

Introduction

Education, playing a more and more important role in human society nowadays, enables students as well as society to make progress by obtaining professional skills, scientific innovation and spiritual strength. One of the most the important parts of Meiji restoration in Japan is education. Meiji restoration promoted significantly social progress, helped Japan develop rapidly and became the only Asian country that kept democracy & independence in the modern times. It also built the foundation of being a powerful developed country for Japan. Nearly all the countries treat education as the foundation of development and their future competitiveness.

Academic performance is one of the outcomes of education that can be observed easily. It is affected by many factors such as gender, school type, the conditions for learning, parental education level, family annual income, etc.

Unfortunately, while testing the effect from a certain intervention to students' academic performance, there is always bias from confounding variables which may influence the estimate of how the intervention we focus on will work on the outcome.

In randomized experiments, the biased estimation can be weakened. Because in terms of the law of large numbers, randomization will average the bias to make both test and control group balanced. However, the intervention assigned to school students can hardly be random. Thus, we will use propensity score matching to simulate randomized experiment by extracting individuals from experimental subjects, to create test and control group based on the comparable covariates. Propensity score can estimate the effect of receiving treatment when it is not suitable to assign treatments to subjects randomly. Propensity score matching refers to the pairing of test and control units with similar values on the propensity score; and possibly other covariates (the characteristics of participants); and the discarding of all unmatched units.

The school-level dataset was acquired from Kaggle. It is used to test whether a certain intervention (lunch type and test preparation course respectively in our case) can really improve students' academic performance in the total score of math, reading and writing. Conducting propensity score matching to eliminate bias from confounding variables (gender, parental education level, race/ethnicity), while investigating if standard lunch and test preparation course made progress in students' reading and mathematics performance.

Eliminating confounding variables can help us make more reasonable conclusion that if an intervention really brings higher academic achievement to students. With accurate conclusion, education department will be able to take more efficient actions to provide students with high quality studying experience.

Description of Dataset

This dataset was obtained from Kaggle for a predictive analysis project. It contains 1000 rows and 9 columns. The columns are gender, race_ethnicity, parental_level_of education, lunch, test_preparation_course, math_score, reading_score and writing_score.

	math_score	reading_score	writing_score
count	1000.00000	1000.000000	1000.000000
mean	66.08900	69.169000	68.054000
std	15.16308	14.600192	15.195657
min	0.00000	17.000000	10.000000
25%	57.00000	59.000000	57.750000
50%	66.00000	70.000000	69.000000
75%	77.00000	79.000000	79.000000
max	100.00000	100.000000	100.000000

Table 1 Description of Scores

Data Cleaning and Wrangling

The steps of data cleaning and wrangling are as follows:

1. Applied *dropna()* function to remove rows containing missing values
2. Dropped rows with value beyond the 3 standard deviations of writing_score column.
Repeated the same process for reading_score and math_score columns respectively.

```
# Remove the score outliers from df
df = df[np.abs(df['math_score']-df['math_score'].mean()) <= (3*df['math_score'].std())]
df = df[np.abs(df['reading_score']-df['reading_score'].mean()) <= (3*df['reading_score'].std())]
df_cleaned = df[np.abs(df['writing_score']-df['writing_score'].mean()) <= (3*df['writing_score'].std())]
df_cleaned
```

Figure 1 Outliers removal

- For those columns whose values are hierarchical, assigned increasing consecutive integers start from 0, to the values with level from low to high. Take parental education level for example, some high school, associate's degree, some college, bachelor, master correspond to the integers 0 to 5.
- For those columns whose values are non-hierarchical, applied function *pd.get_dummies* on the columns to quantify the values and deleted one column of the results to avoid dummy trap.

```
# quantify race_ethnicity
dummies1 = pd.get_dummies(df_cleaned.race_ethnicity)
dummies1
```

Figure 2 Non-hierarchical Column Quantifying

	group A	group B	group C	group D	group E
0	0	1	0	0	0
1	0	0	1	0	0
2	0	1	0	0	0
3	1	0	0	0	0
4	0	0	1	0	0
5	0	1	0	0	0
6	0	1	0	0	0

Table 2 Quantified race_ethnicity

- Merged the newly created columns in step 4 to the dataframe and then deleted the original non-hierarchical columns.

