

Replication 2 Megan Ellertson

Megan Ellertson

4/11/2021

Replication 2

```
## -- Attaching packages ----- tidyverse 1.3.0 --

## v ggplot2 3.3.3      v purrr 0.3.4
## v tibble 3.0.5       v dplyr 1.0.3
## v tidyr 1.1.2        v stringr 1.4.0
## v readr 1.4.0        v forcats 0.5.0

## -- Conflicts ----- tidyverse_conflicts() --
## x dplyr::filter() masks stats::filter()
## x dplyr::lag()     masks stats::lag()

##
## Please cite as:

## Hlavac, Marek (2018). stargazer: Well-Formatted Regression and Summary Statistics Tables.

## R package version 5.2.2. https://CRAN.R-project.org/package=stargazer

##
## Attaching package: 'kableExtra'

## The following object is masked from 'package:dplyr':
##
## group_rows
```

#1. ##Data Prep The following provides the code for the data formatting which is necessary prior to the analysis. Covariates which were to be included in both the Logit and OLS models were included in the first portion (plugged into the “control” dataset). The covariates specific to the Probit (cubed) and OLS (quadratic) were then separated into two different dataframe’s to keep the analysis separate.

```
control <- combo %>%
  mutate(agesq = age^2,
         agecube = age^3,
         educsq = educ*educ,
         u74 = case_when(re74 == 0 ~ 1, TRUE ~ 0),
         u75 = case_when(re75 == 0 ~ 1, TRUE ~ 0),
         interaction1 = educ*re74,
```

Table 1: Logit Model

Predictor	B	SE	t	p
Intercept	-34.65	4.095	-8.46	<0.001
age	2.32	0.353	6.57	<0.001
agesq	-0.06	0.011	-5.73	<0.001
agecube	0.00	0.000	4.84	<0.001
educ	0.72	0.664	1.09	0.278
educsq	-0.03	0.074	-0.43	0.669
educcube	0.00	0.003	-0.30	0.761
marr	-1.54	0.255	-6.04	<0.001
nodegree	0.92	0.345	2.68	0.007
black	3.84	0.266	14.41	<0.001
hisp	1.69	0.408	4.14	<0.001
re74	0.00	0.000	1.74	0.081
re74sq	0.00	0.000	-2.15	0.031
re74cube	0.00	0.000	2.62	0.009
re75	0.00	0.000	1.04	0.300
re75sq	0.00	0.000	-2.33	0.020
re75cube	0.00	0.000	2.45	0.014
u74	2.15	0.415	5.17	<0.001
u75	0.46	0.352	1.30	0.193

```

    re74sq = re74^2,
    re75sq = re75^2,
    interaction2 = u74*hisp)
#Logit related dataframe
nsw_dw_cpscontrol_logit <- control %>%
  mutate(educcube = educ^3,
         re74cube=re74^3,
         re75cube=re75^3)
#OLS related dataframe
nsw_dw_cpscontrol_ols <- control

```

##a & b The following equations are for the logit and OLS versions of the basic equation provided in the mixtape nsw_pscore.do. This expands off the basic logit provided.

The following model is the Logit model which includes through the cubed polynomial of each non-binary variable from the basic equation provided.

$$\text{logit}(\text{treat}) = \text{age} + \text{age}^2 + \text{age}^3 + \text{educ} + \text{educ}^2 + \text{educ}^3 + \text{marr} + \text{nodegree} + \text{black} + \text{hisp} + \text{re74} + \text{re74}^2 + \text{re74}^3 + \text{re75} + \text{re75}^2 + \text{re75}^3 + \text{u74} + \text{u75}$$

```

logit <- glm(treat ~ age + agesq + agecube + educ + educsq
            + educcube + marr + nodegree
            + black + hisp + re74 + re74sq + re74cube
            + re75 + re75sq + re75cube + u74 + u75,
            family = binomial(link =
            "logit"), data = nsw_dw_cpscontrol_logit)

```

Running this Logit regression provides the following output table. Each of the included variables except for the education related variables are statistically significant.

Table 2: OLS Model

Predictor	B	SE	t	p
Intercept	-0.39	0.031	-12.34	<0.001
age	0.03	0.003	11.59	<0.001
agesq	0.00	0.000	-11.11	<0.001
agecube	0.00	0.000	10.44	<0.001
educ	0.01	0.001	5.90	<0.001
educsq	0.00	0.000	-6.39	<0.001
marr	-0.02	0.002	-10.53	<0.001
nodegree	0.02	0.003	8.15	<0.001
black	0.10	0.003	35.97	<0.001
hisp	0.01	0.003	2.37	0.018
re74	0.00	0.000	-0.51	0.609
re74sq	0.00	0.000	1.37	0.170
re75	0.00	0.000	-3.24	0.001
re75sq	0.00	0.000	3.02	0.003
u74	0.04	0.004	11.36	<0.001
u75	0.01	0.004	3.70	<0.001

The following model is the OLS regression and version which will be included along with the Logit model. This includes the covariates from the basic model up to the squared polynomial, and not including the binary covariates with a polynomial version. Note that the cube polynomial of the age covariate is kept in the analysis because it was included in the basic model.

$$ols(treat) = age + age^2 + educ + educ^2 + marr + nodegree + black + hisp + re74 + re74^2 + re75 + re75^2 + u74 + u75 + educ * hisp + e$$

```
ols <- lm(treat ~ age + agesq + agecube + educ + + educsq +
          marr + nodegree+ black + hisp + re74 + re74sq
          + re75 + re75sq + u74 + u75,
          data = nsw_dw_cpscontrol_ols)
```

```
ols %>%
  tidy() %>%
  mutate(
    p.value = scales::pvalue(p.value),
    term = c("Intercept", "age", "agesq", "agecube", "educ",
             "educsq", "marr", "nodegree", "black", "hisp",
             "re74", "re74sq", "re75", "re75sq", "u74",
             "u75")
  ) %>%
  kable(
    caption= "OLS Model",
    col.names = c("Predictor", "B", "SE", "t", "p"),
    digits = c(0, 2, 3, 2, 3)
  )
```

In the OLS version of the model, more of the included variables are statistically significant compared to the Logit model. Only the “re74” covariate is not statistically significant.

C.

To map out the p-scores for the Logit and OLS models, the fitted values of the regressions are taken respectively. How this was done with the syntax is provided below:

```
nsw_dw_cpscontrol_logit <- nsw_dw_cpscontrol_logit %>%
  mutate(pscore = logit$fitted.values)

nsw_dw_cpscontrol_ols <- nsw_dw_cpscontrol_ols %>%
  mutate(pscore = ols$fitted.values)
```

