor(dmar) + dfeduc +
##     dfeduc, data = mom_dt)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -3266.8  -305.5    27.1   355.5  2847.0
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
```



```

## (Intercept)      3143.3496   135.3308   23.227   < 2e-16 ***
## tobacco          -211.5450     4.8451  -43.661   < 2e-16 ***
## factor(stresfip)4   104.9295   424.9725    0.247   0.80498
## factor(stresfip)6   230.1092   212.4899    1.083   0.27885
## factor(stresfip)8   -91.4810   424.9685   -0.215   0.82956
## factor(stresfip)9  -115.4927   315.1802   -0.366   0.71404
## factor(stresfip)10  119.2259   139.3989    0.855   0.39239
## factor(stresfip)11  -27.2348   424.9703   -0.064   0.94890
## factor(stresfip)12  -24.3628   174.8365   -0.139   0.88918
## factor(stresfip)13   93.1325   232.7895    0.400   0.68910
## factor(stresfip)17  304.4973   288.2343    1.056   0.29078
## factor(stresfip)19  -35.0911   424.9799   -0.083   0.93419
## factor(stresfip)21  598.6825   288.2359    2.077   0.03780 *
## factor(stresfip)23  294.5712   585.7923    0.503   0.61506
## factor(stresfip)24  164.6769   140.7645    1.170   0.24205
## factor(stresfip)25  246.5205   288.2296    0.855   0.39239
## factor(stresfip)26  234.9028   355.5717    0.661   0.50885
## factor(stresfip)27  461.8891   585.8024    0.788   0.43042
## factor(stresfip)29  655.4602   585.7972    1.119   0.26318
## factor(stresfip)31  464.4897   585.7795    0.793   0.42781
## factor(stresfip)32  135.0456   585.8041    0.231   0.81768
## factor(stresfip)34  100.3322   135.0549    0.743   0.45754
## factor(stresfip)36   33.3774   143.9878    0.232   0.81669
## factor(stresfip)37   83.7336   212.4957    0.394   0.69355
## factor(stresfip)38 -320.2048   585.8013   -0.547   0.58465
## factor(stresfip)39   67.8776   137.8950    0.492   0.62255
## factor(stresfip)40  170.6424   424.9829    0.402   0.68803
## factor(stresfip)42  149.7446   134.4131    1.114   0.26525
## factor(stresfip)44 -130.7123   585.7846   -0.223   0.82343
## factor(stresfip)45 -450.0991   315.1740   -1.428   0.15327
## factor(stresfip)46    3.9693   585.7958    0.007   0.99459
## factor(stresfip)47  196.0279   268.7782    0.729   0.46580
## factor(stresfip)51 -10.1072   183.1554   -0.055   0.95599
## factor(stresfip)53 -97.4314   355.5682   -0.274   0.78407
## factor(stresfip)54   84.3980   148.7595    0.567   0.57048
## factor(stresfip)55   89.4668   315.1812    0.284   0.77652
## factor(cntocpop)1   28.5759     5.3625     5.329   9.9e-08 ***
## factor(cntocpop)2   45.0628     4.5035    10.006   < 2e-16 ***
## factor(cntocpop)3   38.7955     4.9875     7.779   7.4e-15 ***
## factor(mrace3)2    -210.1045    12.1728  -17.260   < 2e-16 ***
## factor(mrace3)3    -183.7404     5.9943  -30.652   < 2e-16 ***
## dimage              1.6388     0.5012     3.270   0.00108 **
## dmeduc              2.0312     1.0318     1.969   0.04901 *
## factor(dmar)2      -79.6520     4.9422  -16.117   < 2e-16 ***
## dfage              0.7715     0.4154     1.857   0.06329 .
## dfeduc             2.9109     0.9904     2.939   0.00329 **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 570.1 on 114564 degrees of freedom
## Multiple R-squared:  0.0511, Adjusted R-squared:  0.05073
## F-statistic: 137.1 on 45 and 114564 DF, p-value: < 2.2e-16

```

If our selection on observables assumption holds, we can interpret these results as indicating the causal effect of maternal smoking on birth weight. Specifically, in our estimation of (1), we can interpret our results as ceteris paribus, maternal smoking has an ATE of decreasing birth weight by 242.453 grams (a decrease of 7.356 percent from a base of 3296.006 grams - note that the omitted categories are white, large county, foreign residents), which is statistically significant at the 1% level. In our estimation of (2), our results indicate that ceteris paribus, maternal smoking has an ATE of decreasing birth weight by 211.545 grams (a decrease of 6.73 percent from a base of 3143.3496 grams - where the omitted category

is white, large county, married, foreign residents), which is also significant at the 1% level.

Summary Tables: Questions 2b, 2c

