André Filipe Silva - 26005

## MICROECONOMETRICS – Final Exam Delivery

## 1.Panel Data

**1a)**

Research question: <u>How does heavy drinking impact wages?</u>

Literature often shows correlation between drinking problems and poverty. My work will be addressing how heavy alcohol consumption impacts average hourly wage, controlling for a number of covariates that seem relevant and that are included in our dataset.

I will be defining "heavy drinking" as having 6 or more drinks in one sitting. This is given by the variable *drnk6m* in the dataset. Everyone that drinks 6 or more drinks in one sitting per month will be considered a "heavy drinker", even though of course the higher the *drnk6m* variable value, the heavier a drinker they are.

The average hourly wage variable *loghrwage* will be created by using *ln (wgsal / hrswrk).* This will be my labor market outcome – my dependent variable.
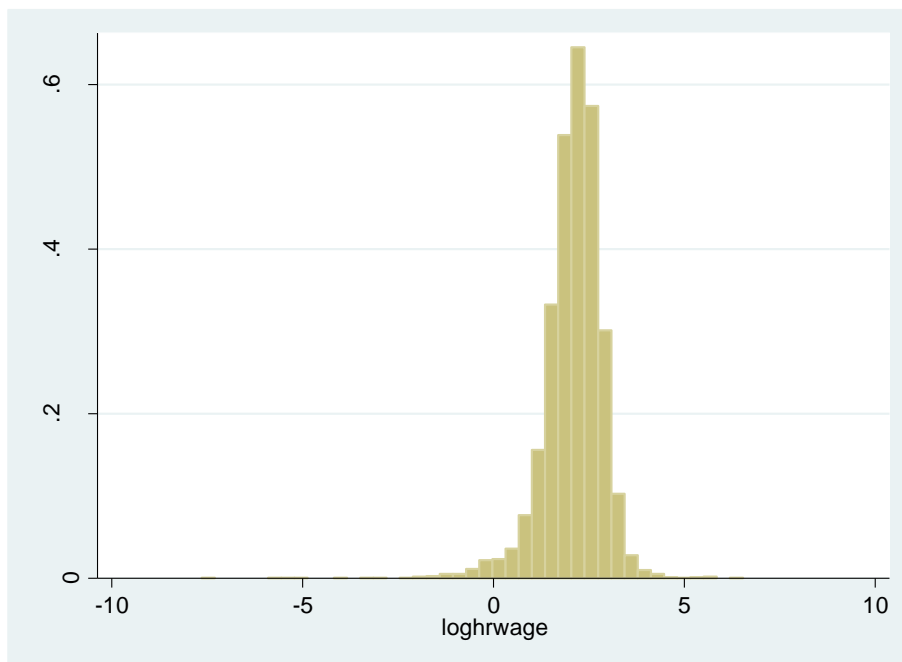
**1b)**

| Freq. | Percent | Cum. | Pattern* |
|-------|---------|--------|----------|
| 6105 | 61.70 | 61.70 | 11 |
| 3190 | 32.24 | 93.95 | 1. |
| 599 | 6.05 | 100.00 | .1 |
| 9894 | 100.00 | | XX |

The table above shows us that 6105 individuals are observed for the two periods, 3190 only for the first, and 599 only for the second. In total, we have 9894 unique individuals, but only 6105 that were observed in both periods. This means we have a very unbalanced panel, unfortunately.

 Defining retention rate as the number of individuals observed for all the time periods, in proportion of all individuals, we get a retention rate of (6105/9894)*100 = 61.7% - this implies an attrition rate of 38.3%.

<u>Description of variable *loghrwage*</u>

| Variable | Obs | Mean | Std. Dev. | Min | Max |
|----------|-----|------|-----------|-----|-----|
| loghrwage | 13,243 | 2.1019 | .7632306 | -7.640123 | 6.547606 |

Looking at the table and histogram, we can see a normally distributed variable with a lot of variation that can be explored – which is a good thing.

Description of *drnk6m* variable

| drnk6m | Overall Freq. | Overall Percent | Between Freq. | Between Percent | Within Percent |
|---|---|---|---|---|---|
| 0 | 10885 | 68.04 | 7692 | 77.74 | 89.55 |
| 1 | 1436 | 8.98 | 1335 | 13.49 | 62.96 |
| 2 | 1786 | 11.16 | 1608 | 16.25 | 65.36 |
| 3 | 827 | 5.17 | 772 | 7.80 | 61.33 |
| 4 | 347 | 2.17 | 336 | 3.40 | 61.31 |
| 5 | 155 | 0.97 | 149 | 1.51 | 60.74 |
| 6 | 563 | 3.52 | 501 | 5.06 | 68.76 |
| Total | 15999 | 100.00 | 12393 | 125.26 | 79.84 |

(n = 9894)

This variable is categorical. There is significant within variation for every category (which is good for the next questions), except for drnk6m==0. The observation count is considerably lower for higher values of this variable, which should be noted. However, there is no reason to suspect that the lower number is not just a result of it being an even of naturally rarer occurrence.