The means of these p-scores for both the Logit and the OLS models are determined for the treated and control groups. The mean of these p-scores among the two groups for each model are provided in the table below. The Logit treated is at about 0.4 while the control group is far lower. The OLS is also skewed far lower than the more centered Logit output.

```
pscore_control_logit <- nsw_dw_cpscontrol_logit %>%
  filter(treat==0) %>%
  pull(pscore) %>%
  mean()
pscore_treat_logit <- nsw_dw_cpscontrol_logit %>%
  filter(treat==1) %>%
  pull(pscore) %>%
  mean()

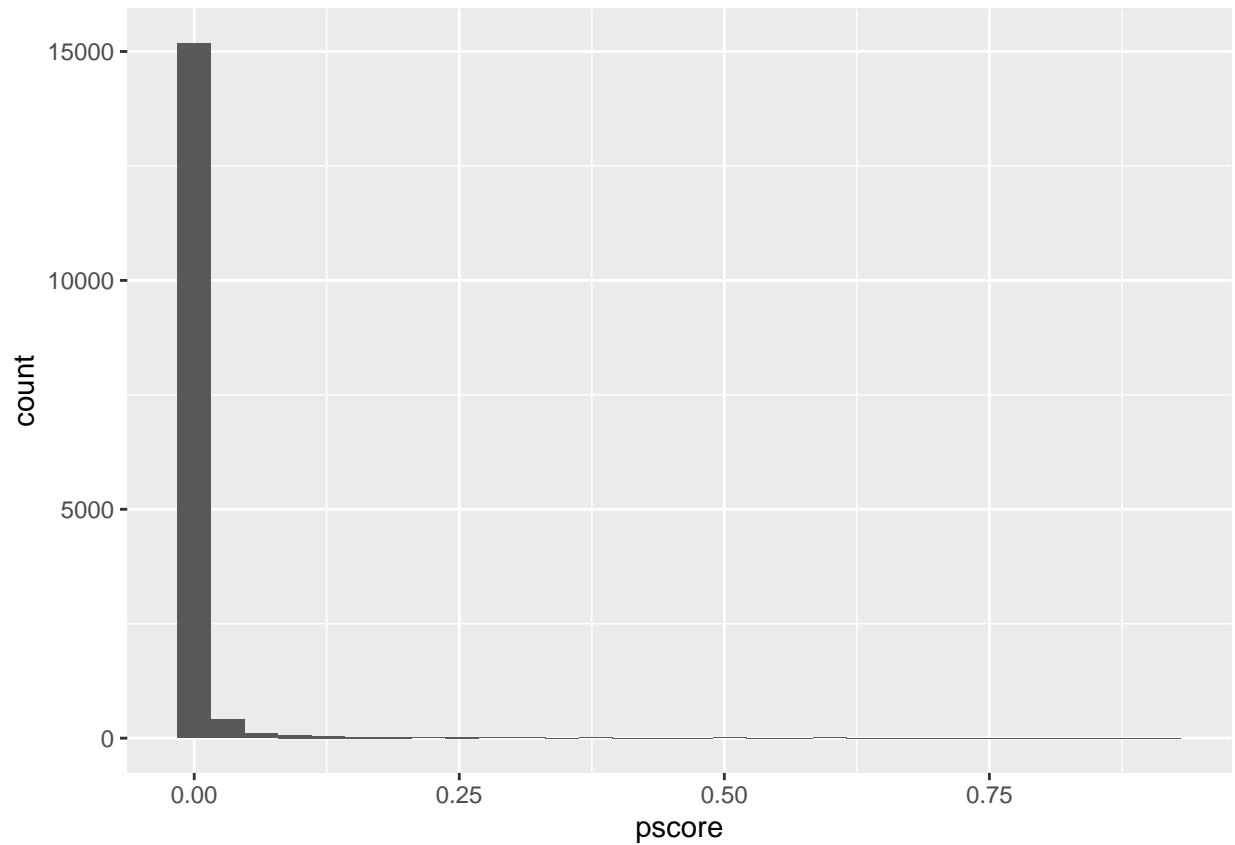
pscore_control_ols <- nsw_dw_cpscontrol_ols %>%
  filter(treat==0) %>%
  pull(pscore) %>%
  mean()
pscore_treat_ols <- nsw_dw_cpscontrol_ols %>%
  filter(treat==1) %>%
  pull(pscore) %>%
  mean()
```

##	Groups	Means
## 1	Logit Control	0.006582229
## 2	Logit Treat	0.431010791
## 3	OLS Control	0.009867421
## 4	OLS Treat	0.147028133

To see further into what the p-score distribution looks like for the treated and control groups of each model the following histograms provide a more easily understood visual representation.

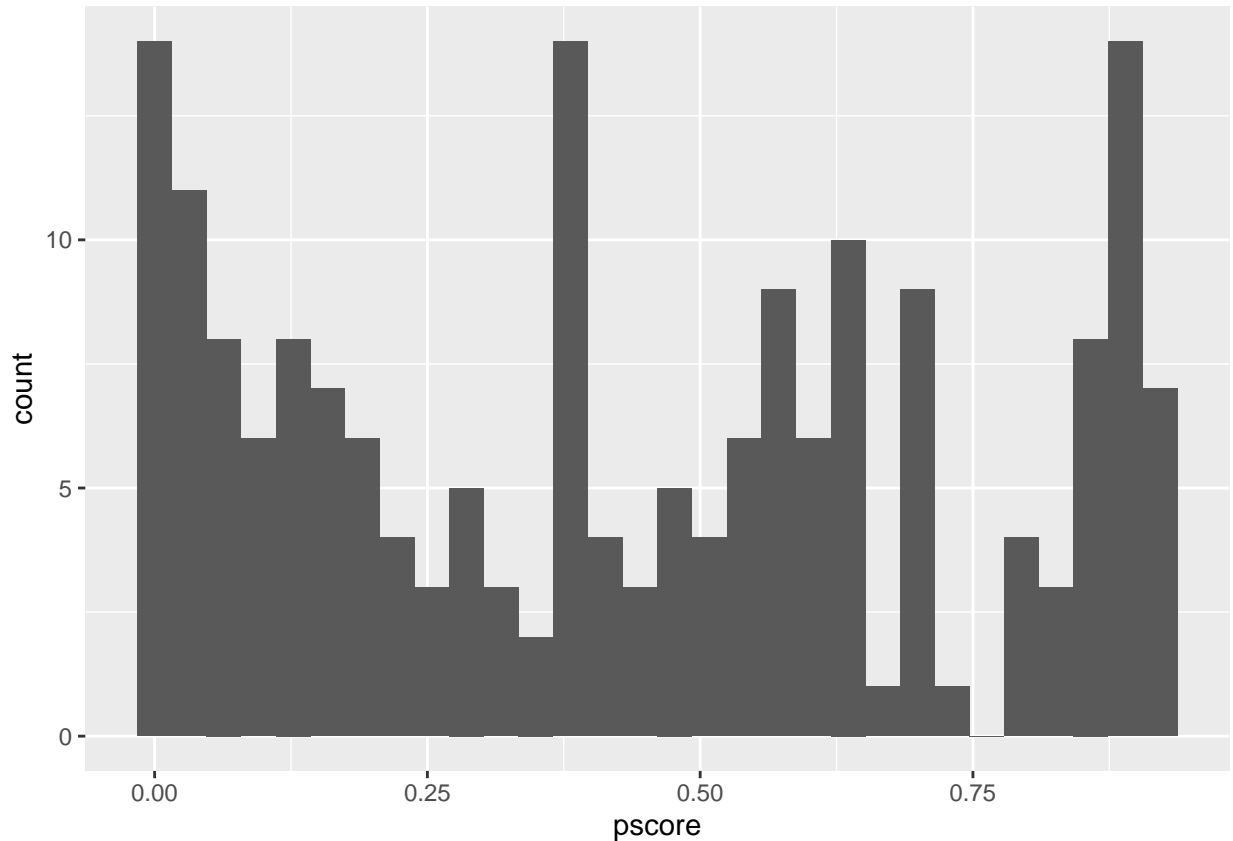
```
nsw_dw_cpscontrol_logit %>%
  filter(treat == 0) %>%
  ggplot() +
  geom_histogram(aes(x = pscore))
```

```
## 'stat_bin()' using 'bins = 30'. Pick better value with 'binwidth'.
```



```
nsw_dw_cpscontrol_logit %>%  
  filter(treat == 1) %>%  
  ggplot() +  
  geom_histogram(aes(x = pscore))
```

```
## 'stat_bin()' using 'bins = 30'. Pick better value with 'binwidth'.
```

