

Impact Evaluation: Assignment 2

Mukhammed Zainidinov

Women as policymakers

1. Import data file

```
library('tidyverse')  
women = read.csv('women.csv')
```

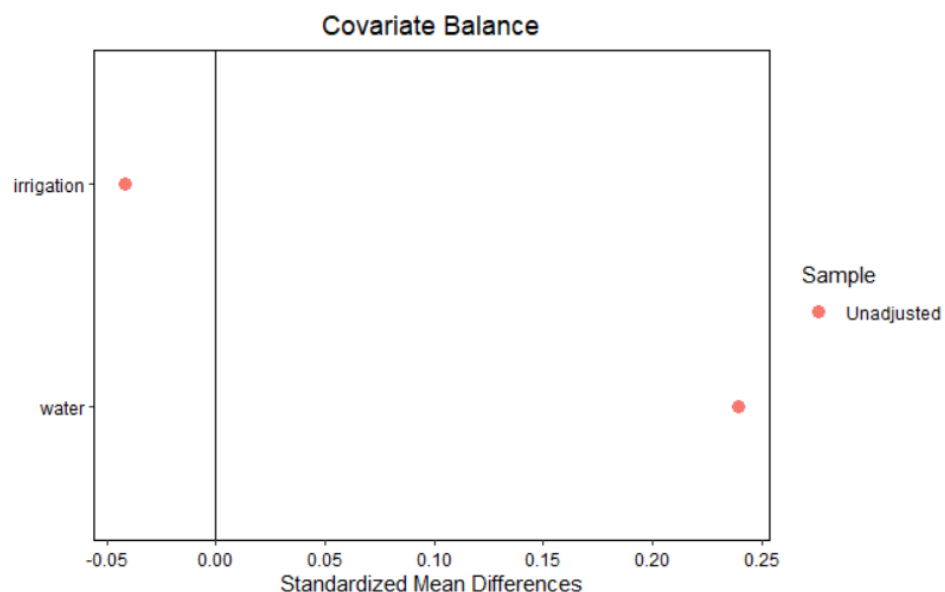
2. Summary statistics

```
glimpse(women)  
summary(women)
```

There are 161 Gram Panchayats in the study. According to the paper, two villages were randomly selected from each Gram Panchayats, hence the variable *village* takes two values – 1 and 2. Variable *reserved* is a treatment variable that takes values 0 for the control group, and 1 for the treated group. Variable *female* is dummy variables for gender. Variables *irrigation* and *water* explain the level of irrigation and access to drinking water, respectively.

3. Check for balance

```
library(cobalt)  
love.plot(reserved ~ irrigation + water, data = women, stars = 'std')
```

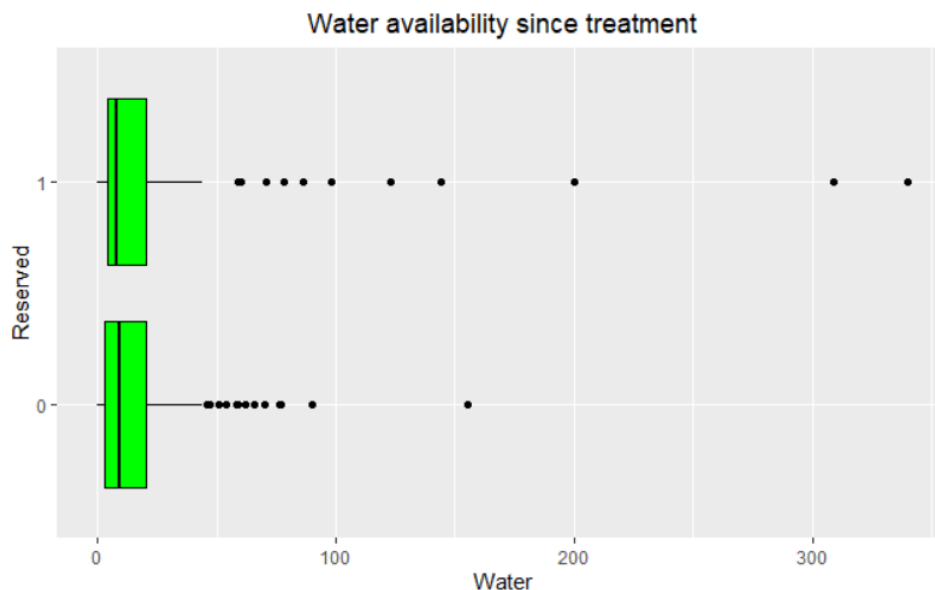


From the graph, we see that there are differences in mean between the control and treated groups meaning that covariates are not balanced. However, mean differences are not statistically significant for both covariates despite a small difference in *irrigation* and a large difference in *water*.