**1c)**

Controlling for individual-specific heterogeneity means we must use a within estimator – the Fixed Effects model is suitable. This model eliminates omitted variable bias due to time-invariant variables. This is good, however it also prevents us from getting estimates on time-invariant variables that we could be interested in. Even knowing that time-invariant variables will not be estimated, I will include them in the original equation of interest, so that it is easy to see which covariates I would be interested in estimating.

$$loghrwage_{it} = \beta_i drnk6m_{it} + \delta_1 health_{it} + \delta_2 logfaminc_{it} + \delta_3 povst_{it} + \delta_4 urate_{it}$$
$$+ \gamma_1 sex_i + \gamma_2 race_i + \gamma_3 afqtrev_i + \gamma_4 depression_i + \mu_{it}$$

With $\mu_{it} = c_i + \epsilon_{it}$

$\beta_i$ is the parameter of interest. In fact, $\beta_i$ is a vector of 6 parameters of interest due to the fact that *drnk6m* is a categorical variable. $\delta_i$ are time-variant control variables. $\gamma_i$ are time-invariant control variables.

Also note that $logfaminc_{it}$ was obtained by transforming the variable $faminc_{it}$.

**Estimation assumptions**:

1 – The attrition in our panel data, responsible for making it unbalanced, is uncorrelated with the idiosyncratic error. This is an essential assumption for estimation – from this it follows there are no sample selection problems and we will not have inconsistent or biased estimators.

2 – Strict Exogeneity: $E[\epsilon_{it}|x_{i1}, \ldots, x_{it}, c_i] = 0$ (note: take $x_{it}$ to account for all the time-variant variables in the model)

3 - $c_i$ is freely correlated with the explanatory variables: $E[x_{it}c_i] \neq 0$

From the list of covariates in our survey, I found these as most adequate. First, the gender-gap in wages is very real and obviously influences our dependent variable. The same goes for race. The remaining covariates' choice should appear obvious – I can't go into further detail as I have a character limit for the answer. I will just talk about the depression covariate as it is trickier to explain. If the dependent variable was something like weeks worked during a year, one would think a higher depression rate would lead to lower number of weeks worked in a year, and the inclusion in such a framework would be more obvious. But my argument for including it when the dependent variable is average hourly wage is that depression is many times a persistent phenomenon (due to poor treatment) and that might affect negatively the kind of jobs one individual gets and thus drive down hourly wages in a significant fashion.

The standard errors that will result from running this model through fixed effects will not be correct. The within transformation implies that there is serial correlation in $\widetilde{\epsilon_{it}}$. In fact, this has an easy intuition: for each individual, their observations over time are necessarily correlated with each other to some extent. Usually, we would cluster by individual to solve this and get the correct standard-errors – however, our panel has just two periods, making that approach impossible. The variance we get from our regression is given by: $\widehat{\sigma_{\widetilde{\epsilon}}^2} = \frac{SSR}{NT-K}$. But we need $\widehat{\sigma_{\epsilon}^2}$.

As a consistent estimate for $\sigma_{\epsilon}^2$ is $\frac{SSR}{N(T-1)-K}$, the difference is considerable when we have a small T, which is the case.

The equation to be estimated, through Fixed-Effects estimation, is obtained by transforming the original equation of interest as follows:

$loghrwage_{it} - \overline{loghrwage_{\iota}}$
$$= \beta_1 \left(drnk6m_{it} - \overline{drnk6m_{\iota}}\right) + \delta_1 \left(health_{it} - \overline{health_{\iota}}\right)$$
$$+ \delta_2 \left(logfaminc_{it} - \overline{logfaminc_{\iota}}\right) + \delta_3 \left(povst_{it} - \overline{povst_{\iota}}\right)$$
$$+ \delta_4 \left(urate_{it} - \overline{urate_{\iota}}\right) + \gamma_1 \left(sex_i - sex_i\right) + \gamma_2 \left(race_i - race_i\right)$$
$$+ \gamma_3 \left(afqtrev_i - afqtrev_i\right) + \gamma_4 \left(depression_i - depression_i\right) + c_i - c_i$$
$$+ \epsilon_{it} - \overline{\epsilon_{\iota}}$$

$\Leftrightarrow loghr\widetilde{wage}_{\iota t} = \beta_i \, drn\widetilde{k6m}_{\iota t} + \delta_1 hea\widetilde{lth}_{\iota t} + \delta_2 logfa\widetilde{minc}_{\iota t} + \delta_3 po\widetilde{vst}_{\iota t} + \delta_4 ura\widetilde{te}_{\iota t} + \widetilde{\epsilon_{\iota t}}$

**1d)**

The random effects estimation would be preferred with the same two initial assumptions explained above for the fixed effects estimator, but with a different third assumption:
3 – Uncorrelated unobserved effect: $E[x_{it}c_i] = 0$

Essentially meaning that $c_i$ is treated as a random variable as opposed to the fixed effects framework where it is treated as a parameter to be estimated for each cross section observation. The random effects estimation is done through GLS, and provided we get consistent estimators, it would allow us to get estimates on the influence of time-invariant variables that in this particular case would be helpful – as there are a lot of interesting time-invariant covariates to include in the model.