```
print(xtable(category_dt, caption = 'Proportion Missing by Group', digits = 2),  
      include.rownames = FALSE, size = "small", comment = FALSE)
```

```
print(xtable(compare_dt, caption = 'Difference in Means Missing v Nonmissing', digits = 2),  
      include.rownames = FALSE, size = "small", comment = FALSE)
```

```
print(xtable(category_sum_dt, caption = 'Proportion Smoking by Group', digits = 2),  
      include.rownames = FALSE, size = "small", comment = FALSE)
```

```
print(xtable(summary_dt, caption = 'Difference in Means Smoker v Nonsmoker', digits = 2),  
      include.rownames = FALSE, size = "small", comment = FALSE)
```

Variable	Proportion Missing
Attendant: M.D.	0.05
Attendant: D.O.	0.06
Attendant: C.N.M.	0.03
Attendant: Other Midwife	0.08
Attendant: Other	0.06
County: 100k-250k	0.05
County: 250k-500k	0.05
County: 500k-1m	0.03
County: 1 million +	0.06
State: Foreign	0.00
State: AZ	0.00
State: CA	0.08
State: CO	0.00
State: CT	0.00
State: DE	0.03
State: DC	0.00
State: FL	0.07
State: GA	0.00
State: IL	0.17
State: IA	0.00
State: KY	0.00
State: ME	0.00
State: MD	0.05
State: MA	0.00
State: MI	0.00
State: MN	0.50
State: MO	0.50
State: NE	0.00
State: NV	0.00
State: NJ	0.04
State: NY	0.18
State: NC	0.00
State: ND	0.00
State: OH	0.06
State: OK	0.33
State: PA	0.05
State: RI	0.00
State: SC	0.20
State: SD	0.00
State: TN	0.14
State: TX	1.00
State: VA	0.12
State: WA	0.00
State: WV	0.06
State: WY	0.00
Race: White	0.04
Race: Black	0.09
Race: Other	0.05

Table 1: Proportion Missing by Group

Variable	Nonmiss means	Nonmiss sd	Miss means	Miss sd	Difference	t-stat	p-value
Hospital	1.02	0.13	1.01	0.11	0.01	4.16	0.00
Mother age	27.76	5.70	27.05	5.97	0.71	8.84	0.00
Mother educ	13.21	2.27	12.51	2.26	0.70	23.16	0.00
Marital status	1.25	0.43	1.44	0.50	-0.19	-28.00	0.00
Prenatal adequacy	1.30	0.55	1.63	0.79	-0.33	-31.58	0.00
Number living child	0.97	1.15	1.24	1.43	-0.27	-14.16	0.00
Number dead or living child	1.99	1.17	2.27	1.47	-0.28	-14.38	0.00
Total live birth or terminations	2.42	1.52	2.81	1.87	-0.39	-15.76	0.00
Birth order	2.41	1.46	2.78	1.74	-0.37	-15.96	0.00
Month prenatal began	2.50	1.33	2.80	1.92	-0.30	-11.79	0.00
Number prenatal visits	11.15	3.52	9.32	4.90	1.84	28.30	0.00
Time since last birth	350.41	362.33	315.97	355.26	34.44	7.23	0.00
Father age	30.06	6.41	29.61	7.04	0.46	4.84	0.00
Father educ	13.28	2.33	12.67	2.29	0.60	19.61	0.00
Gestation	39.15	2.44	38.53	3.42	0.62	13.77	0.00
Child sex	1.49	0.50	1.48	0.50	0.00	0.14	0.89
Birth weight	3373.29	585.17	3191.90	716.95	181.39	19.03	0.00
Number born	1.03	0.17	1.04	0.21	-0.01	-4.17	0.00
One min Apgar	8.12	1.26	7.90	1.57	0.21	10.18	0.00
Five min Apgar	9.01	0.71	8.88	1.03	0.13	9.47	0.00
Anemia	1.99	0.10	1.99	0.12	0.00	2.66	0.01
Cardiac disease	1.99	0.08	1.99	0.09	0.00	0.70	0.49
Lung disease	1.99	0.08	1.99	0.10	0.00	1.57	0.12
Diabetes	1.97	0.16	1.97	0.16	0.00	0.07	0.94
Herpes	1.99	0.08	1.99	0.10	0.00	2.44	0.01
Chron. hypertension	1.99	0.09	1.99	0.10	0.00	1.31	0.19
Preg. hypertension	1.97	0.17	1.97	0.16	-0.00	-2.03	0.04
Previous heavy birth	1.99	0.12	1.99	0.10	-0.00	-2.81	0.01
