# Methods in Epidemiologic Research

## Sample Problems

## Chapter 13 – Confounding

## Preparation

We will carry on with the data that were used in the exercises for Chapter 6, but we will focus on the effects of being obese (bmi>30) or having coronary angioplasty on the risk of dying within 1 year of admission for the MI included in this dataset. All individuals in this dataset were followed for at least 1 year after admission, so there are no withdrawals (censored observations) to worry about.

The first thing you will have to do is generate four new variables.

- Create a 0/1 variable which indicates whether or not the individual died during the 1st year following admission (-died_1yr-).

*If you are using Stata, some code to generate the variable is*
```
codebook surv_mi if died==0
gen died_1yr=(died==1 & surv_mi<365)
```
- Create a variable indicating if the person was obese (bmi>30) or not (-obese-).

*If you are using Stata, some code to generate the variable is*
```
gen obese=bmi>=30 & !missing(bmi)
```
- Create 2 categorical variables for age:

  - a dichotomous (0/1) variable with age split at 75 years (-age_c2-),

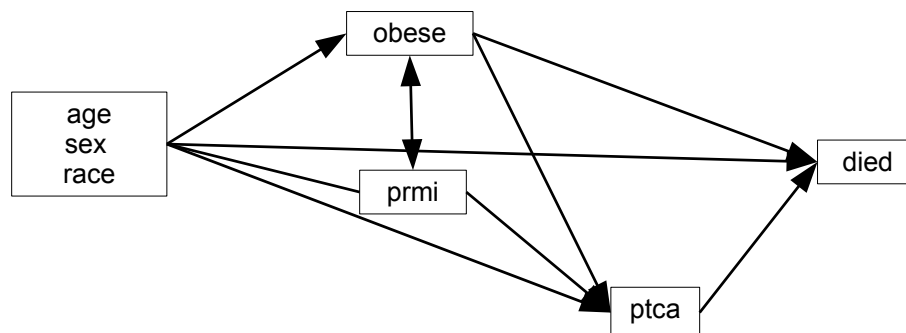  - a 5 category variable with splits at 60, 70, 80 and 90 years (-age_c5-).

*If you are using Stata, some code to generate the variable is*
```
egen age_c2=cut(age), at(0 75 999) icodes
egen age_c5=cut(age), at(0 60 70 80 90 999) icodes
```

## Questions

1. Draw a causal diagram which you think represents the relationships among the following variables: death within 1 year of admission (-diead_1yr-), sex (-sex-), race (-white-), previous MI (-prmi-), cardiac angioplasty (-ptca-), obese (-obese-), and age (-age_c2- or -age_c5-).

*The following would be a reasonable representation of the causal relationship among the variables of interest.*

*Age, sex and race (-white-) have been lumped together as characteristics of the patient which may affect subsequent variables of interest. No time (causal) relationship has been specified between -obese- and -prmi- because they may influence each other (ie being obese may alter the probability that you have had a previous MI, but a previous MI may also alter the possibility that you are now obese). Both -obese- and -prmi- may affect the probability of you having an angioplasty and all 3 likely influence your probability of survival.*

*Based on this diagram, -prmi- is a potential confounder of the relationship between -ptca- and -died-, but -ptca- can not be a confounder of the -obese- to -died- relationship (because it is intervening).*

(a) Evaluate the correlations among all of these variables. (Note – correlations are not designed for use with dichotomous variables but they are a "quick and dirty" way of assessing the relationships among these variables. They certainly should not be used with categorical variables so you should leave -age_c5- out of this correlation matrix)

```
. corr sex-age_c2
(obs=2894)
             |      sex     white      prmi      ptca died_1yr     obese    age_c2
-------------+---------------------------------------------------------------------
         sex |   1.0000
       white |   0.0049    1.0000
        prmi |  -0.0077   -0.0210    1.0000
        ptca |   0.1942   -0.0115   -0.1041    1.0000
    died_1yr |  -0.0738    0.0080    0.1165   -0.3134    1.0000
       obese |   0.0150    0.0155   -0.0593    0.0838   -0.1323    1.0000
      age_c2 |  -0.2496    0.1356    0.0963   -0.3435    0.2903   -0.2203    1.0000
```