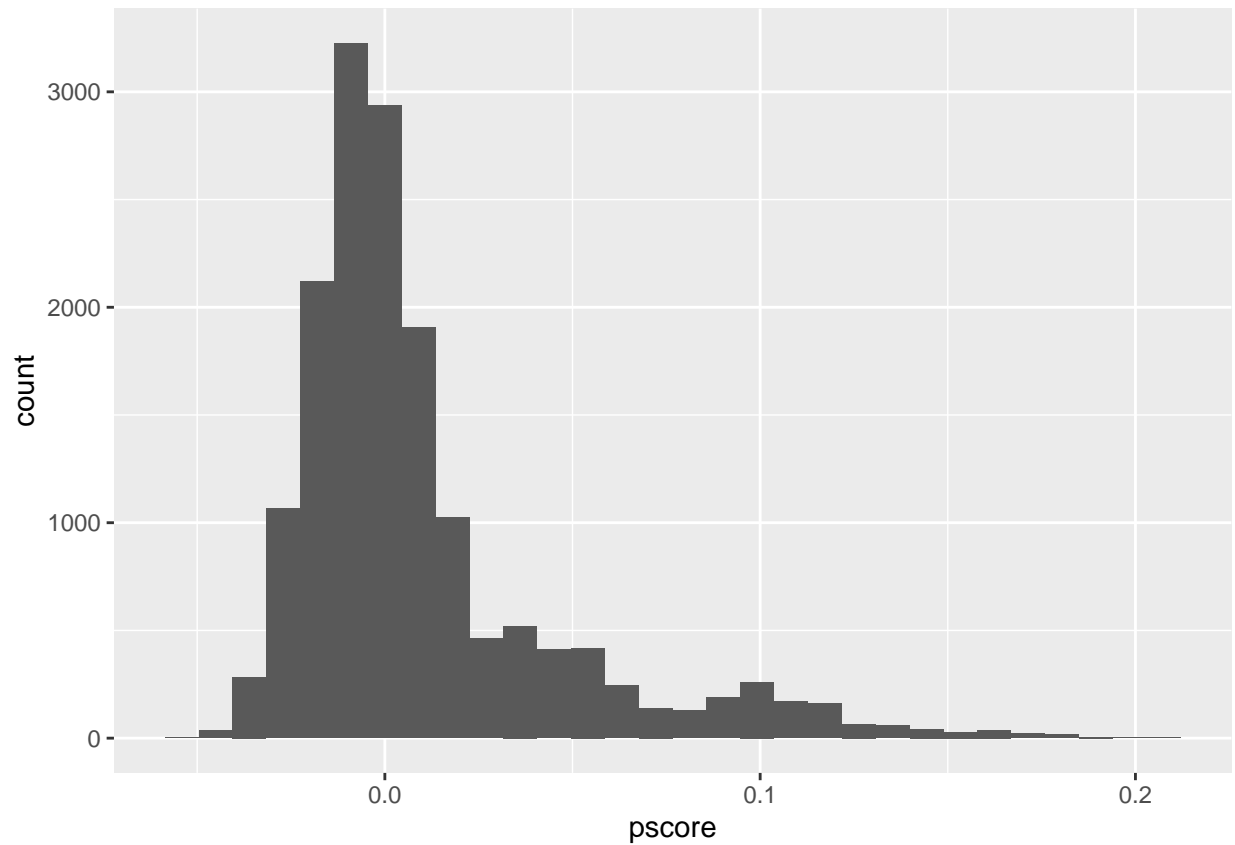


The two histograms above represent first the control p-score distribution and the treated distribution. The control group for the Logit model is highly skewed to the left. There is a vast amount of bunching around 0. Where the minimum value appears to go below 0 but really sit right at 0. The maximum value is around the 0.25. The treatment p-score distribution is far more evenly distributed across the frame. The histogram ranges from just below 0 to 1 and is distributed across the range in a non-uniform manner.

The following provides the histograms for the treated and control OLS model p-scores.

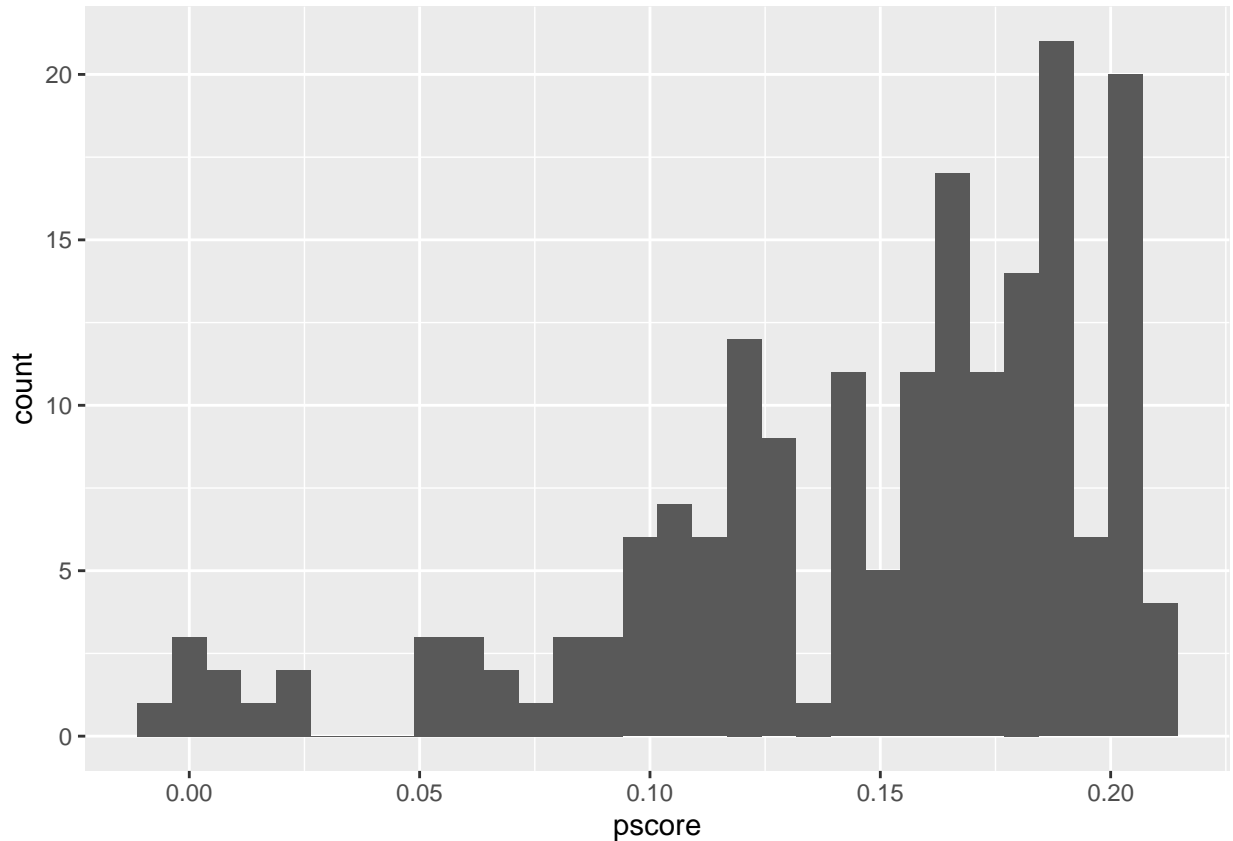
```
nsw_dw_cpscontrol_ols %>%
  filter(treat == 0) %>%
  ggplot() +
  geom_histogram(aes(x = pscore))
```

```
## 'stat_bin()' using 'bins = 30'. Pick better value with 'binwidth'.
```



```
nsw_dw_cpscontrol_ols %>%  
  filter(treat == 1) %>%  
  ggplot() +  
  geom_histogram(aes(x = pscore))
```

```
## 'stat_bin()' using 'bins = 30'. Pick better value with 'binwidth'.
```