Exploratory Data Analysis

1. Intervention Effect Analysis

In the first part of the project, I use standard lunch and reduced lunch to classify the dataset into 2 groups, control and test respectively. Besides, I set total score of the 3 subjects as the dependent variable to study the relation between lunch types and students' academic performance. In this case, the other variables are considered covariates, and propensity score matching is applied to minimize the influence from the covariates such as gender, parental education level. The purpose of this project is to eliminate the selecting bias in contrast test.

```
# create test and control groups by lunch type
test = df_cleaned[df_quant.lunch == "free/reduced"]
control = df_cleaned[df_quant.lunch == "standard"]
test['lunch'] = 1
control['lunch'] = 0
```

Figure 1 Test and Control Groups Classification

Before doing inferential statistics to figure out what the outcome tells us, data wrangling and propensity score matching of the test and control groups are needed. For data wrangling, I follow the steps: (1) Apply *df.dropna()* to remove the missing values in the dataframe *df*. (2) Remove the rows in *df* which contain value beyond the range of 3 standard deviation of either one of the 3 score columns. (3) For the columns whose values can be ranked by number directly, assign consecutive integer to each type of value based on the level from low to high. (4) For those columns that cannot be quantify directly, use one hot encoding and delete one column to avoid dummy variable trap afterwards. (5) Import *Matcher* and classify test and control group by lunch type free/reduced and standard respectively. (6) Set the 3 types of score as dependent variables, lunch types as independent variables and others as covariates. (7) Conduct PSM using *Pymatch* module and we can obtain the pairs of matched test and control groups for inferential statistics.

```
# Matcher
m = Matcher(test, control, yvar="lunch", exclude=['math_score', 'reading_score', 'writing_score'])
```

Formula:
 lunch ~ gender+race_ethnicity+parental_level_of_education+test_preparation_course
 n majority: 643
 n minority: 350

```
# for reproducibility
np.random.seed(20170925)
m.fit_scores(balance=True, nmodels=350)
```