*Age (-age_c2-) is quite highly correlated with most other variables (including -died_1yr-) suggesting it will be an important confounder for a number of variables (including -obese-). Race (-white-) has very low correlation with -died_1yr- so it is unlikely to exert much confounding effect. Etc., etc. ....*

2. Compute and interpret the risk ratio (RR) for being obese.

```
.            cs died_1yr obese
                   | obese                |
                   |  Exposed   Unexposed |      Total
-----------------+----------------------+------------
           Cases |      137          640 |        777
        Noncases |      679         1509 |       2188
-----------------+----------------------+------------
           Total |      816         2149 |       2965
                 |                       |
            Risk |  .1678922     .2978129 |    .2620573
                 |                       |
                 |       Point estimate  |    [95% Conf. Interval]
                 |-----------------------+----------------------
 Risk difference |        -.1299208      |    -.1620377    -.0978039
      Risk ratio |         .5637504      |     .4775357     .6655303
   Prev. frac. ex. |       .4362496      |     .3344697     .5224643
  Prev. frac. pop |         .1200606     |
                 +---------------------------------------------
                        chi2(1) =     51.62   Pr>chi2 = 0.0000
```
*The RR=0.56 suggests that being obese substantially reduced the risk of death in the first year*

*after admission (it almost cut the risk in half). Perhaps this seems counter-intuitive, but keep 2 things in mind. We are looking at the risk of death after an MI, not the risk of having an MI (for which obesity is a known risk factor). Second, the causal diagram suggested that a number of factors MIGHT be confounders and the correlation matrix highlighted age as potential confounder …. so we will evaluate ways for controlling the effects of age.*

3. Use matching to "control" for age when evaluating the effect of obese. Match on the 5-level version of age (-age_c5-). (Note – this takes quite a bit of programming skill, but (Stata) code is provided with the solution set).

*If you are using Stata, some code to perform the matching (and subsequent analysis) is shown below. The code determines the proportion of non-obese which need to be sampled in each age category to match the number of obese in that age category. It then randomly samples the non-obese in each age category.*

```
set seed 94827
tab age_c5 obese
capture drop ao*
egen aogrp=group(age_c5 obese)
bysort aogrp: gen aoN=_N
sort age_c5 id
by age_c5: egen aoN0=min(aoN)
by age_c5: egen aoN1=max(aoN)
gen ao_prp=aoN0/aoN1
     * br age_c5 obese aogrp aoN aoN0 aoN1 ao_prp
gen aorand=uniform()
keep if (obese==1) | (obese==0 & aorand<=ao_prp)
* verify that age distn. approx equal in obese and non-obese groups
     tab age_c5 obese
* evaluate effect of -obese-
     cs died_1yr obese
```

(a) Does controlling for age increase or reduce the apparent effect of being obese?

*The -tab- command above showed us that we now had roughly the same age distribution among obese and non-obese individuals (results not shown). However, this has reduced the number of observations used in the analysis from 2965 to 1677.*

```
. cs died_1yr obese
                 | obese                    |
                 |   Exposed    Unexposed   |      Total
-----------------+--------------------------+------------
          Cases  |       137          200   |        337
       Noncases  |       679          661   |       1340
-----------------+--------------------------+------------
          Total  |       816          861   |       1677
           Risk  |  .1678922      .232288   |   .2009541
                 |                          |
                 |         Point estimate   |    [95% Conf. Interval]
                 |--------------------------+------------------------
 Risk difference |          -.0643959       |   -.1025184   -.0262734
      Risk ratio |           .7227757       |    .5946444    .8785163
  Prev. frac. ex.|           .2772243       |    .1214837    .4053556
  Prev. frac. pop|           .1348927       |
                 +--------------------------------------------------
                            chi2(1) =    10.82  Pr>chi2 = 0.0010
```