However, the orthogonality assumption between $c_i$ and $x_{it}$ is often a very strong one, and sticking to it would most likely lead to endogeneity problems.

**1e)**

```
Fixed-effects (within) regression          Number of obs     =      11,100
Group variable: id                         Number of groups  =       7,667

R-sq:                                       Obs per group:
     within  = 0.2327                                      min =           1
     between = 0.3292                                      avg =         1.4
     overall = 0.3286                                      max =           2

                                            F(10,3423)        =      103.81
corr(u_i, Xb)  = 0.1517                     Prob > F          =      0.0000
```

| loghrwage | Coef. | Std. Err. | t | P>\|t\| | [95% Conf. Interval] | |
|---|---|---|---|---|---|---|
| drnk6m | | | | | | |
| 1 | -.0047489 | .0297106 | -0.16 | 0.873 | -.0630013 | .0535035 |
| 2 | -.0470142 | .0296082 | -1.59 | 0.112 | -.1050657 | .0110372 |
| 3 | -.0795864 | .0398664 | -2.00 | 0.046 | -.1577508 | -.001422 |
| 4 | -.0469254 | .0591661 | -0.79 | 0.428 | -.1629298 | .069079 |
| 5 | -.0112692 | .081398 | -0.14 | 0.890 | -.1708627 | .1483244 |
| 6 | .1441169 | .05616 | 2.57 | 0.010 | .0340063 | .2542275 |
| | | | | | | |
| health | .0958279 | .0487168 | 1.97 | 0.049 | .0003111 | .1913448 |
| logfaminc | .3682156 | .0144592 | 25.47 | 0.000 | .3398661 | .396565 |
| povst | -.0743844 | .0423623 | -1.76 | 0.079 | -.1574424 | .0086736 |
| urate | .0286555 | .0037557 | 7.63 | 0.000 | .0212918 | .0360192 |
| _cons | .7125999 | .0533659 | 13.35 | 0.000 | .6079676 | .8172323 |

```
sigma_u  |  .60240638
sigma_e  |  .47442376
    rho  |  .6171958   (fraction of variance due to u_i)
```

```
F test that all u_i=0: F(7666, 3423) = 1.92          Prob > F = 0.0000
```

Focusing on the relation between heavy drinking and hourly wage, we get the signs we expected for the coefficients on the categorical variable *drnk6m*, except for when drnk6m==6. A general interpretation is: on average, drinking heavily has a negative impact on hourly wage, ceteris paribus. We would expect that the coefficients would become more negative as the number of times per month an individual drinks heavily increases, however that is only true until drnk6m==4, which is when the individual drank heavily 6 or 7 times in the past month. At a 5% significance level, the only statistically significant results are for drnk6m==3 and drnk6m==6. Again, this is not what one would expect from the start – especially the sign for drnk6m==6. One possible explanation might be that the attrition in our survey data is not random at all, leading to sample selection issues and inconsistent estimates. Also, we should keep in mind that the standard errors are not correct, as has been explained above.

**1f)**

$$loghr\widetilde{wage}_{it} = \beta_i(\widetilde{drnk6m}_{it} * sex_i) + \delta_1\widetilde{health}_{it} + \delta_2 log\widetilde{faminc}_{it} + \delta_3\widetilde{povst}_{it}$$
$$+ \delta_4\widetilde{urate}_{it} + \widetilde{\epsilon}_{it}$$

```
Fixed-effects (within) regression          Number of obs     =     11,100
Group variable: id                         Number of groups  =      7,667

R-sq:                                       Obs per group:
     within  = 0.2444                                    min =          1
     between = 0.2129                                    avg =        1.4
     overall = 0.2114                                    max =          2

                                            F(16,3417)        =      69.08
corr(u_i, Xb)   = -0.4102                   Prob > F          =     0.0000
```