Fitting Models on Balanced Samples: 350\350
 Average Accuracy: 53.52%

Figure 2 Propensity Score Matching

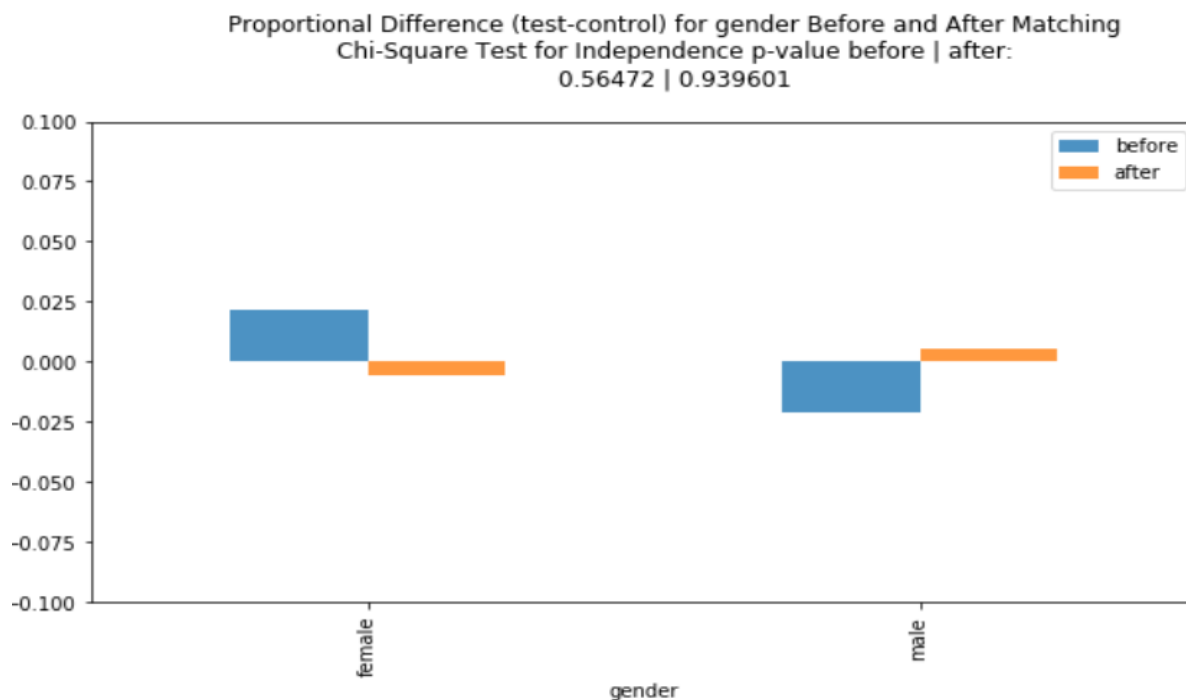


Figure 3 Proportion Difference for Gender

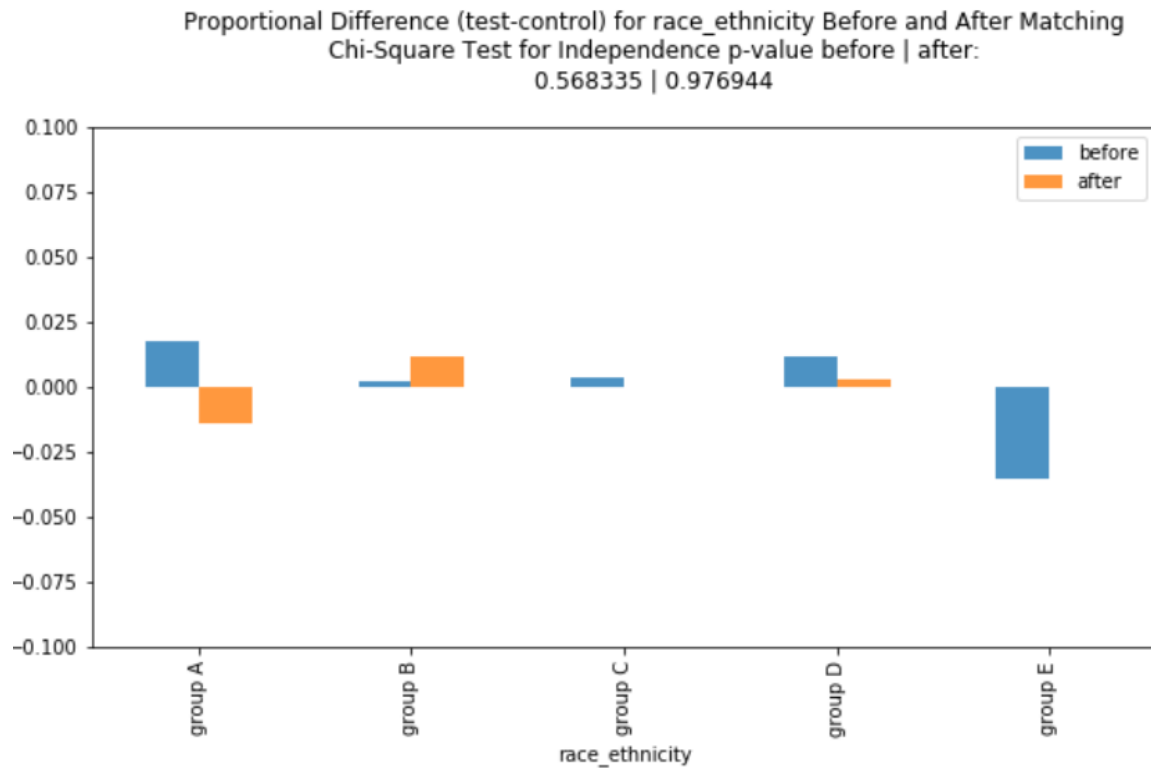


Figure 4 Proportional Difference for Race

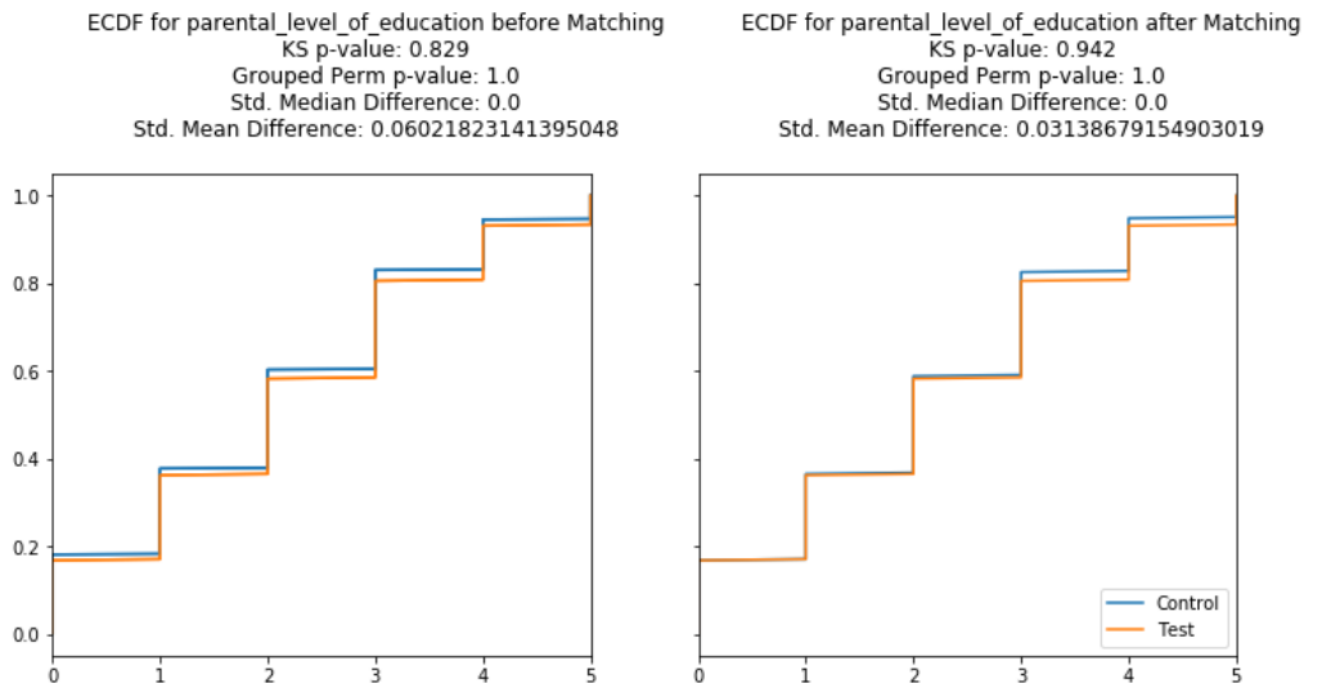


Figure 5 Assessment of Propensity Score Matching

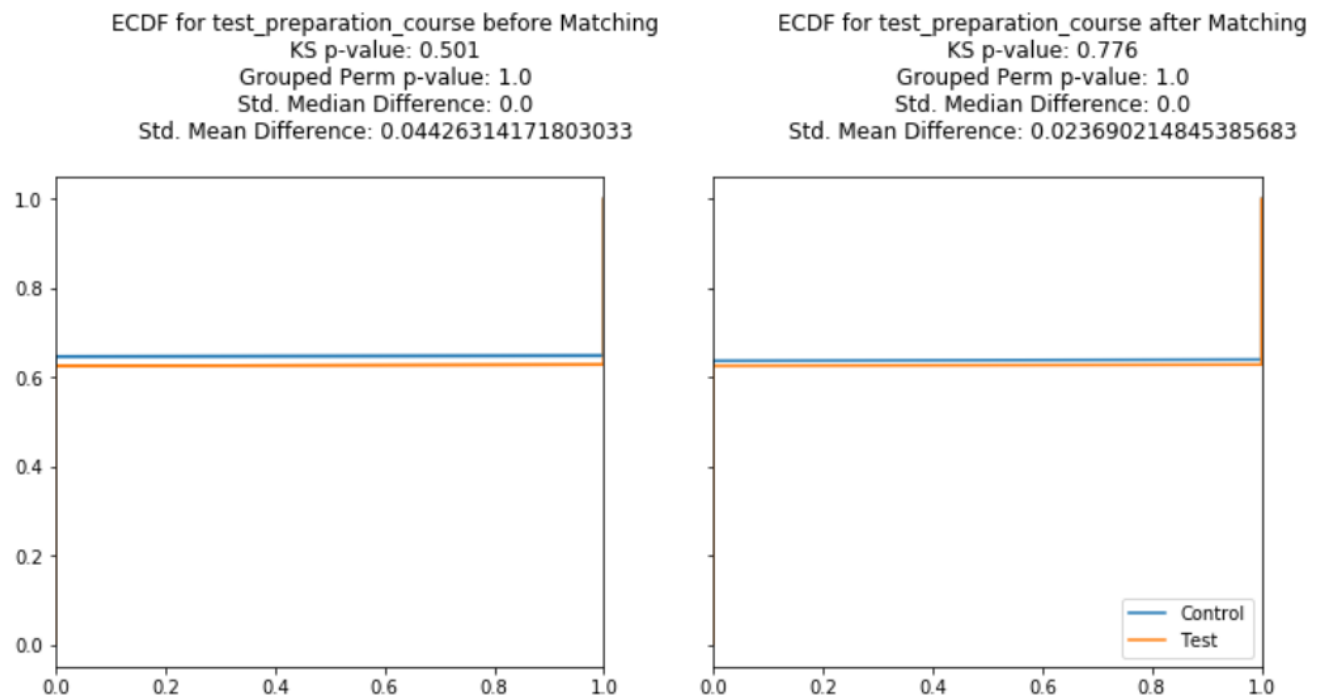


Figure 6 Assessment of Propensity Score Matching

In inferential statistic part, there are two methods to estimate the outcome.

1. Obtain the mean total score, standard deviation of the group with standard lunch and the group with reduced lunch respectively. μ_1 , μ_2 , σ_1 and σ_2 . Besides, obtain $\mu_1 - \mu_2$ and standard deviation of mean difference σ_{1-2} .