The first histogram represents the control distribution. This is centered around 0 and has a loosely normal shaped distribution. There are far more values below 0 in this histogram than the prior histograms. The minimum appears to be around -0.05 and the maximum is just below 0.2. The second histogram is the treatment group distribution. This distribution appears to be skewed to the right side. The maximum values appear to be just above 0.2 (around 0.21) while the minimum is just below 0. The distribution is more weighted toward the higher end of the distribution with a peak around 0.18. ## d. The p-scores which are less than 0.1 and greater than 0.9 can then be dropped and the histograms recalibrated with the trimmed p-scores. The following will go through the same process as in part c. Additionally, the total number of observations in this trimmed dataset for both Logit and OLS models will be stored for later to use in the weighting for the ATT. Once the trimmed datasets are calculated they will be used for the rest of the problems. The first chunk of code trims the data and stores the number of observations (n) for each dataset.

```
nsw_dw_cpscontrol_logit <- nsw_dw_cpscontrol_logit %>%
  filter(!(pscore >= 0.9)) %>%
  filter(!(pscore <= 0.1))

Nl <- nrow(nsw_dw_cpscontrol_logit)

nsw_dw_cpscontrol_ols <- nsw_dw_cpscontrol_ols %>%
  filter(!(pscore >= 0.9)) %>%
  filter(!(pscore <= 0.1))

No <- nrow(nsw_dw_cpscontrol_ols)
```

The following chunk provides the code for recalculating the p-score means for the control and treatment group for OLS and Logit models.


```

pscore_control_logit2 <- nsw_dw_cpscontrol_logit %>%
  filter(treat==0) %>%
  pull(pscore) %>%
  mean()
pscore_treated_logit2 <- nsw_dw_cpscontrol_logit %>%
  filter(treat==1) %>%
  pull(pscore) %>%
  mean()

pscore_control_ols2 <- nsw_dw_cpscontrol_ols %>%
  filter(treat==0) %>%
  pull(pscore) %>%
  mean()
pscore_treated_ols2 <- nsw_dw_cpscontrol_ols %>%
  filter(treat==1) %>%
  pull(pscore) %>%
  mean()

```

The means are then provided in the following table (as previously provided). All of the means are a bit larger than the prior means from the first table. This change makes sense given the data now will not conclude the lower values which was weighting down the distribution in majority of the groups. Only the treatment group of the Logit model had a distribution which was near the higher end of the trimming (0.9), considering that and the many negative or close to 0 p-score's are gone, then the increase in these values is justified.

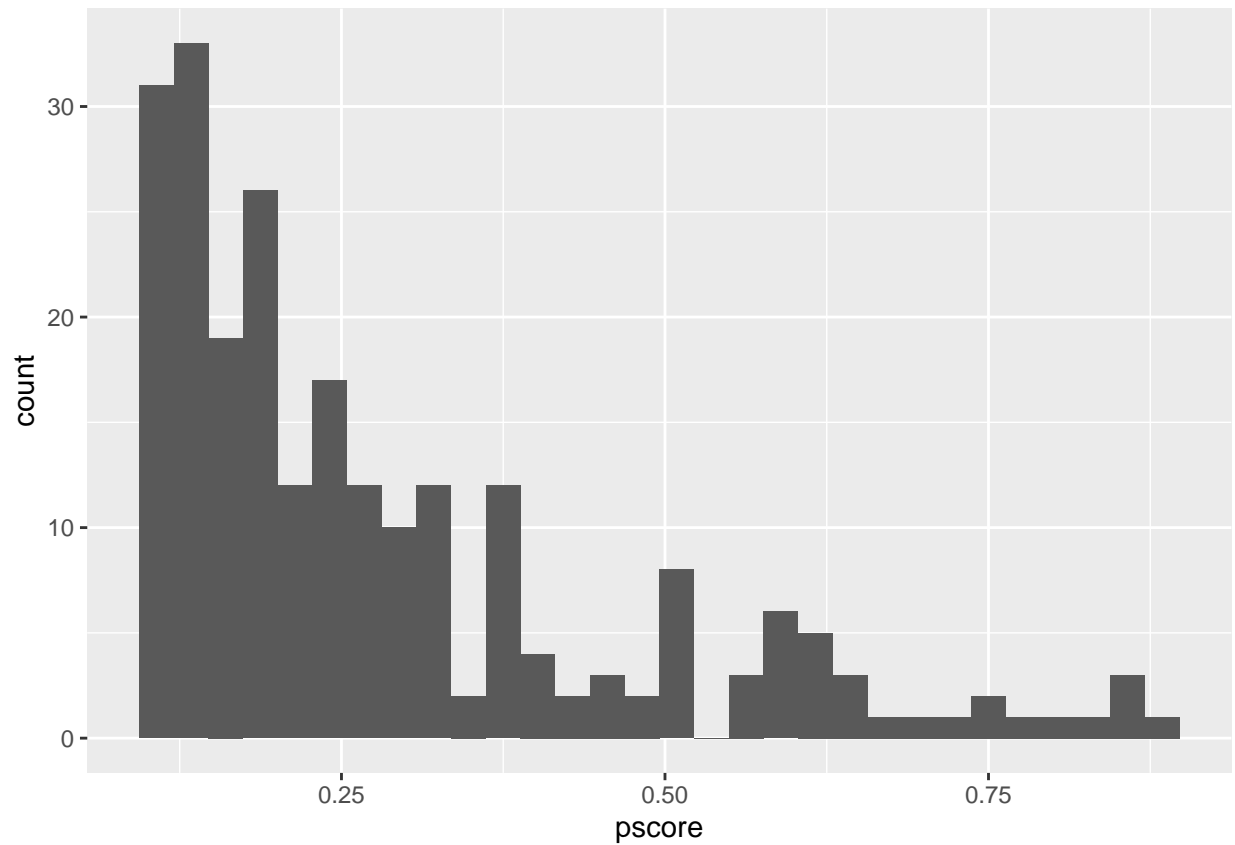
##	Groups	Means
## 1	Logit Control	0.2825155
## 2	Logit Treat	0.5112145
## 3	OLS Control	0.1249108
## 4	OLS Treat	0.1635168

```

nsw_dw_cpscontrol_logit %>%
  filter(treat == 0) %>%
  ggplot() +
  geom_histogram(aes(x = pscore))

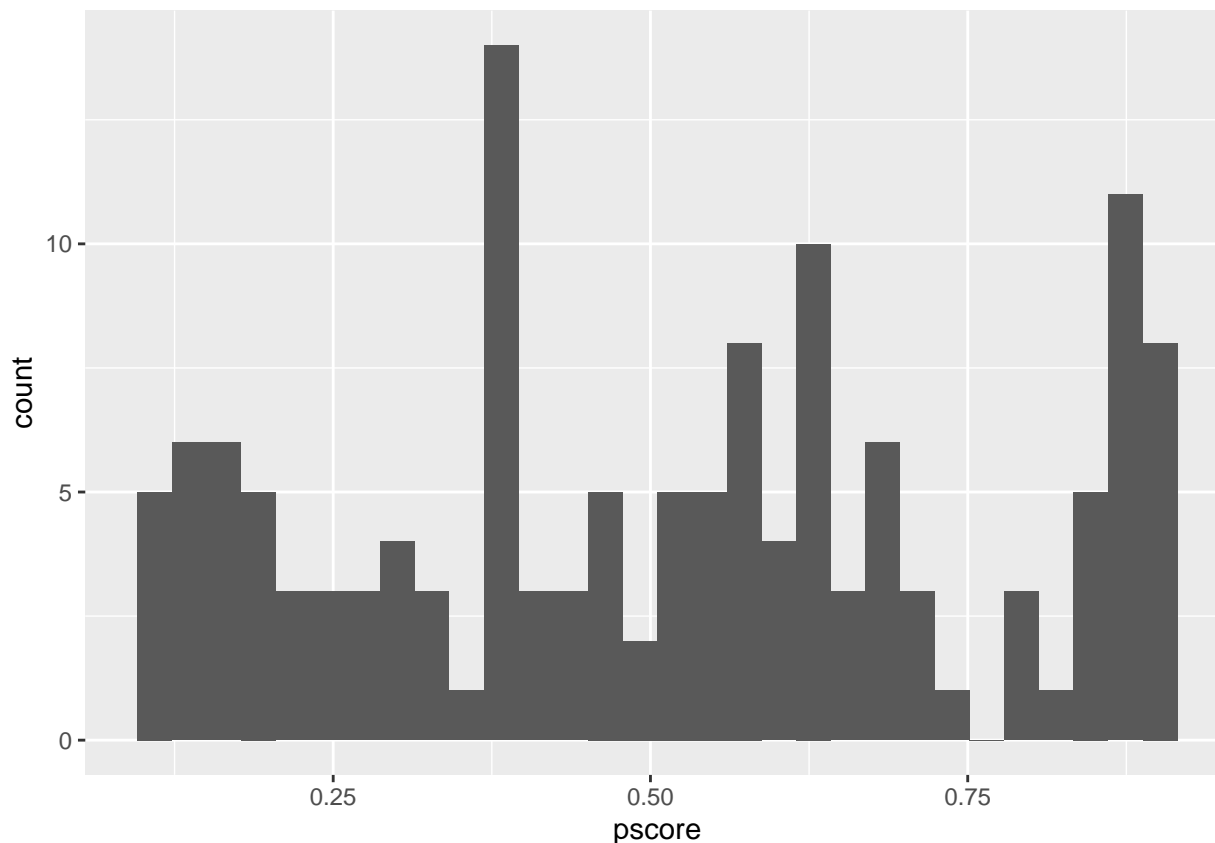
```