| loghrwage | Coef. | Std. Err. | t | P>\|t\| | [95% Conf. | Interval] |
|---|---|---|---|---|---|---|
| drnk6m#sex | | | | | | |
| 0 2 | -1.05168 | .1459217 | -7.21 | 0.000 | -1.337782 | -.7655773 |
| 1 1 | -.0368812 | .0372941 | -0.99 | 0.323 | -.1100023 | .0362398 |
| 1 2 | -.9964904 | .149022 | -6.69 | 0.000 | -1.288672 | -.7043091 |
| 2 1 | -.0351456 | .0362783 | -0.97 | 0.333 | -.1062749 | .0359837 |
| 2 2 | -1.14589 | .1498968 | -7.64 | 0.000 | -1.439787 | -.8519937 |
| 3 1 | -.091506 | .0466562 | -1.96 | 0.050 | -.1829829 | -.0000291 |
| 3 2 | -1.07669 | .1511127 | -7.13 | 0.000 | -1.37297 | -.7804095 |
| 4 1 | -.0180536 | .0672477 | -0.27 | 0.788 | -.1499033 | .1137961 |
| 4 2 | -1.234331 | .1799174 | -6.86 | 0.000 | -1.587088 | -.8815747 |
| 5 1 | -.0610987 | .0902118 | -0.68 | 0.498 | -.2379732 | .1157758 |
| 5 2 | -.9821262 | .233894 | -4.20 | 0.000 | -1.440712 | -.5235399 |
| 6 1 | -.0018617 | .0611681 | -0.03 | 0.976 | -.1217914 | .1180681 |
| 6 2 | 0 | (omitted) | | | | |
| | | | | | | |
| health | .0899802 | .0485332 | 1.85 | 0.064 | -.0051769 | .1851374 |
| logfaminc | .3678086 | .0143732 | 25.59 | 0.000 | .3396277 | .3959895 |
| povst | -.0859014 | .0421693 | -2.04 | 0.042 | -.1685809 | -.0032219 |
| urate | .0289692 | .0037336 | 7.76 | 0.000 | .0216489 | .0362895 |
| _cons | 1.216435 | .0856663 | 14.20 | 0.000 | 1.048473 | 1.384397 |
| | | | | | | |
| sigma_u | .7048164 | | | | | |
| sigma_e | .47120171 | | | | | |
| rho | .69110753 | (fraction of variance due to u_i) | | | | |

```
F test that all u_i=0: F(7666, 3417) = 1.85                 Prob > F = 0.0000
```

The model presented is the same as the one for question 1c), but with a dummy variable for sex, where sex==1 stands for males and sex==2 stands for females, interacting with the *drnk6m* variable.

The first column on the *drnk6m* coefficients stands for the values of the categorical variable (1 through 6). The second column stand for the gender – 2 if female, 1 if male. So, the row with "1 1", for example, provides us the coefficient for a male individual that drank heavily once in the past month.

The coefficients on this model with the interaction variables show clear differences by gender on average hourly wages. The inclusion of gender effects analysis "fixes" the coefficient sign issues we had found in the previous answer. Here, the influence of heavy drinking on average hourly wage is always negative, both for males and females. There is always statistical significance for female heavy drinking. For men, however, we only find a statistically significant drop in hourly salary on the "31" row, that signifies being a male and drinking heavily 4 or 5 times in the last month.

To end, I would like to point out that although it is very likely that we get differences across gender, it seems to me unlikely that heavy drinking does not affect in a statistically significant manner the hourly wages of men, meaning that my estimates are most likely inconsistent – meaning the assumption of no sample selection may be too strong.
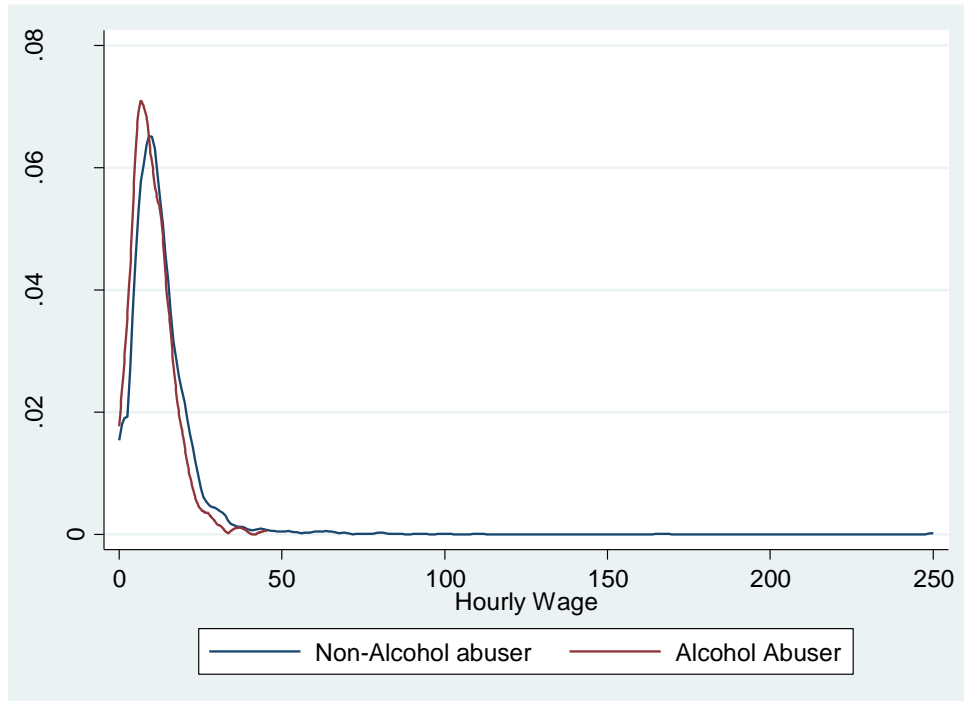
**1g)**

