

## An Explanation for CEM Weights

[Gary King](#)

25 March 2012

([j.mp/CEMweights](http://j.mp/CEMweights))

**Goal:** With a [CEM](#) matched sample, the goal is to estimate the ATT (the average treatment effect on the treated). To compute the ATT, we can literally compute the causal effect for each of the treated units in the sample and then average them all. We can do this easily without weights in an easy-to-understand, but involved, procedure. What then is the point of the weights? They enable us to use a calculation trick that makes it easy to estimate the ATT in a weighted least squares regression program without the involved procedure. In this brief document, I'll show both ways. (For articles on CEM, matching, and causal inference, see [this link](#).)

**Example:** Let  $s$ =stratum (within each of which, all the coarsened covariates match),  $Y$ =dependent variable,  $T$ =treatment (1=treated, 0=control). Here's an example with 1 treated and 2 controls in the first stratum and 2 treateds and 2 controls in the second:

$s$	$Y$	$T$
1	1	1
1	2	0
1	3	0

2	4	1
2	5	1
2	6	0
2	7	0

To compute the causal effect for the first treated unit in the first stratum, we take the value of  $Y$  for that unit (which is 1) and subtract the average value of  $Y$  for the control units in that stratum (which is  $(2+3)/2=2.5$ ). So then the causal effect for the first unit is  $1-2.5=-1.5$ . This is just arithmetic; no magical statistics. It should be very intuitive: just a *difference in means within the stratum* that CEM defines as having units that are all essentially the same (with respect to background covariates).

To compute the causal effect for the 2 treated units in the second stratum, we take the average value of  $Y$  for the two treated units and subtract the average value of  $Y$  for the two control units:  $4.5 - 6.5 = -2$ . Again, it's the difference in means within the stratum.

What then is the ATT in the entire sample? We just average the causal effects for the 3 treated units, one of which is in the first stratum and two of which are in the second. Within each stratum, we have only one causal effect for all the observations. Thus, the calculation is to average the causal effects with one from the first stratum (at -1.5) and two from the second (at -2). The result is:  $(-1.5 -2 -2)/3 = -1.833$ .

**The Role of the weights:** The ATT for the entire sample is the goal of the whole analysis. If you'd prefer to compute this as done above, without the weights, that's perfectly fine. However, it's usually easier to use the weights, even if they are harder to understand directly. The best way to understand them is merely as a trick that enables you to use a standard weighted least squares program to do what would otherwise be a somewhat involved calculation.

Weights in CEM are defined as  $W = 1$  for treated units and  $(m_C/m_T)*W_s$  for control units -- where  $m_C$  and  $m_T$  are the numbers of controls and treateds in the sample, and  $m_T^s$  and  $m_C^s$  are the number of treateds and controls in stratum  $s$ , and  $W_s = (m_T^s)/(m_C^s)$ . Here's the same example with the addition of  $W_s$  and  $W$ :

**s Y T Ws W**

```
1 1 1 1 1
1 2 0 1/2 (4/3)(1/2)=2/3
1 3 0 1/2 (4/3)(1/2)=2/3

2 4 1 1 1
2 5 1 1 1
2 6 0 1 4/3
2 7 0 1 4/3
```

The weights have two components.  $Ws$  is the *unnormalized weights*, the part that varies over strata so that the sum of the control units equals the number of treated units. We can see this in the example: For stratum 1, the number of treated units is 1, and the sum of  $Ws$  for the control units equals  $1/2+1/2=1$ . For stratum 2, the number of treated units is 2 and the sum of  $Ws$  for the control units is also 2.

We could use  $Ws$  as the weights except that the sum of the weights must always equal the number of observations in the sample; if this isn't true, the standard errors will be incorrect (we'd be lying to the computer program, telling it we have a different number of observations than we really do). And in the example,  $\text{sum}(Ws) = 6$ , but  $n=7$ . The normalization factor,  $(m_C/m_T)$ , fixes this problem by taking all values of  $Ws$  and scaling them up to the total  $n$ ; in the example the normalization factor is  $4/3$  and must be multiplied by every

unnormalized weight. And here we can see that  $\sum_{i=1}^n W_i = n = 7$ , as should be the case.

So now that we have the (normalized) CEM weights,  $W$ , we can use a standard WLS program to compute the ATT in the entire sample, without having to go through the easy-to-understand but involved procedure I used above. So here's an example in R:

```
y=c(1,2,3,4,5,6,7)
t=c(1,0,0,1,1,0,0)
w=c(1,2/3,2/3,1,1,4/3,4/3)
lm(y~t,weight=w)
```

```
Call:
lm(formula = y ~ t, weights = w)
```

```
Coefficients:
(Intercept)      t
      5.167      -1.833
```

Which is the identical number (-1.833) as we got above with the simple-to-understand but more involved procedure.

The same logic above can be applied to adjusting for pre-treatment covariates. Here it's very easy with the weights to merely run the weighted least squares regression with the additional covariates.