*The RR has now increased from 0.56 (a 44% reduction in risk attributable to obesity) to 0.72 (a 28% reduction in risk). This suggests that part of the preventive effect of obesity was simply due to the fact that obesity was much more common in young MI patients than in older ones and the improved survival was attributable to their younger age.*

4. Now go back to the full dataset and analytically control (ie Mantel-Haenszel procedure) for the effect of -sex- on the effect -obese-.

```
.           cs died_1yr obese, by(sex)
            gender |       RR        [95% Conf. Interval]   M-H Weight
-----------------+-------------------------------------------------
              0 |     .6061196     .4850101    .7574709      88.01955
              1 |     .5254684     .4107681    .6721969      87.67829
-----------------+-------------------------------------------------
          Crude |     .5637504     .4775357    .6655303
   M-H combined |     .5658723     .4795729    .6677015
-----------------------------------------------------------------
Test of homogeneity (M-H)       chi2(1) =     0.713  Pr>chi2 = 0.3986
```

(a) Is -sex- an important confounder. If so, why? If not, why not?

*Sex appears to have no confounding effects because the crude and M-H adjusted RR are virtually identical. This is because -sex- has very weak associations with both -obese- and -died_1yr-. The following code carries out the analyses to demonstrate this (examining the results is left up to you).*

```
cs obese sex
cs died_1yr sex if obese==0
```

5. Control for the effects of age in the analysis of -obese-.

(a) First use the dichotomous version of age (-age_c2-). Is -age_c2- an important confounder. If so, why? If not, why not?

```
. cs died_1yr obese, by(age_c2)

            age_c2 |       RR        [95% Conf. Interval]   M-H Weight
-----------------+-------------------------------------------------
              0 |     .6980282     .525116     .9278775      52.98507
              1 |     .7297297     .5966241    .892531       84.6831
-----------------+-------------------------------------------------
          Crude |     .5637504     .4775357    .6655303
   M-H combined |     .7175286     .608283     .8463944
-----------------------------------------------------------------
Test of homogeneity (M-H)       chi2(1) =     0.063  Pr>chi2 = 0.8019
```

*As was noted in Question 3, the apparent effect of -age_c2- has been reduced from a RR of 0.56 to 0.72. This is because -age_c2- is strongly (but negatively) associated with -obese- (RR=0.47) and strongly positively associated with -died_1yr- (RR=2.7) among non-obese individuals. (Note, we carry out this last analysis using only non-obese individuals to ensure that we are measuring a true effect of age, not one that is influenced by the connection between age and obesity.)*

(b) Next use the 5 category version of age (-age_c5-). Does controlling the effect of age with the 5-category version have a bigger or smaller effect than the 2-category version? Why?

```
. cs died_1yr obese, by(age_c5)
            age_c5 |       RR        [95% Conf. Interval]   M-H Weight
```

```
----------------+--------------------------------------------------------------
            0 |    .7278119       .413731      1.280325       13.67699
            1 |    .6066986      .3696826      .9956736       19.51245
            2 |    .7537901      .5522742      1.028836       38.09211
            3 |    .8011782      .6269348      1.023849       49.68293
            4 |    1.023142      .6664021      1.570853       9.769565
----------------+--------------------------------------------------------------
        Crude |    .5637504      .4775357      .6655303
  M-H combined |    .7672557      .6527059      .9019089
----------------+--------------------------------------------------------------
Test of homogeneity (M-H)        chi2(4) =      2.760   Pr>chi2 = 0.5988
```