The model presented cannot have a causal interpretation. In order to have causal interpretation, we must have counterfactuals, because we must be comparing similar individuals. In this case, we observe the labor market outcomes for heavy drinkers, but do not observe the outcome for the same individuals in a context in which they would not be heavy drinkers. Or the other way around: we observe labor market outcomes for non-heavy drinkers, but do not observe labor market outcomes for non-heavy drinkers in a context when they would be heavy drinkers. The absence of the counterfactual can be overcome by various methods, however our simple fixed effects method does not assure causality. One of the main threats to the identification strategy is time-variant heterogeneity among individuals – the fixed-effects framework controls for time-invariant heterogeneity only. Since we need to get a counterfactual, we need to account for this kind of heterogeneity and here we are not comparing similar individuals because of this issue. There can also be the problem of reverse causality: low hourly wages may cause heavy drinking when some other characteristics are accounted for (which is common in "poverty traps").

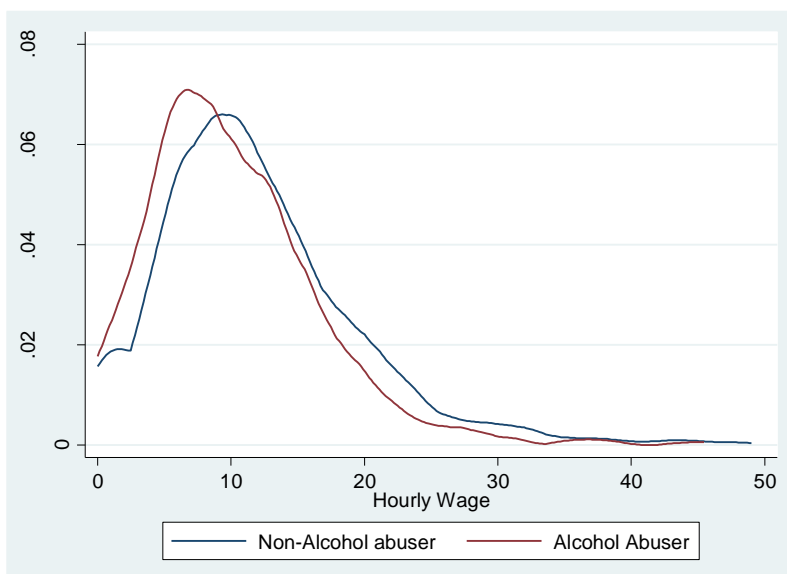## 2. Determining causal effects using propensity scores

**2a)**

I will show first the original graphic (including outliers).



Commenting first on the skewness of the graph. One thing is clear: very high hourly wages are not reachable for alcohol abusers - at least on the data we have, that is the way it seems. This makes sense intuitively, as alcohol abusers probably have a harder time being productive and keeping high-salary jobs. The "outliers" only exist for the non-alcohol abusing group.

Now, dropping the outliers.

The density curve peaks at a lower hourly wage value for the alcohol abusers. The concentration is more towards lower hourly wages than for non-alcohol abusers (density curve region leftwards compared to the non-alcohol abusers' curve). Again, this is to be expected. Alcohol abusers most likely have trouble keeping high-salary jobs, because of decreased ability to be productive, thus making them relatively less-paid than the non-alcohol abuser counterparts.

**2b)**

```
 ALCOHOLIC
ABUSER (=1;
      =0
otherwise)        Freq.      Percent        Cum.

        0        4,233        89.93        89.93
        1          474        10.07       100.00

    Total        4,707       100.00
```

Before providing sample means, it is important to note that we do not have many observations for the group of alcohol abusers. That is going to be important later on because it may take statistical power out of our estimations.

```
alcoholic |  perday      sex      race    afqtrev  depres~n    health     povst     urate    faminc

        0 | 2.51193  1.473187  2.38885  44.16702  3.475785  .0654382  .0930782  7.281006  39829.45
        1 | 6.016878 1.204641  2.28481  31.87342  4.417722  .0696203  .1455696  7.348101  23750.39
```

I only included one variable for the alcohol use because, obviously, being tagged as alcoholic means you have very different characteristics in terms of most (if not all) of the alcohol consumption variables. The *perday* variable suffices to illustrate my point. To ensure comparability across the two groups, we must match the participants over a set of common characteristics – and the alcohol use variables are essentially "self-selection".