'stat_bin()' using 'bins = 30'. Pick better value with 'binwidth'.



```
nsw_dw_cpscontrol_logit %>%  
  filter(treat == 1) %>%  
  ggplot() +  
  geom_histogram(aes(x = pscore))
```

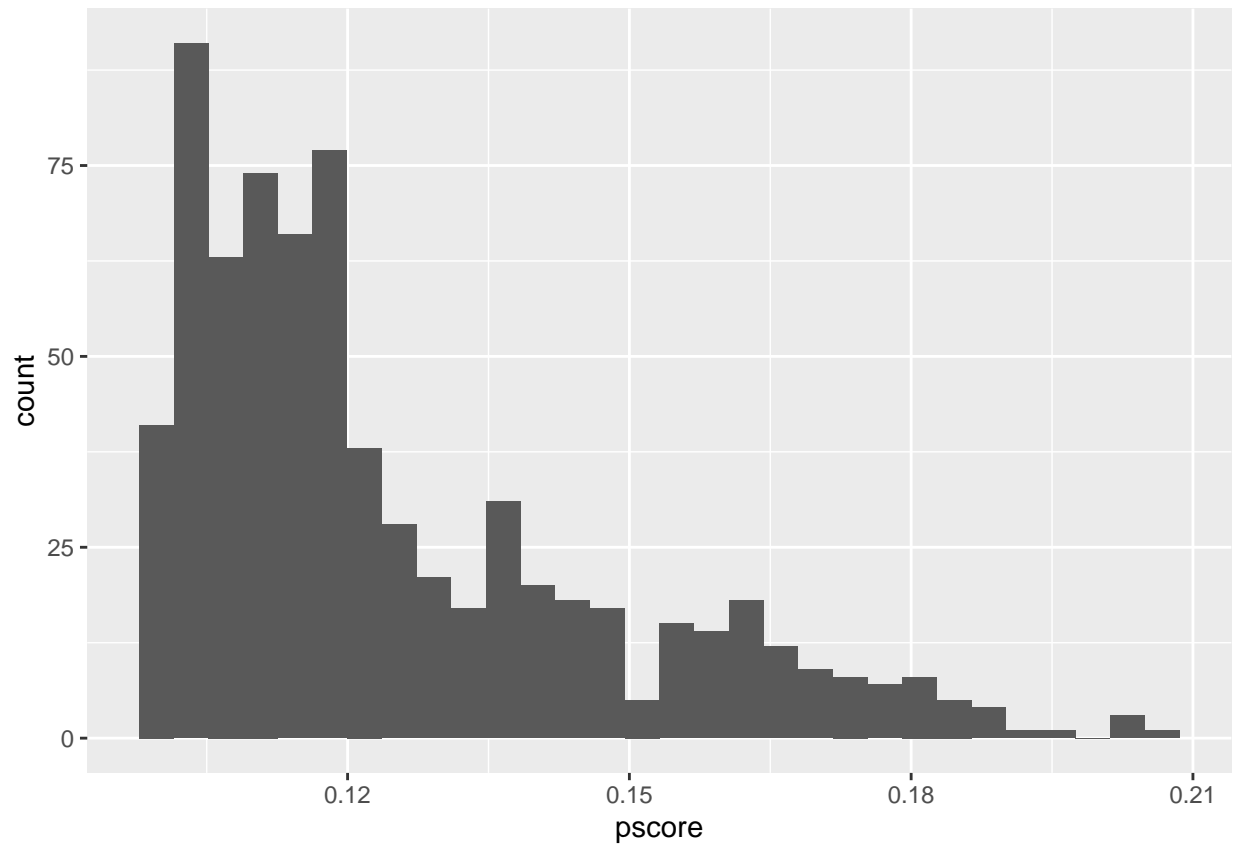
```
## 'stat_bin()' using 'bins = 30'. Pick better value with 'binwidth'.
```



The histograms above represent the trimmed version of the p-score distributions which were generated by the Logit model. This changes the distribution of the histograms from the untrimmed version of the data. Now the minimum values for both the control and treatment groups are 0.1 and maximums of 0.9. For the treatment group the distribution looks very similar. Just values at the very ends of the fairly even distribution have been cut therefore it would make sense that the histogram looks similar. However, the control group distribution is significantly different than the untrimmed version. This is due to the fact that the previous was surrounding right where the cutoff of the trimming exists. Therefore a lot of items were dropped. The distribution is skewed in the same direction but less extreme.