2. Assume significance level at $\alpha = 5\%$, null hypothesis $\mu_1 - \mu_2 = 0$ and alternative hypothesis: $\mu_1 - \mu_2 > 0$. Since $(\mu_1 - \mu_2) = 18.15$, which is out of 95% CI, I reject null hypothesis in favor of alternative hypothesis that $\mu_1 > \mu_2$.
3. Import *ttest_ind* from *scipy.stats* and conduct t test based on μ_1 , μ_2 , σ_1 and σ_2 . The small result also shows that I should reject null hypothesis in favor of alternative hypothesis.

```

from scipy.stats import ttest_ind
# Method 2 t test
lunch = np.random.normal(miu1,sigma1,1000)
no_lunch = np.random.normal(miu2,sigma2,1000)
t_test = ttest_ind(lunch,no_lunch)
t_test
# p_value < alpha. We can reject Ho in favor of Ha

```

Figure 7 Method 2 of Inferential Statistics

After hypothesis test, calculate correlation by importing *linregress* from *scipy.stats* and we can see the correlation r-value between lunch type and total score is -0.216, which means when lunch value is 1 (test group: free/reduced lunch), students get lower score, and when lunch type is 0, students get higher score. we want to compare the difference between means of total score of these two groups. The means and standards deviation can both be obtained from our *total_score* columns. Thus, hypothesis test on difference of means can be an effective way and t test is chosen in our case. The steps are as follows: (1) From *scipy.stats* import *ttest_ind*. (2) Calculate means and standard deviations of test and control groups. (3) Set lunch by assigning *miu1*, *sigma1* and size 1000 to *np.random.normal* respectively. Set *no_lunch* by filling in *np.random.normal* with *miu2*, *sigma2*, 1000. (4) Conduct t test and obtain p-value. (5) Compare p-value with significant level 5%.