Because treatment is assigned randomly, *irrigation* and *water* should be balanced between the treated and control group, i.e. they should have the same mean.

4. Boxplot

```
ggplot(women, aes(y = water, x = factor(reserved))) +  
  geom_boxplot(fill = "green", color = "black") + coord_flip() +  
  labs(x = "Reserved", y = "Water") +  
  ggtitle("Water availability since treatment") +  
  theme(plot.title = element_text(hjust = 0.5))
```



The boxplot shows a lower median for the treated group. However, there are several outliers for the treated group that are greater than 200.

5. Regression analysis

```
mod_irr = lm(irrigation ~ reserved, data = women)  
mod_water = lm(water ~ reserved, data = women)  
library(texreg)  
screenreg(list(mod_irr, mod_water))
```

The treatment effect		
=====		
	Dependent variable:	

	irrigation	water
	(1)	(2)

reserved	-0.369	9.252**
	(1.122)	(3.948)
Constant	3.388***	14.738***
	(0.650)	(2.286)

Observations	322	322
R2	0.0003	0.017
Adjusted R2	-0.003	0.014
Residual Std. Error (df = 320)	9.506	33.446
F Statistic (df = 1; 320)	0.108	5.493**
=====		
Note:	*p<0.1; **p<0.05; ***p<0.01	

The treatment effect on *irrigation* is not statistically significant. That's why we cannot say whether there is the effect or no effect of reserving the GP head seat for women on irrigation. The treatment effect on *water* is statistically significant at a 95% confidence level and it is positive. We can conclude that the effect of reserving GP head seats for women positively affects drinking water availability by 9.252 units.

6. If the treatment was assigned randomly, then we shouldn't care about other covariates. When we do a random assignment, i.e. toss a coin, if it is head the head seat is reserved for woman, if it is tail the subject falls into control group, it doesn't care about the current state of the village, for example, level of irrigation, access to clean water, infrastructure, climate, number of population, average income per villager, etc. In this way, randomization removes omitted variable bias and unobserved variable issues.

If the treatment was not assigned randomly, then we should think about what covariates to include or how to change our model because R-squared is very low meaning that only very few variations in independent variables are explained by the treatment variable.

STAR experiment

1. Import data

```
star = read.csv('star.csv')
```

2. Average total score by small and regular class

```
library('dplyr')

star %>%

  filter(small == 1) %>%

  select(totalscore) %>%

  summary(star)

star %>%

  filter(regular == 1) %>%

  select(totalscore) %>%

  summary(star)
```

Summary statistics of the total score for small class shows that the average total score is 931.9, while for regular class average total score is 918. Students studying in small classes are expected to get higher scores than those studying in regular classes.