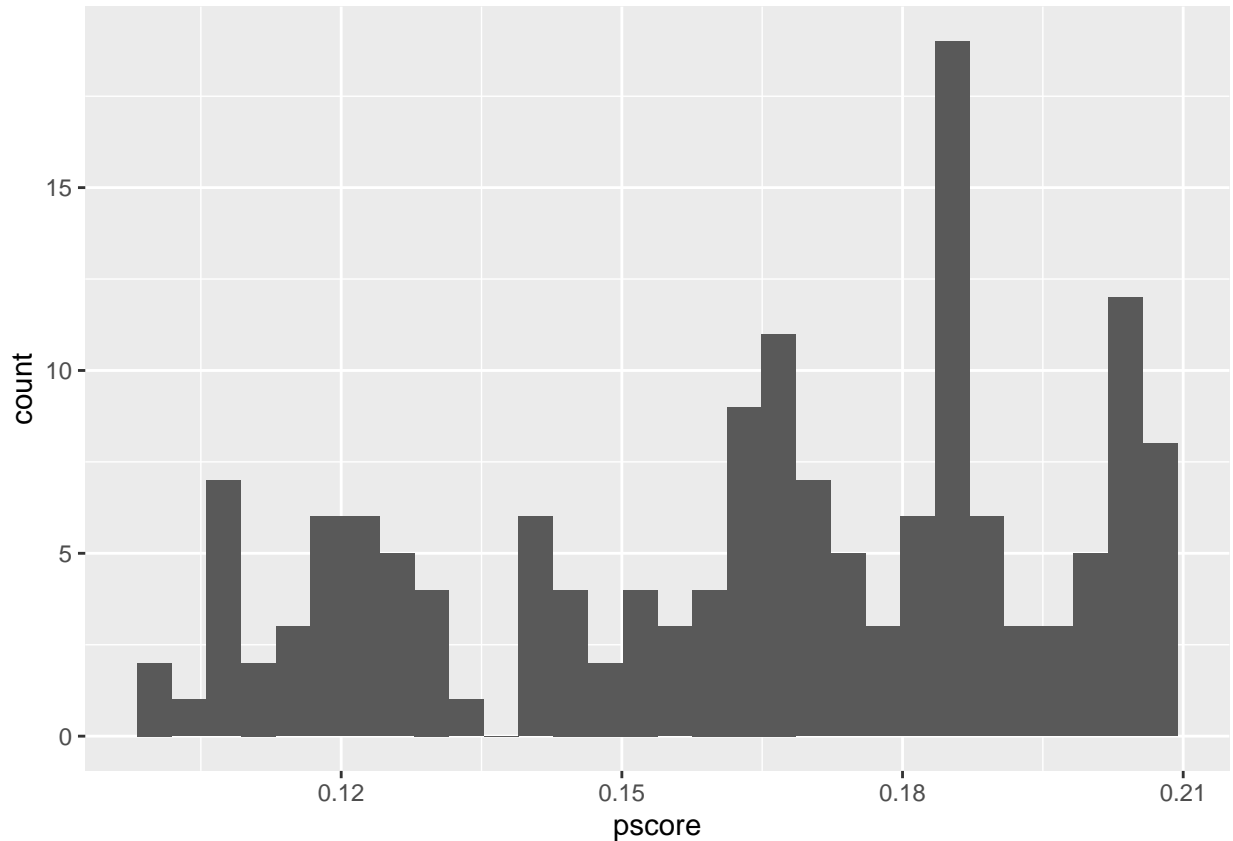
```
nsw_dw_cpscontrol_ols %>%  
  filter(treat == 0) %>%  
  ggplot() +  
  geom_histogram(aes(x = pscore))
```

```
## 'stat_bin()' using 'bins = 30'. Pick better value with 'binwidth'.
```



```
nsw_dw_cpscontrol_ols %>%  
  filter(treat == 1) %>%  
  ggplot() +  
  geom_histogram(aes(x = pscore))
```

```
## 'stat_bin()' using 'bins = 30'. Pick better value with 'binwidth'.
```



The histograms above represent the trimmed OLS control and treatment group p-score distributions. The control group now has a minimum of 0.1 and a maximum of just below 0.21. This distribution is skewed at the lower end. This makes sense with the process of the trimming. Much of the data was at or below the 0.1 cutoff. Therefore, much of the values in the control group were dropped with the trimming. The shape of the distribution remains the same. The treatment group is now has a minimum value of 0.1 and a maximum value of just below 0.21 similar to the control group. The untrimmed version was skewed to the higher end, while this distribution is a little more spread across the range. #2 The first differencing compares the outcomes of the trimmed data provides the ATE output. This provides the first difference.

```
nsw_dw_cpscontrol_logit %>%
  filter(treat == 1) %>%
  summary(re78)
```

```
##      data_id      treat      age      educ
## Length:139      Min.   :1      Min.   :17.00      Min.   : 4.00
## Class :character 1st Qu.:1      1st Qu.:20.00      1st Qu.: 9.00
## Mode  :character Median :1      Median :25.00      Median :11.00
##                      Mean  :1      Mean  :25.31      Mean  :10.27
##                      3rd Qu.:1      3rd Qu.:28.00      3rd Qu.:11.50
##                      Max.   :1      Max.   :48.00      Max.   :15.00
##      black      hisp      marr      nodegree
## Min.   :0.0000      Min.   :0.00000      Min.   :0.0000      Min.   :0.0000
## 1st Qu.:1.0000      1st Qu.:0.00000      1st Qu.:0.0000      1st Qu.:0.5000
## Median :1.0000      Median :0.00000      Median :0.0000      Median :1.0000
## Mean   :0.9424      Mean   :0.03597      Mean   :0.1439      Mean   :0.7482
## 3rd Qu.:1.0000      3rd Qu.:0.00000      3rd Qu.:0.0000      3rd Qu.:1.0000
```

```
## Max. :1.0000 Max. :1.00000 Max. :1.0000 Max. :1.0000
## re74 re75 re78 agesq
## Min. : 0 Min. : 0 Min. : 0 Min. : 289.0
## 1st Qu.: 0 1st Qu.: 0 1st Qu.: 0 1st Qu.: 400.0
## Median : 0 Median : 0 Median : 4056 Median : 625.0
## Mean : 1815 Mean : 1108 Mean : 6545 Mean : 683.8
## 3rd Qu.: 0 3rd Qu.: 1682 3rd Qu.: 9621 3rd Qu.: 784.0
## Max. :35040 Max. :11537 Max. :60308 Max. :2304.0
## agecube educsq u74 u75
## Min. : 4913 Min. : 16.0 Min. :0.0000 Min. :0.0000
## 1st Qu.: 8000 1st Qu.: 81.0 1st Qu.:1.0000 1st Qu.:0.0000
## Median : 15625 Median :121.0 Median :1.0000 Median :1.0000
## Mean : 19839 Mean :109.4 Mean :0.7626 Mean :0.6403
## 3rd Qu.: 21952 3rd Qu.:132.5 3rd Qu.:1.0000 3rd Qu.:1.0000
## Max. :110592 Max. :225.0 Max. :1.0000 Max. :1.0000
## interaction1 re74sq re75sq interaction2
## Min. : 0 Min. :0.000e+00 Min. : 0 Min. :0.00000
## 1st Qu.: 0 1st Qu.:0.000e+00 1st Qu.: 0 1st Qu.:0.00000
## Median : 0 Median :0.000e+00 Median : 0 Median :0.00000
## Mean : 20033 Mean :2.644e+07 Mean : 5333070 Mean :0.03597
## 3rd Qu.: 0 3rd Qu.:0.000e+00 3rd Qu.: 2830599 3rd Qu.:0.00000
## Max. :490561 Max. :1.228e+09 Max. :133092455 Max. :1.00000
## educcube re74cube re75cube pscore
## Min. : 64 Min. :0.000e+00 Min. :0.000e+00 Min. :0.1063
## 1st Qu.: 729 1st Qu.:0.000e+00 1st Qu.:0.000e+00 1st Qu.:0.3180
## Median :1331 Median :0.000e+00 Median :0.000e+00 Median :0.5276
## Mean :1194 Mean :5.791e+11 Mean :3.478e+10 Mean :0.5112
## 3rd Qu.:1530 3rd Qu.:0.000e+00 3rd Qu.:4.763e+09 3rd Qu.:0.6887
## Max. :3375 Max. :4.302e+13 Max. :1.535e+12 Max. :0.8985
```

```
meanllog <- nsw_dw_cpscontrol_logit %>%
  filter(treat == 1) %>%
  pull(re78) %>%
  mean()

nsw_dw_cpscontrol_logit$y1 <- meanllog

nsw_dw_cpscontrol_logit %>%
  filter(treat == 0) %>%
  summary(re78)
```

```
## data_id treat age educ
## Length:234 Min. :0 Min. :16.00 Min. : 2.00
## Class :character 1st Qu.:0 1st Qu.:19.00 1st Qu.: 9.00
## Mode :character Median :0 Median :24.00 Median :11.00
## Mean :0 Mean :25.67 Mean :10.37
## 3rd Qu.:0 3rd Qu.:31.00 3rd Qu.:12.00
## Max. :0 Max. :55.00 Max. :16.00
## black hisp marr nodegree
## Min. :0.0000 Min. :0.00000 Min. :0.0000 Min. :0.0000
## 1st Qu.:1.0000 1st Qu.:0.00000 1st Qu.:0.0000 1st Qu.:0.0000
## Median :1.0000 Median :0.00000 Median :0.0000 Median :1.0000
## Mean :0.9231 Mean :0.05556 Mean :0.2265 Mean :0.6453
## 3rd Qu.:1.0000 3rd Qu.:0.00000 3rd Qu.:0.0000 3rd Qu.:1.0000
```

```
## Max. :1.0000 Max. :1.00000 Max. :1.0000 Max. :1.0000
## re74 re75 re78 agesq
## Min. : 0 Min. : 0.0 Min. : 0 Min. : 256.0
## 1st Qu.: 0 1st Qu.: 0.0 1st Qu.: 0 1st Qu.: 361.0
## Median : 0 Median : 177.2 Median : 2404 Median : 576.0
## Mean : 2433 Mean : 1709.0 Mean : 4821 Mean : 720.2
## 3rd Qu.: 3489 3rd Qu.: 2896.3 3rd Qu.: 7976 3rd Qu.: 961.0
## Max. :22322 Max. :13004.9 Max. :25565 Max. :3025.0
## agecube educsq u74 u75
## Min. : 4096 Min. : 4.0 Min. :0.0000 Min. :0.0000
## 1st Qu.: 6859 1st Qu.: 81.0 1st Qu.:0.0000 1st Qu.:0.0000
## Median : 13824 Median :121.0 Median :1.0000 Median :0.0000
## Mean : 22115 Mean :112.4 Mean :0.5598 Mean :0.4658
## 3rd Qu.: 29791 3rd Qu.:144.0 3rd Qu.:1.0000 3rd Qu.:1.0000
## Max. :166375 Max. :256.0 Max. :1.0000 Max. :1.0000
## interaction1 re74sq re75sq interaction2
## Min. : 0 Min. : 0 Min. : 0 Min. :0.00000
## 1st Qu.: 0 1st Qu.: 0 1st Qu.: 0 1st Qu.:0.00000
## Median : 0 Median : 0 Median : 31803 Median :0.00000
## Mean : 24443 Mean : 22658379 Mean : 9272917 Mean :0.02991
## 3rd Qu.: 35041 3rd Qu.: 12173655 3rd Qu.: 8388694 3rd Qu.:0.00000
## Max. :267863 Max. :498268545 Max. :169127434 Max. :1.00000
## educcube re74cube re75cube pscore
## Min. : 8 Min. :0.000e+00 Min. :0.000e+00 Min. :0.1001
## 1st Qu.: 729 1st Qu.:0.000e+00 1st Qu.:0.000e+00 1st Qu.:0.1416
## Median :1331 Median :0.000e+00 Median :5.774e+06 Median :0.2217
## Mean :1261 Mean :2.845e+11 Mean :6.500e+10 Mean :0.2825
## 3rd Qu.:1728 3rd Qu.:4.248e+10 3rd Qu.:2.430e+10 3rd Qu.:0.3666
## Max. :4096 Max. :1.112e+13 Max. :2.199e+12 Max. :0.8771
## y1
## Min. :6545
## 1st Qu.:6545
## Median :6545
## Mean :6545
## 3rd Qu.:6545
## Max. :6545
```

```
mean0log <- nsw_dw_cpscontrol_logit %>%
  filter(treat == 0) %>%
  pull(re78) %>%
  mean()

nsw_dw_cpscontrol_logit$y0 <- mean0log

atelog <- unique(nsw_dw_cpscontrol_logit$y1 - nsw_dw_cpscontrol_logit$y0)
```

