

# Ec533: Labour Economics for Research Students

## Regression versus matching

Jörn-Steffen Pischke

LSE

October 29, 2009

You want to estimate the effect of treatment on the treated

$$\delta_{TOT} \equiv E[Y_{1i} - Y_{0i} | D_i = 1].$$

Using the law of iterated expectations

$$\delta_{TOT} = E\{E[Y_{1i} | X_i, D_i = 1] - E[Y_{0i} | X_i, D_i = 1] | D_i = 1\}$$

Under conditional independence (CIA):

$$E[Y_{0i} | X_i, D_i = 0] = E[Y_{0i} | X_i, D_i = 1].$$

# The matching estimand

Using this we get

$$\begin{aligned}\delta_{TOT} &= E \{ E [Y_{1i} | X_i, D_i = 1] - E [Y_{0i} | X_i, D_i = 0] | D_i = 1 \} \\ &= E \{ E [y_i | X_i, D_i = 1] - E [y_i | X_i, D_i = 0] | D_i = 1 \} \\ &= E [\delta_X | D_i = 1],\end{aligned}$$

where

$$\delta_X \equiv E [y_i | X_i, D_i = 1] - E [y_i | X_i, D_i = 0]$$

is an  $X$ -specific difference in means at covariate value  $X_j$ .

With discrete  $X_j$ , the matching estimand is

$$\delta_M = \sum_x \delta_x P(X_j = x | D_i = 1),$$

where  $P(X_j = x | D_i = 1)$  is the probability mass function for  $X_j$  given  $D_i = 1$ .

# Rewriting the matching estimand

Using Bayes Rule

$$P(X_i = x | D_i = 1) = \frac{P(D_i = 1 | X_i = x) \cdot P(X_i = x)}{P(D_i = 1)}.$$

the matching estimand can be written as

$$\begin{aligned} \delta_M &= \sum_x \delta_x P(X_i = x | D_i = 1) \\ &= \frac{\sum_x \delta_x P(D_i = 1 | X_i = x) P(X_i = x)}{\sum_x P(D_i = 1 | X_i = x) P(X_i = x)} \end{aligned}$$

Suppose we run instead the regression

$$y_i = \sum_x d_{ix} \beta_x + \delta_R D_i + \varepsilon_i,$$

where  $d_{ix}$  is a dummy that indicates  $X_i = x$ ,  $\beta_x$  is a regression-effect for  $X_i = x$ , and  $\delta_R$  is the regression estimand. Note that this regression model allows a separate parameter for every value taken on by the covariates.

It can be shown that  $\delta_R$  can be written as

$$\delta_R = \frac{\sum_x \delta_x [P(D_i = 1|X_i = x)(1 - P(D_i = 1|X_i = x))] P(X_i = x)}{\sum_x [P(D_i = 1|X_i = x)(1 - P(D_i = 1|X_i = x))] P(X_i = x)}$$

# Regression vs. matching

So the matching estimand

$$\delta_M = \sum_x \delta_x \left[ \frac{P(D_i = 1|X_i = x)P(X_i = x)}{\sum_x P(D_i = 1|X_i = x)P(X_i = x)} \right]$$

and the regression estimand

$$\delta_R = \sum_x \delta_x \left[ \frac{[P(D_i = 1|X_i = x)(1 - P(D_i = 1|X_i = x))] P(X_i = x)}{\sum_x [P(D_i = 1|X_i = x)(1 - P(D_i = 1|X_i = x))] P(X_i = x)} \right]$$

differ by the weights they use to combine the covariate specific treatment effects  $\delta_x$ .

# Regression vs. matching in words

- Matching uses weights which depend on

$$P(D_i = 1|X_i = x)$$

i.e. the fraction of treated observations in a covariate cell (or the mean of  $D_i$ ). This is larger in cells where there are many treated observations (makes sense as we want the effect of treatment on the treated).

- Regression uses weights which depend on

$$P(D_i = 1|X_i = x)(1 - P(D_i = 1|X_i = x))$$

i.e. the variance of  $D_i$  in the covariate cell. This weight is largest in cells where there are half treated and half untreated observations. This also makes sense because these cells will produce the lowest variance estimates of  $\delta_x$ . If all the  $\delta_x$  are the same, the most efficient estimand uses the lowest variance cells most heavily. This is what regression does.