3. To generate boxplot to compare total score by small and regular class size, first we have to remove observations for regular class with aide.

```
star = star[!(star$aide == 1),] #drop observation if aide == 1

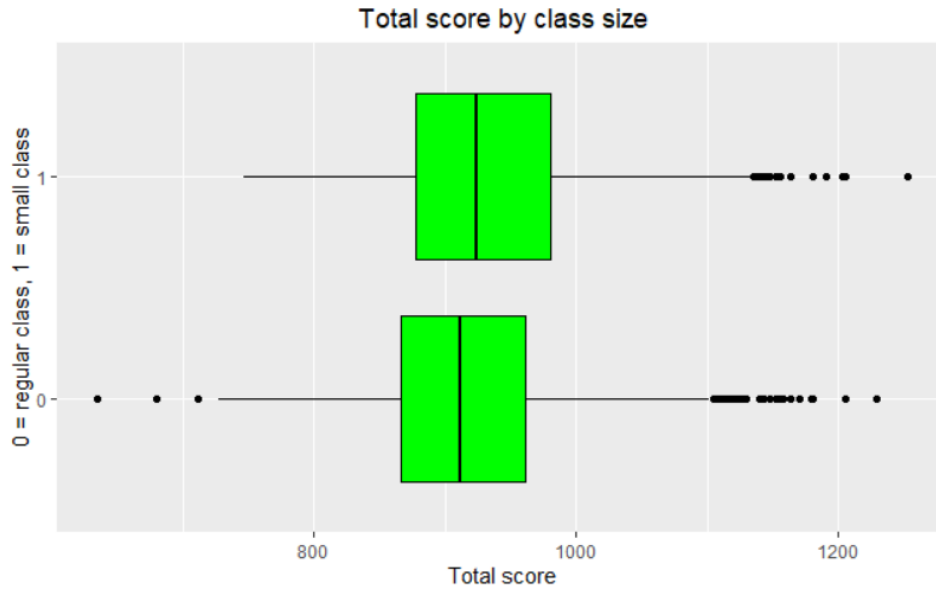
ggplot(star,aes(y = totalscore, x = factor(small))) +

  geom_boxplot(fill = "green", color="black") + coord_flip() +

  labs(x = "0 = regular class, 1 = small class", y = "Total score")+

  ggtitle("Total score by class size") +

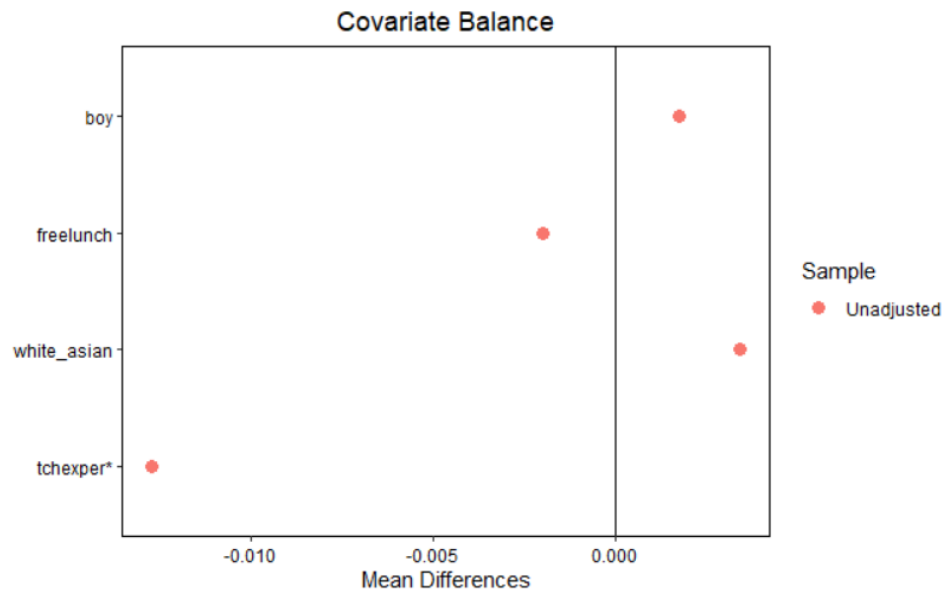
  theme(plot.title = element_text(hjust = 0.5))
```



There is a shift to the right in the distribution for the small classes on total scores. Median of small class is greater than the median of regular class. The interquartile range, i.e. middle 50% of the total score, for the small class is larger than for the regular class. Three outliers below the lower whisker for the regular class might contribute to a lower average total score for the regular class.

4. Covariate balance for treated and control groups

```
love.plot(small ~ boy + freelunch + white_asian + tchexper,
          data = star, stars = 'std')
```



The mean of the variables *boy*, *freelunch*, and *white_asian* is different for the treatment and control groups, but differences are not statistically significant. Mean differences for *tchexper* are large and statistically significant.

5. Formal test of balance

```
mod_star_check <- lm(small ~ boy + white_asian + tchexper + freelunch,
  data = star)

summary(mod_star_check)
```

Results of the OLS show that coefficients on all covariates are not statistically significant meaning that treatment and control groups do not differ on these covariates.

6. Effect of small classes on test scores

```
mod1 = lm(totalscore ~ small, data = star)
mod2 = lm(totalscore ~ small + white_asian + freelunch, data = star)
mod3 = lm(totalscore ~ small + white_asian + tchexper +
  freelunch + schurban, data = star)
stargazer(mod1, mod2, mod3,
  title="Effect of small classes on test scores",
  type = "text", out='regression.rtf')
```

Effect of small classes on test scores

Dependent variable:			
	(1)	(2)	(3)
small	13.899*** (2.447)	13.789*** (2.355)	13.866*** (2.351)
white_asian		12.041*** (2.807)	9.500*** (3.260)
tchexper			0.727*** (0.208)
freelunch		-34.511*** (2.616)	-33.803*** (2.631)
schurban			-3.004 (3.178)
Constant	918.043*** (1.667)	926.191*** (3.138)	921.897*** (4.195)
Observations	3,743	3,743	3,743
R2	0.009	0.082	0.085
Adjusted R2	0.008	0.081	0.084
Residual Std. Error	74.651 (df = 3741)	71.850 (df = 3739)	71.734 (df = 3737)
F Statistic	32.273*** (df = 1; 3741)	111.377*** (df = 3; 3739)	69.869*** (df = 5; 3737)
Note: *p<0.1; **p<0.05; ***p<0.01			

When we regress *totalscore* on *small*, the coefficient of *small* is 13.899 and it is statistically significant. This coefficient tells us the treatment effect. On average, students in small classes tend to have higher scores by 13.899 points than those who are studying in regular classes.