**Difference in sample means tests:**

perday

**Two-sample t test with equal variances**

| | obs1 | obs2 | Mean1 | Mean2 | dif | St_Err | t_value | p_value |
|---|---|---|---|---|---|---|---|---|
| perday by alcoholic | 4233 | 474 | 2.512 | 6.017 | -3.505 | .111 | -31.8 | 0 |

sex

**Two-sample t test with equal variances**

| | obs1 | obs2 | Mean1 | Mean2 | dif | St_Err | t_value | p_value |
|---|---|---|---|---|---|---|---|---|
| sex by alcoholic: ~1 | 4233 | 474 | 1.473 | 1.204 | .269 | .024 | 11.3 | 0 |

### race

**Two-sample t test with equal variances**

|  | obs1 | obs2 | Mean1 | Mean2 | dif | St_Err | t_value | p_value |
|---|---|---|---|---|---|---|---|---|
| race by alcoholic:~1 | 4233 | 474 | 2.389 | 2.285 | .104 | .037 | 2.8 | .005 |

### afqtrev

**Two-sample t test with equal variances**

|  | obs1 | obs2 | Mean1 | Mean2 | dif | St_Err | t_value | p_value |
|---|---|---|---|---|---|---|---|---|
| afqtrev by alcohol~1 | 4233 | 474 | 44.167 | 31.874 | 12.294 | 1.398 | 8.8 | 0 |

### depression

**Two-sample t test with equal variances**

|  | obs1 | obs2 | Mean1 | Mean2 | dif | St_Err | t_value | p_value |
|---|---|---|---|---|---|---|---|---|
| depression by alco~1 | 4233 | 474 | 3.476 | 4.418 | -.942 | .19 | -4.95 | 0 |

### health

**Two-sample t test with equal variances**

|  | obs1 | obs2 | Mean1 | Mean2 | dif | St_Err | t_value | p_value |
|---|---|---|---|---|---|---|---|---|
| health by alcoholi~1 | 4233 | 474 | .066 | .07 | -.004 | .012 | -.35 | .728 |

### povst

**Two-sample t test with equal variances**

|  | obs1 | obs2 | Mean1 | Mean2 | dif | St_Err | t_value | p_value |
|---|---|---|---|---|---|---|---|---|
| povst by alcoholic~1 | 4233 | 474 | .093 | .145 | -.052 | .015 | -3.65 | .001 |

### urate

**Two-sample t test with equal variances**

|  | obs1 | obs2 | Mean1 | Mean2 | dif | St_Err | t_value | p_value |
|---|---|---|---|---|---|---|---|---|
| urate by alcoholic~1 | 4233 | 474 | 7.281 | 7.348 | -.067 | .133 | -.5 | .614 |

### faminc

**Two-sample t test with equal variances**

|  | obs1 | obs2 | Mean1 | Mean2 | dif | St_Err | t_value | p_value |
|---|---|---|---|---|---|---|---|---|
| faminc by alcoholi~1 | 4233 | 474 | 39829.45 | 23750.39 | 16079.06 | 1930.543 | 8.35 | 0 |

As we can see from the performed t-tests, the only variables for which there are no significant differences among groups are health problems and unemployment.