In the second part of the project, I classify test and control groups by *test_preparation_course* (complete:1, none:0) and repeat the whole process mentioned above. In this case, we can see the p-value is significantly smaller than significance level 5%, which implies we should reject null hypothesis in favor of alternative hypothesis that the test group (test preparation course completed) has higher mean total score. The r-value is 0.30, meaning that taking the course can affect the total score positively.

In a word, the results of the project show that standard lunch and test preparation course can help students perform better in academic field.

2. Population Targeting Analysis

From the analysis above, we already know that both standard lunch type and test preparation course can provide students' academic performance with positive effects. However, should we target a certain group of people to make our intervention more efficient? This is definitely an important

and practical topic needed to be answered. In this project, I chose students with standard lunch and those with free/reduced lunch as two different groups. The objective is to find out that after receiving test preparation course, which group of students can have a more significant progress on the total score, and whether we should emphasize implementing the intervention to a certain group to have a cost-efficient outcome.

To achieve this analysis, the steps are as follows:

1. Split the data frame into 4 groups: free/reduced lunch without test preparation course, free/reduced lunch with course, standard lunch without course, standard lunch with course.
2. Extract a sample of 131 units (the least length among the 4 group in step 1) for the 4 groups.
3. Calculate the difference of total score between free/reduced lunch with course and free/reduced lunch without course, as well as the difference between standard lunch with course and standard lunch without course. Now, we have the progress both for students with free/reduced lunch and students with standard lunch after giving them test preparation course.
4. Calculate the difference between the progress of students with standard lunch and progress of students with free/reduced lunch. Average the difference. The value is 1.893, which means the former have a slightly better progress than the later.
5. Conduct Z test to find out if the difference between progresses is significant. We obtain the p-value larger than alpha value 0.05. Thus, we fail to reject null hypothesis, which means the difference between progresses is not significant.

```
# difference of progresses between the 2 groups
progress_diff = progress_std - progress_free
progress_diff

array([ 91,  72, -69, -108,   3, -53, -24, 116, -76, -31, -18,
       -57,  90, -84, -120, 125, -25, -15, 139, -85,  49, -14,
       -98, 110, -42, -55,  79,  43,  -4,  81, 112,  40, -89,
       -34, -121, -126, -12,  69,  10, -134,  36,  66,  99,  26,
      -108, 194,  94, -12, -160, -92,  79,  32,  48,  12, -104,
        41,  74,  -1,  27, -11, -73, -27, -73, -37,  32, -66,
       -62, -124, -222,  79, -70,  82,  13, 119,   9, -166,  11,
        97,  31, -52, -10,  42,  56, -61, -55,  71, -138, -89,
        78, -32, -84, -20, -119, 107,  75,   8,  85, -101,  28,
       -25, -65, -21, 123, -20, -26,   9, 158,  -3, -77,  40,
        14, -142,   3,  88, -64,  32,  60, 135, -69,   4, -20,
        61, -27, 122, -82, -73, -47,  71,  -4,  -3, -39],
      dtype=int64)
```

Figure 8 Differences between Progresses

Conclusion

After conduct propensity score matching to eliminate bias when we select experimental subjects, and applying inferential statistics to compare the difference of students' academic performance in different conditions, we are able to draw a conclusion that standard lunch and test preparation course can help students perform better in the academic field. It would be a good choice for the government to invest more in food in the underdeveloped regions. It is not necessary to consider regional financial situation while investing in test preparation course, because the effect of this intervention does not vary whether students have quality food or not.

Future work

Although we've found out that standard lunch and test preparation course can affect academic performance, we don't know how much these variables can improve the performance. In the future work, we can focus on performance prediction base on machine learning skills such as linear regression and logistic regression. It will help government make fully use of limited fund to provide students with better education more efficiently.