When we add more independent variables, the treatment effect does not differ much. It varies between 13.789 and 13.899.

Labor training

Import data

```
l = read.csv('lalonge.csv')
```

Effect of training program with and without covariates

```
mod4 = lm(re78 ~ treat, data = l)
mod5 = lm(re78 ~ treat + age + education + married, data = l)
mod6 = lm(re78 ~ treat + age + education + married +
          black + nodegree + re74 + re75, data = l)
stargazer(mod4, mod5, mod6,
          title="Effect of labor training program on earnings",
          type = "text", out='regression.rtf')
```

```
Effect of labor training program on earnings
=====
Dependent variable:
-----
(1)                re78                (3)
-----
treat              1,794.342***          1,655.255***          1,670.982***
                  (632.853)             (633.255)             (635.949)
age                43.350                55.125
                  (44.859)             (45.196)
education          391.018**             393.723*
                  (174.572)            (226.362)
married            151.591                -133.444
                  (852.107)            (877.316)
black              -2,245.301***
                  (841.727)
nodegree           -64.728
                  (1,001.674)
re74               0.082
                  (0.077)
re75               0.053
                  (0.135)
Constant           4,554.801***          -499.380             893.297
                  (408.046)            (2,106.928)          (3,212.735)
-----
Observations       445                   445                   445
R2                 0.018                   0.032                   0.055
Adjusted R2        0.016                   0.023                   0.037
Residual Std. Error 6,579.542 (df = 443)  6,555.299 (df = 440)  6,506.089 (df = 436)
F Statistic        8.039*** (df = 1; 443)  3.595*** (df = 4; 440)  3.160*** (df = 8; 436)
=====
Note:                *p<0.1; **p<0.05; ***p<0.01
```


When we run a regression with only the treatment variable, we get a statistically significant coefficient. The training program, on average, increased earnings by 1,794 units.

If we include more covariates, the effect of the training program decreases.

Age and education groups

```
summary(l$age)
```

```
summary(l$education)
```

	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
age	17.00	20.00	24.00	25.37	28.00	55.00
education	3.0	9.0	10.0	10.2	11.0	16.0

The minimum age in data is 17 years old. At the age of 21, a person legally becomes an adult. That's why the first group will be between 17 and 21. I think the majority complete their bachelor's or master's degrees by the age of 25. So, the second group will be from 21 to 26. The next group will be between 26 and 40, which is middle age. The maximum age in data is 55, so the last group will be between 41 and 55. Groups by age will be as follows:

- 17 – 20
- 21 – 25
- 26 – 40
- 41 – 55

I would like to group years of education according to the education system of the United States. The minimum years of education in data is 3 years and the maximum one is 16 years. So groups will be as follows:

- 3 – 6
- 7 – 8
- 9 – 12
- 13 – 16

There is no necessity to coarsen marriage status and race further because they are dummy variables and take only two values.

Matching

```
library(cem)

agecut = c(0, 20.5, 25.5, 40.5, 55.5)

educut = c(0, 6.5, 8.5, 12.5, 16.5)

mat = cem(treatment = "treat", data = 1, drop = "re78",
          cutpoints = list(education=educut, age=agecut))
```

Effect of training

```
est1 = att(mat, re78 ~ treat, data = 1)

est2 = att(mat, re78 ~ treat + age + education + married, data = 1)

est3 = att(mat, re78 ~ treat + age + education + married +
            black + nodegree + re74 + re75, data = 1)
```

Treatment effect			
	est1	est2	est3
SATT point estimate	2008.237673	2030.775763	1989.539342
p.value	0.006695	0.005831	0.006259

The effect of labor training on earnings is 2008.24 units without covariates. When we add variables *age*, *education*, and *married* the treatment effect increases to 2030.78 units. But if we continue to add more variables such as *black*, *nodegree*, *re74*, and *re75*, the treatment effect decreases to 1989.54 units.

If we compare the treatment effect with and without matching, the effect of the training program is higher when we do matching.