These significant differences in most variables are precisely what shows us that we need to resort to matching methods in order to get causal effects – we must find, or more specifically match, "similar" individuals in our sample so that we can thus obtain the treatment effect. Simply comparing outcomes on both groups without matching techniques would not lead to us to a causal treatment effect, because we would essentially be comparing "apples with oranges".
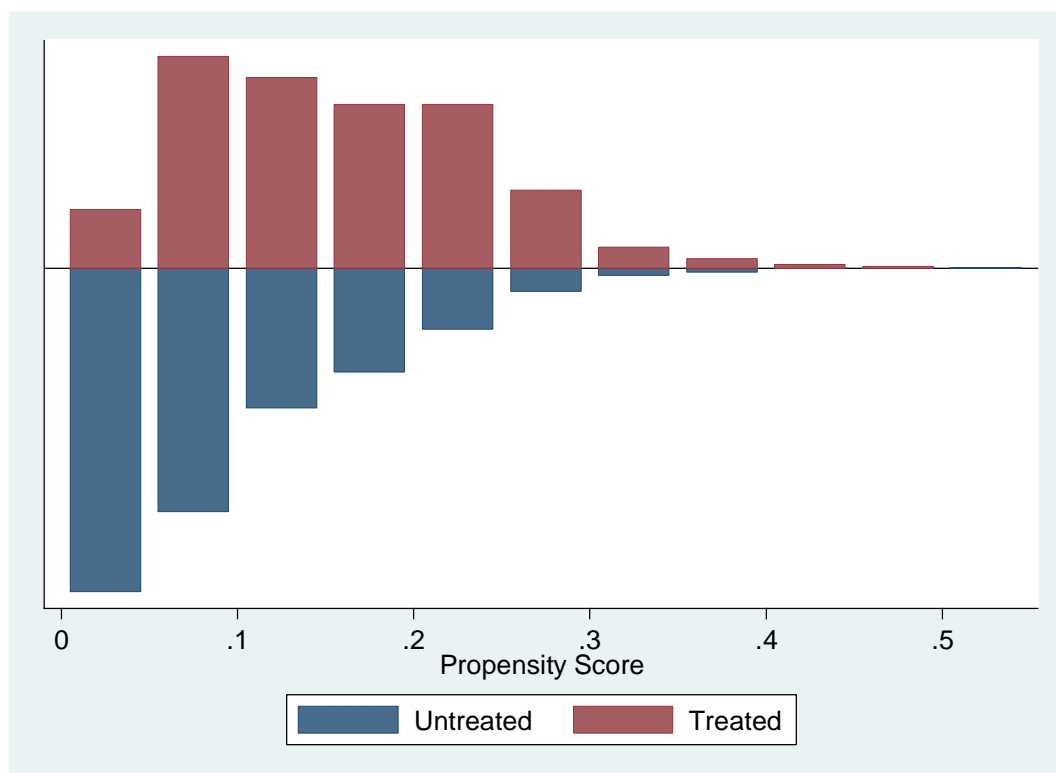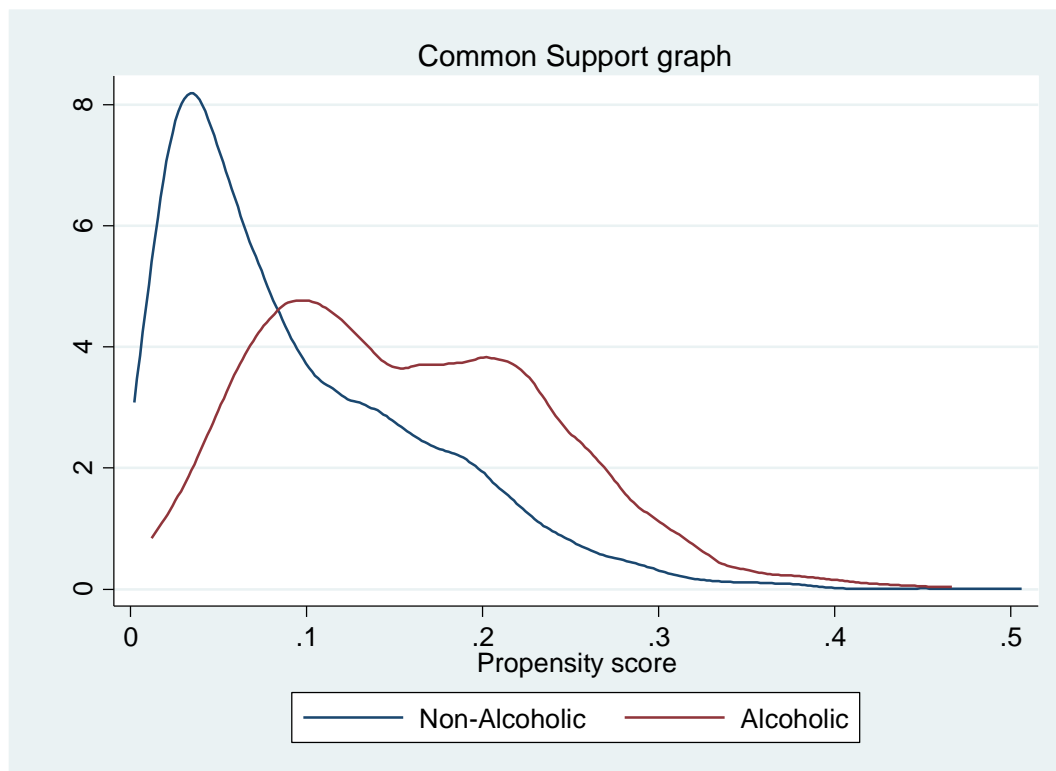
**2c)**

I will be using for the estimation of the propensity score all of the covariates that I provided sample means above, except for any alcohol use variables (such as the *perday*). My justification for not using alcohol related variables can be found in the previous question.

The first table reports the estimate by logit model of the predicted probability of being alcoholic based on, describing from first to last, the covariates sex, race, depression, health, intelligence percentile test, a dummy for the individual's family being below the poverty line, unemployment rate, and family income.

```
Logistic regression                              Number of obs    =      4,705
                                                 LR chi2(9)       =     297.98
                                                 Prob > chi2      =     0.0000
Log likelihood = -1388.2003                      Pseudo R2        =     0.0969
```

| alcoholic | Coef. | Std. Err. | z | P>|z| | [95% Conf. Interval] | |
|---|---|---|---|---|---|---|
| sex | -1.431874 | .1243996 | -11.51 | 0.000 | -1.675692 | -1.188055 |
| race | .0756575 | .0727772 | 1.04 | 0.299 | -.0669831 | .2182981 |
| afqtrev | -.0101642 | .0021331 | -4.76 | 0.000 | -.014345 | -.0059833 |
| depression | .0691237 | .0123938 | 5.58 | 0.000 | .0448323 | .0934151 |
| DMISS_depression | 0 | (omitted) | | | | |
| health | -.5147786 | .2073307 | -2.48 | 0.013 | -.9211393 | -.1084178 |
| DMISS_health | 0 | (omitted) | | | | |
| povst | .0260859 | .1746457 | 0.15 | 0.881 | -.3162134 | .3683853 |
| DMISS_povst | -.2463439 | .1614669 | -1.53 | 0.127 | -.5628132 | .0701254 |
| urate | .0010986 | .0181798 | 0.06 | 0.952 | -.0345331 | .0367304 |
| DMISS_urate | 0 | (omitted) | | | | |
| faminc | -.0000144 | 2.89e-06 | -5.00 | 0.000 | -.0000201 | -8.78e-06 |
| DMISS_faminc | 0 | (omitted) | | | | |
| _cons | .1591593 | .2913187 | 0.55 | 0.585 | -.4118149 | .7301334 |

Common Support graph



Note that the density curves graph stops at .5, as does the histogram, so it might be misleading if one neglects to check the range of the x-axis. I chose to plot it in this way because it gives us more graphical detail on the part that we will actually work with next, when applying the propensity score matching method.