*The M-H adjusted RR has now risen to 0.77 (even less of an effect of obesity than we have seen in any previous analyses. The reason for the difference in the estimates based on -age_c2- and -age_c5- probably relates to how effectively these two variables "account for" age. -age_c2- is a very crude measure of age (over or under 75 yrs). -age_c5- reflects a persons age a bit better (10 yr increments) and consequently probably does a better job of removing the confounding effect of age. We could probably do an even better job of removing the confounding effects of age if we treated it as a continuous variable (ie used the person's actual age), but this requires some techniques which you will learn in Chapters 15 and 16.*

6.  We are now going to switch to evaluating the effects of angioplasty (-ptca-) on the risk of death.

    (a)  First compute the crude RR for the effect of -ptca- on -died_1yr-. Interpret the result.

```
. cs died_1yr ptca
                 | coronary angioplasty  |
                 |  Exposed    Unexposed |       Total
-----------------+-----------------------+------------
          Cases |      157          620  |         777
       Noncases |     1223          965  |        2188
-----------------+-----------------------+------------
          Total |     1380         1585  |        2965
                 |                        |
           Risk |  .1137681     .3911672 |    .2620573
                 |                        |
                 |      Point estimate    |    [95% Conf. Interval]
                 |------------------------+------------------------
Risk difference |        -.2773991       |   -.3066884    -.2481098
     Risk ratio |         .2908427       |    .2479511     .3411538
  Prev. frac. ex. |        .7091573       |    .6588462     .7520489
 Prev. frac. pop |         .3300631       |
                 +----------------------------------------------------
                     chi2(1) =    293.55  Pr>chi2 = 0.0000
```

*The risk ratio of 0.29 suggests that having angioplasty reduces the risk of death by a factor of almost 4.*

    (b)  Now, control for age (-age_c2-). Is age a confounder? If so, why? If not, why not?

```
. cs died_1yr ptca, by(age_c2)
        age_c2 |       RR        [95% Conf. Interval]   M-H Weight
---------------+------------------------------------------------------
            0 |    .3504941      .2689943      .4566868      79.29104
            1 |    .3965193      .3213446      .4892803      139.3606
---------------+------------------------------------------------------
        Crude |    .2908427      .2479511      .3411538
  M-H combined |    .3798289      .3220408      .4479867
```

```
-------------------------------------------------------------------------
Test of homogeneity (M-H)        chi2(1) =      0.515  Pr>chi2 = 0.4729
```
*The RR has now risen from 0.29 to 0.38, suggesting a smaller beneficial effect from angioplasty (although it is still substantial, and significant (CI does not include 1.0)). The reason for this change in the RR is that -age_c2- is strongly (negatively) associated with -ptca- (RR=0.46) and strongly (positively) associated with died_1yr (RR=2.05 in individuals without angioplasty). Results of these additional analyses are not shown, but the Stata code to perform them is:*
```
cs ptca age_c2
cs died_1yr age_c2 if ptca==0
```

(c) Next, control for whether or not the individual had had a previous MI (-prmi-). Is -prmi- a confounder? Is there evidence of interaction? What do these results tell you about the effect of -ptca-? If there is interaction, is there evidence of antagonism or synergism?

```
.           cs died_1yr ptca, by(prmi)
     previous MI |       RR        [95% Conf. Interval]    M-H Weight
-----------------+-------------------------------------------------------
             0 |    .2373567      .1925501    .2925899       197.1953
             1 |     .457232      .3594987    .5815351        85.99521
-----------------+-------------------------------------------------------
         Crude |    .2908427      .2479511    .3411538
   M-H combined |    .3041253      .2599939    .3557476
-------------------------------------------------------------------------
Test of homogeneity (M-H)        chi2(1) =     16.437  Pr>chi2 = 0.0001
```
*On the surface, it appears that -prmi- is not a confounder because the crude and M-H adjusted RRs are very close. However, it is evident that the effect of -ptca- is very different in individuals without -prmi-=0 (RR=0.24) compared with individuals with -prmi-=1 (RR=0.46). Angioplasty has a much greater protective effect in individuals who have not had a previous MI. The difference between the two RR is statistically significant (M-H test of homogeneity P-value is 0.0001).*

*This is evidence of antagonism. The effect of -ptca- by itself (ie when -prmi-=0) is greater than its effect when -prmi-=1.*

(d) Evaluate interaction on the additive scale. (Note – to do this you will need to determine the risk difference (RD) separately for individuals with prmi=0 and prmi=1).