Previous preterm	1.99	0.12	1.98	0.16	0.01	5.16	0.00
Tobacco use	1.84	0.37	1.57	0.50	0.27	41.11	0.00
Number cigarettes	1.91	5.30	3.94	7.42	-2.03	-18.74	0.00
Alcohol use	1.99	0.10	1.63	0.48	0.36	56.48	0.00
Number drinks	0.03	0.62	0.16	1.47	-0.13	-5.33	0.00
Weight gain	30.36	11.88	30.78	13.14	-0.43	-1.71	0.09

Table 2: Difference in Means Missing v Nonmissing

Variable	Proportion Smoking
Attendant: M.D.	0.16
Attendant: D.O.	0.15
Attendant: C.N.M.	0.18
Attendant: Other Midwife	0.02
Attendant: Other	0.19
County: 100k-250k	0.20
County: 250k-500k	0.16
County: 500k-1m	0.11
County: 1 million +	0.15
State: Foreign	0.11
State: AZ	0.00
State: CA	0.00
State: CO	0.00
State: CT	0.00
State: DE	0.10
State: DC	0.00
State: FL	0.19
State: GA	0.33
State: IL	0.20
State: IA	0.00
State: KY	0.00
State: ME	0.00
State: MD	0.09
State: MA	0.00
State: MI	0.00
State: MN	0.00
State: MO	0.00
State: NE	0.00
State: NV	0.00
State: NJ	0.07
State: NY	0.13
State: NC	0.00
State: ND	0.00
State: OH	0.23
State: OK	0.00
State: PA	0.16
State: RI	0.00
State: SC	0.00
State: SD	0.00
State: TN	0.00
State: VA	0.05
State: WA	0.00
State: WV	0.20
State: WY	0.25
Race: White	0.16
Race: Black	0.17
Race: Other	0.03

Table 3: Proportion Smoking by Group

Variable	Nonsmoker means	Nonsmoker sd	Smoker means	Smoker sd	Difference	t-stat	p-value
Hospital	1.02	0.14	1.00	0.06	0.02	28.15	0.00
Mother age	28.06	5.67	26.17	5.61	1.88	41.56	0.00
Mother educ	13.44	2.30	11.99	1.63	1.46	102.72	0.00
Marital status	1.21	0.41	1.48	0.50	-0.27	-70.08	0.00
Prenatal adequacy	1.28	0.53	1.41	0.63	-0.14	-27.41	0.00
Number living child	0.93	1.13	1.15	1.22	-0.22	-22.60	0.00
Number dead or living child	1.95	1.15	2.18	1.27	-0.23	-22.98	0.00
Total live birth or terminations	2.36	1.48	2.74	1.67	-0.39	-29.06	0.00
Birth order	2.35	1.42	2.73	1.60	-0.38	-29.95	0.00
Month prenatal began	2.45	1.28	2.75	1.51	-0.30	-25.11	0.00
Number prenatal visits	11.25	3.45	10.63	3.84	0.63	20.52	0.00
Time since last birth	358.71	364.07	306.63	349.76	52.08	18.33	0.00
Father age	30.27	6.34	28.96	6.65	1.31	24.60	0.00
Father educ	13.49	2.37	12.13	1.67	1.37	94.00	0.00
Gestation	39.17	2.39	39.05	2.71	0.13	5.88	0.00
Child sex	1.49	0.50	1.48	0.50	0.00	1.04	0.30
Birth weight	3411.62	579.73	3171.14	572.08	240.48	51.98	0.00
Number born	1.03	0.18	1.02	0.15	0.01	5.26	0.00
One min Apgar	8.12	1.26	8.10	1.27	0.02	1.71	0.09
Five min Apgar	9.01	0.71	9.01	0.71	0.00	0.03	0.98
Anemia	1.99	0.10	1.99	0.12	0.00	4.61	0.00
Cardiac disease	1.99	0.08	1.99	0.08	-0.00	-1.50	0.13
Lung disease	1.99	0.08	1.99	0.10	0.00	3.80	0.00
Diabetes	1.97	0.16	1.97	0.16	-0.00	-0.02	0.98
Herpes	1.99	0.08	1.99	0.08	0.00	0.99	0.32
Chron. hypertension	1.99	0.09	1.99	0.08	-0.00	-2.05	0.04
Preg. hypertension	1.97	0.18	1.98	0.14	-0.01	-10.51	0.00
Previous heavy birth	1.98	0.12	1.99	0.09	-0.01	-9.16	0.00
Previous preterm	1.99	0.11	1.98	0.15	0.01	10.37	0.00
Number cigarettes	0.00	0.00	11.96	7.47	-11.96	-216.57	0.00
Alcohol use	2.00	0.07	1.97	0.18	0.03	21.83	0.00
Number drinks	0.01	0.25	0.14	1.44	-0.13	-11.77	0.00
Weight gain	30.52	11.56	29.47	13.45	1.05	9.92	0.00

Table 4: Difference in Means Smoker v Nonsmoker