

Empirical Project 1

The Oregon Health Plan Experiment (OHP)

Due at midnight (Pacific Time) on Thursday, April 29, 2019

In this empirical project, you will analyze experimental data from the Oregon Health Plan Experiment (OHP) which we discussed in class. In particular, you will analyze real data from in-person interviews conducted in the Portland, Oregon, metropolitan area. The interviews were conducted about 25 months after the OHP lottery. The interview included detailed questionnaires on insurance coverage, health care use, health status, inventory of medications. The interview also contains questions that help assess depression and health-related quality of life. Finally, performance of anthropometric and blood-pressure measurements is taken, and dried blood spots are obtained. For more information about OHP and the in-person interviews, see [here](#) and [here](#).

Instructions

Please submit your Empirical Project on Canvas. Your submission should be a single PDF file containing three parts:

1. A 4-6 page research summary (double spaced and including references, graphs, and tables)
2. A copy of do-file with your STATA code or an .R script file with your R code
3. A copy of the log file (or screen output) of your STATA or R output

Specific questions to address in your research summary

1. Explain the difference between the variables *treatment* and *ohp_all_ever_survey*. Explain why *treatment* is the treatment variable (D_i), rather than *ohp_all_ever_survey*.
2. Provide evidence that the OHP lottery really did randomly assign individuals to treatment and control groups. Similar to Table 1 in [Taubman et al. \(2014\)](#), please create a nicely formatted table that reports means of 4 to 6 relevant characteristics for individuals in the *control group*.

Note: Part of this question is to get you to think about which variables should be balanced in a randomized experiment. You need to read carefully through all the variables in the dataset (documentation attached at the end of this file) and decide which 4 to 6 you will summarize.

3. For each of the variables you summarized above, calculate:
 - (i) the difference between the mean in the treatment group and the mean in the control group;

(ii) the standard error for the difference in means.

Add these as columns two and three to the table you started in question 2.

4. Is the balance table consistent with individuals having been randomly assigned to treatment group and control groups? Why or why not?
5. Estimate the compliance rate for the OHP experiment. That is, what is the effect of being assigned to the treatment group on the probability of being enrolled in Medicaid?

Hint: For this question and question 7, you can use the same regression as in question 3, just changing the dependent variable.

6. What is the intent-to-treat (ITT) effect of the OHP experiment on health outcomes? Please create a nicely formatted table that reports ITT estimates on 4 to 6 relevant health outcomes. Again, part of this question is to get you to think about which 4 to 6 variables could be used as health outcome variables.
7. What is the “treatment on the treated” effect (ATET) of the OHP experiment, i.e. the effect among those who applied for Medicaid? Estimate it for every health outcome you chose in question 6 and provide some intuition for the calculation of this estimate.
8. Do you have to worry about attrition bias in analyzing this data? Explain why or why not.
9. Suppose that you are submitting these results to a general interest journal such as *Science* for publication. Write an abstract of 200 or fewer words describing what you have found in your analysis of the OHP data, similar to the abstract in [Taubman et al. \(2014\)](#).

DATA DESCRIPTION, FILE: *ohp.dta*

The data consist of $n = 12,229$ individuals involved in OHP. Each individual is assigned a scrambled ID to protect privacy.

Variable Definitions in *ohp.dta*

Variable	Definition
<i>person_id</i>	Scrambled individual identifier
<i>household_id</i>	Scrambled household identifier
<i>weight_total_inp</i>	Survey weights (inverse probability weighting)
<i>treatment</i>	1 if OHP lottery winner, 0 otherwise
<i>age_inp</i>	Age
<i>bp_sar_inp</i>	Systolic blood pressure, average of three consecutive readings
<i>chl_inp</i>	Total cholesterol (dried blood spot test)
<i>dep_dx_post_lottery</i>	Diagnosed with depression after the lottery
<i>dep_dx_pre_lottery</i>	Diagnosed with depression before the lottery
<i>dia_dx_post_lottery</i>	Diagnosed with diabetes after the lottery
<i>dia_dx_pre_lottery</i>	Diagnosed with diabetes before the lottery
<i>doc_num_mod_inp</i>	Num. of doctor's visits, truncated at 2*99th percentile
<i>edu_inp</i>	Education: highest completed (1 = less than high school; 2 = high school diploma; 3 = post high school, not 4-year college; 4 = 4-year college degree or more)
<i>gender_inp</i>	1 if female
<i>hbp_dx_post_lottery</i>	Diagnosed with hypertension after the lottery
<i>hbp_dx_pre_lottery</i>	Diagnosed with hypertension before the lottery
<i>hispanic_inp</i>	Hispanic/Latino
<i>itvw_english_inp</i>	Interviewed in English
<i>numhh_list</i>	Number of people in household on lottery list
<i>ohp_all_ever_survey</i>	1 if ever enrolled in Medicaid
<i>race_black_inp</i>	Race/Ethnicity is Black
<i>race_nwother_inp</i>	Race/Ethnicity is Non-White Other
<i>race_white_inp</i>	Race/Ethnicity is White
<i>rx_num_mod_inp</i>	Number of prescription medications currently taking

Example R Commands

R command	Description
<pre>*Subset data ohp_cntrl <- subset(ohp, treatment ==0, select = c(xvar1, xvar2, xvar3)) *Report summary statistics summary(ohp_cntrl)</pre>	Subsets the data to the variables <i>xvar1</i> , <i>xvar2</i> , and <i>xvar3</i> for observations with <i>treatment</i> equal to 0. Reports summary statistics for this data frame.
<pre>#Install and load sandwich and lmtest packages install.packages("sandwich") install.packages("lmtest") library(sandwich) library(lmtest) #Regression with homoskedasticity-only standard errors mod1 <- lm(yvar~zvar1+wvar, data = ohp) summary(mod1) #Report heteroskedasticity robust standard errors coeftest(mod1, vcov = vcovHC(mod1, type="HC1"))</pre>	Estimates multivariate regression of <i>yvar</i> on an intercept, <i>zvar</i> , and <i>wvar</i> , with heteroskedasticity-robust standard errors.
<pre>*Method 1 mod1 <-lm(yvar ~ zvar+wvar, data = ohp) mod2 <-lm(xvar ~ zvar+wvar , data = ohp) mod1\$coef[2]/mod2\$coef[2] *Method 2 install.packages("AER") library(AER) tot <- ivreg(yvar ~ xvar + wvar zvar + wvar ,data=ohp) summary(tot, vcov = sandwich, df = Inf, diagnostics = TRUE)</pre>	<p>These commands show how to estimate the ratio between the coefficient on <i>zvar</i> from two different regressions that have all the same right-hand side variables, and different dependent variables.</p> <p>The first method refers to the elements of the vector <code>coef</code> that contains the relevant coefficients and displays the ratio between them.</p> <p>The second method uses the <code>ivreg</code> command to estimate the ratio in one step.</p>
<pre>ohp\$difference <- ohp\$xvar1 - ohp\$xvar2</pre>	Generates a new variable that equals the difference between <i>xvar1</i> and <i>xvar2</i>