The following is for the Logit model.

```
## [1] 1724.773
```

```
nsw_dw_cpscontrol_ols %>%
  filter(treat == 1) %>%
  summary(re78)
```

```
##      data_id      treat      age      educ
## Length:157      Min.   :1      Min.   :17.00      Min.   : 4.00
## Class :character 1st Qu.:1      1st Qu.:20.00      1st Qu.: 9.00
## Mode  :character Median :1      Median :25.00      Median :11.00
##                      Mean   :1      Mean   :25.97      Mean   :10.29
##                      3rd Qu.:1      3rd Qu.:29.00      3rd Qu.:12.00
##                      Max.    :1      Max.    :48.00      Max.    :16.00
##      black      hisp      marr      nodegree
## Min.   :0.0000      Min.   :0.00000      Min.   :0.000      Min.   :0.0000
## 1st Qu.:1.0000      1st Qu.:0.00000      1st Qu.:0.000      1st Qu.:0.0000
## Median :1.0000      Median :0.00000      Median :0.000      Median :1.0000
## Mean   :0.9809      Mean   :0.01274      Mean   :0.172      Mean   :0.7389
## 3rd Qu.:1.0000      3rd Qu.:0.00000      3rd Qu.:0.000      3rd Qu.:1.0000
## Max.   :1.0000      Max.   :1.00000      Max.   :1.000      Max.   :1.0000
##      re74      re75      re78      agesq
## Min.   : 0.0      Min.   : 0      Min.   : 0      Min.   : 289.0
## 1st Qu.: 0.0      1st Qu.: 0      1st Qu.: 0      1st Qu.: 400.0
## Median : 0.0      Median : 0      Median : 4033      Median : 625.0
## Mean   : 1958.6      Mean   : 1375      Mean   : 6348      Mean   : 726.8
## 3rd Qu.: 989.3      3rd Qu.: 1666      3rd Qu.: 9599      3rd Qu.: 841.0
## Max.   :35040.1      Max.   :25142      Max.   :60308      Max.   :2304.0
##      agecube      educsq      u74      u75
## Min.   : 4913      Min.   : 16      Min.   :0.0000      Min.   :0.0000
## 1st Qu.: 8000      1st Qu.: 81      1st Qu.:0.0000      1st Qu.:0.0000
## Median : 15625      Median :121      Median :1.0000      Median :1.0000
## Mean   : 22022      Mean   :110      Mean   :0.7261      Mean   :0.6306
## 3rd Qu.: 24389      3rd Qu.:144      3rd Qu.:1.0000      3rd Qu.:1.0000
## Max.   :110592      Max.   :256      Max.   :1.0000      Max.   :1.0000
##      interaction1      re74sq      re75sq      interaction2
## Min.   : 0      Min.   :0.000e+00      Min.   : 0      Min.   :0.00000
## 1st Qu.: 0      1st Qu.:0.000e+00      1st Qu.: 0      1st Qu.:0.00000
## Median : 0      Median :0.000e+00      Median : 0      Median :0.00000
## Mean   : 21095      Mean   :2.736e+07      Mean   : 11960537      Mean   :0.01274
## 3rd Qu.: 10012      3rd Qu.:9.787e+05      3rd Qu.: 2775933      3rd Qu.:0.00000
## Max.   :490561      Max.   :1.228e+09      Max.   :632132244      Max.   :1.00000
##      pscore
## Min.   :0.1007
## 1st Qu.:0.1414
## Median :0.1677
## Mean   :0.1635
## 3rd Qu.:0.1866
## Max.   :0.2081
```

```
mean1ols <- nsw_dw_cpscontrol_ols %>%
  filter(treat == 1) %>%
  pull(re78) %>%
  mean()

nsw_dw_cpscontrol_ols$y1 <- mean1ols

nsw_dw_cpscontrol_ols %>%
  filter(treat == 0) %>%
  summary(re78)
```



```
##      data_id      treat      age      educ      black
## Length:713      Min.    :0      Min.    :16.0      Min.    : 0.00      Min.    :0.0000
## Class :character 1st Qu.:0      1st Qu.:23.0      1st Qu.: 9.00      1st Qu.:1.0000
## Mode  :character Median :0      Median :29.0      Median :11.00      Median :1.0000
##                      Mean    :0      Mean    :30.7      Mean    :10.97      Mean    :0.9972
##                      3rd Qu.:0      3rd Qu.:37.0      3rd Qu.:12.00      3rd Qu.:1.0000
##                      Max.    :0      Max.    :55.0      Max.    :18.00      Max.    :1.0000
##      hisp      marr      nodegree      re74
## Min.    :0      Min.    :0.0000      Min.    :0.0000      Min.    : 0.0
## 1st Qu.:0      1st Qu.:0.0000      1st Qu.:0.0000      1st Qu.: 311.5
## Median :0      Median :0.0000      Median :1.0000      Median : 7421.7
## Mean    :0      Mean    :0.4516      Mean    :0.5203      Mean    : 9674.4
## 3rd Qu.:0      3rd Qu.:1.0000      3rd Qu.:1.0000      3rd Qu.:17737.2
## Max.    :0      Max.    :1.0000      Max.    :1.0000      Max.    :25862.3
##      re75      re78      agesq      agecube
## Min.    : 0      Min.    : 0.0      Min.    : 256      Min.    : 4096
## 1st Qu.: 0      1st Qu.: 325.1      1st Qu.: 529      1st Qu.: 12167
## Median : 5799      Median : 8378.7      Median : 841      Median : 24389
## Mean    : 8937      Mean    :10476.5      Mean    :1048      Mean    : 39374
## 3rd Qu.:16859      3rd Qu.:19046.4      3rd Qu.:1369      3rd Qu.: 50653
## Max.    :25244      Max.    :25564.7      Max.    :3025      Max.    :166375
##      educsq      u74      u75      interaction1
## Min.    : 0.0      Min.    :0.0000      Min.    :0.0000      Min.    : 0
## 1st Qu.: 81.0      1st Qu.:0.0000      1st Qu.:0.0000      1st Qu.: 2939
## Median :121.0      Median :0.0000      Median :0.0000      Median : 77838
## Mean    :127.1      Mean    :0.2342      Mean    :0.2539      Mean    :106063
## 3rd Qu.:144.0      3rd Qu.:0.0000      3rd Qu.:1.0000      3rd Qu.:176334
## Max.    :324.0      Max.    :1.0000      Max.    :1.0000      Max.    :439659
##      re74sq      re75sq      interaction2      pscore
## Min.    : 0      Min.    : 0      Min.    :0      Min.    :0.1000
## 1st Qu.: 97047      1st Qu.: 0      1st Qu.:0      1st Qu.:0.1078
## Median : 55081677      Median : 33626719      Median :0      Median :0.1172
## Mean    :179485834      Mean    :162181377      Mean    :0      Mean    :0.1249
## 3rd Qu.:314609691      3rd Qu.:284241752      3rd Qu.:0      3rd Qu.:0.1378
## Max.    :668859612      Max.    :637236856      Max.    :0      Max.    :0.2071
##      y1
## Min.    :6348
## 1st Qu.:6348
## Median :6348
## Mean    :6348
## 3rd Qu.:6348
## Max.    :6348
```

```
mean0ols <- nsw_dw_cpscontrol_ols %>%
  filter(treat == 0) %>%
  pull(re78) %>%
  mean()

nsw_dw_cpscontrol_ols$y0 <- mean0ols

ateols <- unique(nsw_dw_cpscontrol_ols$y1 - nsw_dw_cpscontrol_ols$y0)
```