From the graphs we can see that we have a good common support region – there is a significant overlap in P(X) across participants and non-participants. Observations with a propensity score above 0.4 will have to be dropped, as they are outside the common support region.

Although it may appear as a little odd that we have to drop observations for a wide range of P(X), given the variables at hand I believe it is not surprising. We have people tagged as alcoholic and people tagged as non-alcoholic. Intuitively, it does not appear reasonable that non-alcoholics would score very high on the propensity score of being alcoholic. As such, dropping a significant region of the observations would be inevitable in any set of data that deals with this type of issue.

**2d)**

```
. psmatch2 alcoholic, radius caliper(0.0001) outcome(hrwage) pscore(pscore1)
```

| Variable | Sample | Treated | Controls | Difference | S.E. | T-stat |
|---|---|---|---|---|---|---|
| hrwage | Unmatched | 10.0289726 | 12.6249464 | -2.59597376 | .556427871 | -4.67 |
| | ATT | 10.5621821 | 12.0813313 | -1.5191492 | .632158637 | -2.40 |

Note: S.E. does not take into account that the propensity score is estimated.

| psmatch2: Treatment assignment | psmatch2: Common support | | Total |
|---|---|---|---|
| | Off suppo | On suppor | |
| Untreated | 0 | 3,781 | 3,781 |
| Treated | 84 | 317 | 401 |
| Total | 84 | 4,098 | 4,182 |

I chose the Radius matching method, as we have significant differences between treated and non-treated sample means, and there not a lot of observations for the treated. As such, I figured that if matching on just the nearest neighbor, or n-nearest neighbors, we would not get good results as the nearest neighbor seems likely very far away. It is true, however, that the rule for the radius is a bit arbitrary. I experimented with other radiuses in order to try to get a good balance between the observations on common support and the robustness of the matching.

I should note that I also ran bootstrapping to obtain more valid standard errors, but the differences found were not significant compared to the table above. As such, I have not included that table in my final work, but the results can be found on my log file.

The results show a significant difference in the Average treatment on the treated effect (ATT), and it goes the way we predicted: the "treated" are alcoholic, and their average hourly wage is lower than the "non-treated".

In order to validate the results obtained, I will check the balancing of the covariates used through the *pstest* command. The resulting table is reproduced below.

| Variable | Mean Treated | Control | %bias | t-test t | p>\|t\| | V(T)/ V(C) |
|---|---|---|---|---|---|---|
| sex | 1.183 | 1.2006 | -4.1 | -0.56 | 0.574 | 0.93 |
| race | 2.3912 | 2.3325 | 7.6 | 0.97 | 0.335 | 1.05 |
| afqtrev | 38.533 | 38.325 | 0.8 | 0.10 | 0.923 | 0.91 |
| depression | 3.2303 | 3.3976 | -4.3 | -0.57 | 0.568 | 0.84 |
| health | .0347 | .04079 | -3.1 | -0.40 | 0.688 | . |
| povst | .0694 | .07593 | -2.5 | -0.32 | 0.752 | . |
| urate | 7.2729 | 7.2407 | 1.2 | 0.16 | 0.874 | 1.03 |
| faminc | 29229 | 27703 | 4.4 | 0.74 | 0.458 | 1.23 |

* if variance ratio outside [0.80; 1.25]

| Ps R2 | LR chi2 | p>chi2 | MeanBias | MedBias | B | R | %Var |
|---|---|---|---|---|---|---|---|
| 0.003 | 2.22 | 0.974 | 3.5 | 3.6 | 11.8 | 1.05 | 0 |

The balancing is good for all covariates, as can be seen from the t-tests' p-values. This is great news. It certainly helps in validating the estimates found above, and lending credibility to the ATT results found. The overall matching performance is satisfactory.

As for causal interpretation, I would not conclude in favor of it. The problem with matching methods is that different methods can lead to vastly different results. For us to be able to conclude for causal effects would require additional robustness checks, such as comparing results across different matching methods.

**2e)**

| stats | ATE_psw | ATT_psw |
|---|---|---|
| mean | -2.264761 | -.9943998 |

The ATT is very similar to the one found through the Propensity Score Matching estimation. The ATE is also similar. The fact that the values found here are close to the ones found through the deployment of the Propensity Score Matching technique lend credibility to the method used before. The differences between the ATE and the ATT show that compliance is far from 100%. However, it is a good sign that the values found are close to the ones found before.