*Details of the anlyses are not presented, but the code to carry them out was:*
```
cs died_1yr ptca if prmi==0
cs died_1yr ptca if prmi==1
```
*To evaluate interaction on the additive scale, we are interested in whether or not the RD for -ptca- differs between the 2 levels of -prmi-. When -prmi-=0, the RD is -0.28 while when -prmi-=1, the RD is -0.24. The very limited difference suggests that there is no interaction on the additive scale (at both levels of -prmi-, -ptca- reduces the risk of death by ~26 cases per 100 patients). It is common to observe that there is antagonism on the multiplicative scale when there is no interaction on the additive scale.*

7. We will now switch to using propensity scores to help us evaluate the effects of -ptca-. Because angioplasty was not randomly assigned to individuals, it is very likely that those who had an angioplasty while in the hospital were different in some (or many) ways from those who did not get an angioplasty. Consequently, a simple comparison of ptca=0 vs ptca=1 is probably not appropriate. We would like to account for differences in the two groups in terms of: -sex- -age_c5-

-white- -obese- -prmi-. The coding required to carry out the propensity analyses is a bit complicated, but the question we are addressing is very similar to the one used in Examples 13.10 – 13.12, so you can find Stata code for these analyses with the code for all of the analyses in Chapter 13.

(a) First eliminate all records with any missing values for -ptca- or -sex- -age_c5- -white- -obese- -prmi-. (These will have to be ignored in the analysis anyway). Then compute propensity scores and determine if they meet the property of being balanced. When doing this, limit the computations to observations which fall in the region of "common support"

*The Stata code for doing the necessary calculations is:*

```
* propensity scores
* remove all obs. with missing values for any predictors
keep if !missing(ptca) & !missing(sex) & !missing(age_c5) & !missing(white) ///
    & !missing(obese) & !missing(prmi)
    * this is done solely because the help file for this
    * program suggests it should be done
* first compute unconditional association between -died_1yr- and -ptca-
    cs died_1yr ptca
* compute propensity scores (based on -sex- -age_c5- -white- -obese- -prmi-)
    capture drop ps
    capture drop block
    xi:pscore ptca sex i.age_c5 white obese prmi, pscore(ps) logit comsup ///
        blockid(block) numblo(20)
```

*The results suggest that the scores do not achieve the desirable property of being balanced. However, closer inspection of the results (using the -detail- option) to find variables/blocks not balanced reveals that only -prmi- is unbalanced and this happens in only 1 of 16 blocks. Even in that block, it is only barely declared unbalanced (P=0.009 with default threshold set to 0.01). Consequently, it appears that to a very large degree, the property of "balancedness" has been met and we will proceed with these propensity scores.*