The following is for the OLS model.

```
## [1] -4128.57
```

#3 The following is the code syntax for the Logit model weighted with normalized and non normalized ATT calculations. The equation provided for the Abadie (2005) estimations is the non-normalized weights. This also uses the trimmed data.

```
nsw_dw_cpscontrol_logit <- nsw_dw_cpscontrol_logit %>%
  mutate(d1 = treat/pscore,
         d0 = (1-treat)/(1-pscore))

s1 <- sum(nsw_dw_cpscontrol_logit$d1)
s0 <- sum(nsw_dw_cpscontrol_logit$d0)
# Non-Normalized Weights
nsw_dw_cpscontrol_logit <- nsw_dw_cpscontrol_logit %>%
  mutate(y1 = treat * re78/pscore,
         y0 = (1-treat) * re78/(1-pscore),
         ht = y1 - y0)
#Normalized Weights
nsw_dw_cpscontrol_logit <- nsw_dw_cpscontrol_logit %>%
  mutate(y1 = (treat*re78/pscore)/(s1/N1),
         y0 = ((1-treat)*re78/(1-pscore))/(s0/N1),
         norm = y1 - y0)

nsw_dw_cpscontrol_logit %>%
  pull(ht) %>%
  mean()
```

```
## [1] 1744.356
```

```
nsw_dw_cpscontrol_logit %>%
  pull(norm) %>%
  mean()
```

```
## [1] 1635.115
```

The first value is the non-normalized weight ATT, which is the version of interest. But the normalized weights are also provided.

The following is the OLS model of the normalized and non-normalized weighted ATT calculation. Abadie (2005) is the non-normalized version.

```
nsw_dw_cpscontrol_ols <- nsw_dw_cpscontrol_ols %>%
  mutate(d1 = treat/pscore,
         d0 = (1-treat)/(1-pscore))

s1 <- sum(nsw_dw_cpscontrol_ols$d1)
s0 <- sum(nsw_dw_cpscontrol_ols$d0)

#Non-normalized weights
nsw_dw_cpscontrol_ols <- nsw_dw_cpscontrol_ols %>%
  mutate(y1 = treat * re78/pscore,
         y0 = (1-treat) * re78/(1-pscore),
         ht = y1 - y0)
```

```

#Normalized weights
nsw_dw_cpscontrol_ols <- nsw_dw_cpscontrol_ols %>%
  mutate(y1 = (treat*re78/pscore)/(s1/No),
         y0 = ((1-treat)*re78/(1-pscore))/(s0/No),
         norm = y1 - y0)

nsw_dw_cpscontrol_ols %>%
  pull(ht) %>%
  mean()

```

```
## [1] -2322.159
```

```

nsw_dw_cpscontrol_ols %>%
  pull(norm) %>%
  mean()

```

```
## [1] -3938.138
```

The first is the OLS ATT (trimmed data) with non-normalized weights, as specified by Abadie (2005) and the second is the ATT with the normalized weights. The outcome of interest is the first from the non-normalized weights.

The estimates determined by this analysis are slightly different than the ones generated by Dr. Cunningham in the Mixtape due to the fact that this model is a little more complicated in that it includes additional covariates. However the estimates are close and match the direction (i.e. the logit model is positive). The OLS model provides negative estimates. This indicates that the Logit model is more effective in the analysis of this situation.