(b) What is the range of propensity scores computed? Are there many observations outside the region of common support? Recompute the crude association between -died_1yr- and -ptca- using only the observations that fall in the region of common support (for comparison with later results).

```
. codebook ps if comsup==1

---------------------------------------------------------------------------
ps                                                 Estimated propensity score
---------------------------------------------------------------------------
                type:  numeric (double)
               range:  [.05890814,.76068548]        units:  1.000e-09
       unique values:  74                        missing .:  0/2892
                mean:  .466423
           std. dev:  .202759
         percentiles:      10%       25%       50%       75%       90%
                       .220397   .304169   .455606   .667969   .755354


. tab block, sum(ps)


            |   Summary of Estimated propensity
  Number of |              score
      block |       Mean   Std. Dev.        Freq.
```

```
-----------+---------------------------------
        2 |    .07610165    .01017475          191
        3 |    .11781061    .0010385            37
        4 |    .18538051    .00191535           13
        5 |    .22338842    .00814728          167
        6 |    .28031572    .005633            313
        7 |     .314057     .01626669          145
        8 |    .37156658    .01067522          315
        9 |    .42835919    .01028557          255
       10 |    .46710229    .01503808          127
       11 |    .53452504    .00669521          260
       12 |    .58107947    .01354079          193
       13 |    .62512865    .00636353           68
       14 |    .67564253    .0103768           395
       15 |    .71881572    .0026028            58
       16 |    .75832771    .00265154          355
-----------+---------------------------------
    Total |    .46642269    .20275856         2892
```

*The propensity scores that fall within the region of common support range from 0.059 to 0.761. It turns out that there were only 2 observations which fell outside the region of common support and these were 2 individuals that had very low propensity scores (<0.058) with the lowest value observed in individuals who had an angioplasty being 0.059. (This was determined by visually inspecting the data.)*

*Note: If the crude effect of angioplasty is determined after excluding those 2 individuals outside the region of common support, the RD is -0.275 and the OR is 0.202. (We have switched to using OR because in part (e) we will fit some logistic regression models and they produce OR - which we will want to compare with the crude estimate).*

(c) Use the propensity scores to carry out nearest neighbour matching. Does this change the estimate of the effect of -ptca- by much?

```
. attnd died_1yr ptca ,  comsup pscore(ps) detail
… some output omitted

Average outcome of the matched treated
    Variable |       Obs        Mean    Std. Dev.        Min        Max
-------------+----------------------------------------------------------
    died_1yr |      1349     .1141586    .3181218          0          1

Average outcome of the matched controls
    Variable |       Obs      Weight        Mean   Std. Dev.       Min        Max
-------------+------------------------------------------------------------------
    died_1yr |      1528  1348.99995    .2641203    .441008          0          1

ATT estimation with Nearest Neighbor Matching method
(random draw version)
Analytical standard errors
-------------------------------------------------------
n. treat.   n. contr.        ATT     Std. Err.          t
-------------------------------------------------------
     1349        1528      -0.150        0.020     -7.588
-------------------------------------------------------
```

*The risk of death in the treated (-ptca-=1) individuals was 0.114 while in the matched controls (note treated individuals were matched with >1 control) was 0.264 resulting in a RD (also known*

*as ATT – average treatment effect in treated individuals) of -0.150 (ie treatment reduced the risk of death by 15 cases per 100 patients). This is much less than the crude RD of -0.275.*

(d) Use the propensity scores to carry out a stratified analysis (based on the blocks created by the process of generating the propensity scores). Does this change the estimate of the effect of -ptca- by much.

```
. atts died_1yr ptca, pscore(ps) blockid(block) comsup detail
… most output omitted

ATT estimation with the Stratification method
Analytical standard errors
-----------------------------------------------------------
n. treat.   n. contr.        ATT   Std. Err.            t
-----------------------------------------------------------
     1349        1543     -0.157       0.015      -10.128
-----------------------------------------------------------
```

*The ATT is -0.157, which is very close to what was observed by using propensity scores for matching.*

(e) Fit logistic models which control for the factors of interest by:

*Note – This question uses logistic regression models which are covered in Chapter 16. If you are not familiar with these types of models, you can skip this section.*

   i.  including propensity scores as a continuous predictor in the model,

```
. logit died_1yr ptca ps if comsup, or
… some output omitted

Logistic regression                              Number of obs   =       2892
                                                 LR chi2(2)      =     454.09
                                                 Prob > chi2     =     0.0000
Log likelihood = -1433.4488                      Pseudo R2       =     0.1367
-----------------------------------------------------------------------------
   died_1yr | Odds Ratio   Std. Err.      z    P>|z|     [95% Conf. Interval]
------------+----------------------------------------------------------------
       ptca |   .3100939   .0331103   -10.97   0.000     .2515396    .3822789
         ps |   .0485518   .0122562   -11.98   0.000     .0296027    .0796304
      _cons |   1.991847   .2123742     6.46   0.000     1.616216     2.45478
-----------------------------------------------------------------------------
```

*After adjustment for propensity scores the OR rises to 0.31 (from 0.20) suggesting that "propensity for treatment" was a confounder.*

   ii.  including propensity scores as a categorical predictor (blocks) in the model, and

```
. logit died_1yr ptca i.block if comsup, or        /* OR = 0.311 */
… some output omitted

Logistic regression                              Number of obs   =       2892
                                                 LR chi2(15)     =     503.13
                                                 Prob > chi2     =     0.0000
Log likelihood = -1408.9265                      Pseudo R2       =     0.1515
-----------------------------------------------------------------------------
   died_1yr | Odds Ratio   Std. Err.      z    P>|z|     [95% Conf. Interval]
------------+----------------------------------------------------------------
```

```
        ptca |    .3110766    .0334658    -10.85   0.000     .2519385    .3840963
             |
       block |
           3 |    1.583558    .5866795      1.24   0.215     .7660887    3.273321
           4 |    .5685746    .3598982     -0.89   0.372     .1644338    1.966001
… some output omitted
          15 |    .1937204    .0967477     -3.29   0.001     .0727896    .5155629
          16 |    .0949469    .0285414     -7.83   0.000     .0526753    .1711413
             |
       _cons |     1.05361    .1542291      0.36   0.721      .790824    1.403719
-------------------------------------------------------------------------------
```

*The OR=0.311 which is very similar to what was seen with propensity scores added as actual values.*

> iii. including the original predictors (-sex- -age_c5- -white- -obese- -prmi-) directly in the model instead of the propensity scores. Do these generally produce similar results? Are the results much different from what you get by ignoring this set of variables?

```
. logit died_1yr ptca sex i.age_c5 white obese prmi if comsup, or
… some output omitted

Logistic regression                              Number of obs   =       2892
                                                 LR chi2(9)      =     507.09
                                                 Prob > chi2     =     0.0000
Log likelihood = -1406.9497                      Pseudo R2       =     0.1527
-------------------------------------------------------------------------------
    died_1yr | Odds Ratio   Std. Err.      z    P>|z|     [95% Conf. Interval]
-------------+-----------------------------------------------------------------
        ptca |    .3027872    .0325903    -11.10   0.000     .2451991    .3739004
         sex |    1.198485    .1154802      1.88   0.060     .9922355    1.447605
             |
       age_c5 |
           1 |    1.977847    .3927339      3.43   0.001     1.340209    2.918856
           2 |    3.139902    .5582562      6.44   0.000     2.216039    4.448921
           3 |     5.62758    .9752171      9.97   0.000     4.006959    7.903664
           4 |    6.667302    1.424016      8.88   0.000     4.386809    10.13332
             |
       white |     .824577    .1201627     -1.32   0.186     .6197102     1.09717
       obese |    .6826442    .0794924     -3.28   0.001     .5433427    .8576596
        prmi |    1.407338    .1394872      3.45   0.001     1.158863    1.709089
       _cons |    .1808426     .038537     -8.03   0.000     .1191001    .2745929
-------------------------------------------------------------------------------
```

*Not surprisingly, using the original covariates for adjustment, instead of propensity scores produces a very similar result (OR=0.30). Given the large number of observations and the relatively small number of covariates adjusted for, the 2 approaches yield similar results.*

*It makes virtually no difference if this final analysis is limited to the region of common support or not, because there were only 2 observations outside this region (results not shown). If the number of observations outside the region of common support was much larger, then the discrepancy between the propensity score approach and adjusting for covariates directly would likely be larger.*