

11. Matching Estimators

ISS5096 || ECI

Jaewon (“Jay-one”) Yoo

National Tsing Hua University

It's a Match!

You and MHE have liked each other!



Pixiz

Credit: Created using Pixiz (<https://en.pixiz.com/template/it-s-a-Match-Tinder-mockery-with-customizable-text-3150>)

Where are we? Where are we going?

- The recurring question of this course: **where do we find good controls?**
 - Units *randomized* to receive control (RCT)
 - Units with *similar* values of covariates (SOO / reg. adj.)
 - Units with *opposite* value of some instrument (IV)
 - Units *near a discontinuity* in treatment assignment (RDD, W13)
 - Exploit two possible sources of variation for identification:
 - Exploit **cross-sectional** variation in treatment.
 - Exploit variation in treatment **within a unit over time** (before/after)
- Can we make our identification strategies work *better*?
 - \rightsquigarrow **matching** or **weighting**
 - Today: refine the SOO branch by pruning or re-weighting observational controls.

1/ Matching Estimators

The Problem with Regression

- Causal inference is all about comparing **counterfactuals**, like the ATT:

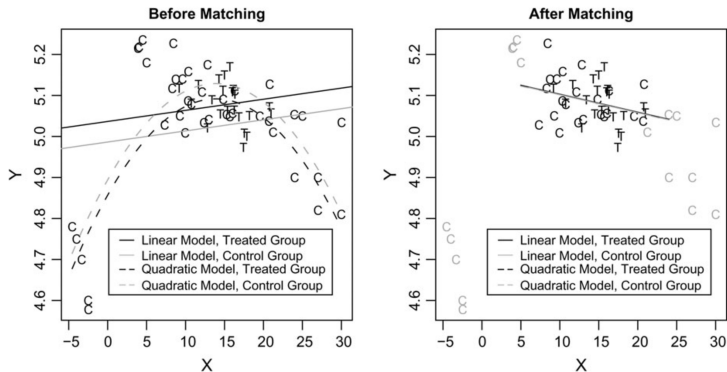
$$\tau_{\text{ATT}} = \mathbb{E}[Y_i(1) - Y_i(0) \mid D_i = 1]$$

- Recall the **imputation** estimators with regression (ECI W5).

$$\hat{\tau}_{\text{reg}} = \frac{1}{n_1} \sum_{i=1}^n D_i (Y_i - \hat{\mu}_0(\mathbf{x}_i))$$

- Common solution: use a parametric model for $\hat{\mu}_0(\mathbf{x}_i)$.
 - For example, could assume it is linear: $\mu_0(\mathbf{x}) = \mathbf{x}'\beta$
 - Regression, MLE, Bayes, etc.
 - But this model might be wrong \rightsquigarrow wrong causal estimates.

Model Dependence



Source: *Figure 1* in Ho, Daniel E., Kosuke Imai, Gary King, and Elizabeth A. Stuart. "Matching as nonparametric preprocessing for reducing model dependence in parametric causal inference." *Political analysis* 15, no. 3 (2007): 199-236.

What is Matching?

- **Matching** is a nonparametric imputation estimator:

$$\hat{\tau}_m = \frac{1}{n_1} \sum_{i=1}^n D_i \left(Y_i - \frac{1}{M} \sum_{j \in \mathcal{J}_M(i)} Y_j \right)$$

- $\mathcal{J}(i)$ are the set of M closest control units to i in terms of \mathbf{X}_i
 - Matching? = pruning bad controls (King, 2020)
- Matching has strong advantages:
 1. Reduces dependence of estimates on parametric models.
 2. Reduces model-based extrapolation.
 3. Makes counterfactual comparisons more transparent.
- What matching isn't: a solution for selection on unobservables.
 - Matching is an **estimation** technique, not an identification strategy.

Types of Matching

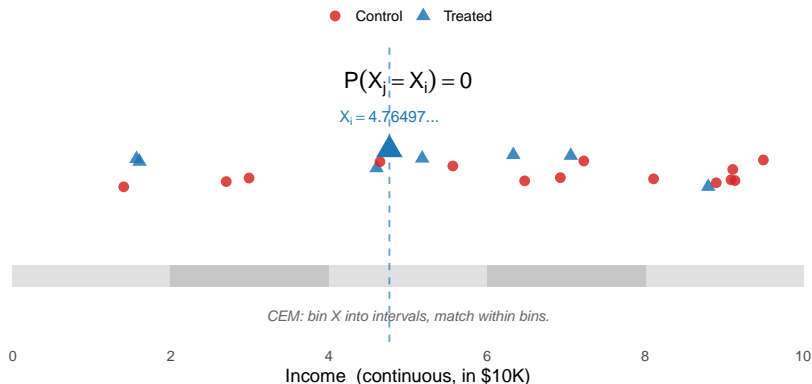
- Assumptions:
 - No unmeasured confounders: $D_i \perp\!\!\!\perp (Y_i(0), Y_i(1)) \mid \mathbf{X}_i$
 - Overlap/positivity: $0 < \mathbb{P}(D_i = 1 \mid \mathbf{X}_i = \mathbf{x}) < 1$
- **Exact matching**: choose matches that have the same value of \mathbf{X}_i .
 - $\mathcal{J}_M(i)$ is a random set of M control units with $\mathbf{X}_j = \mathbf{X}_i$
 - Covariate distribution in treated and matched controls exactly the same:

$$\begin{aligned}\widehat{\mathbb{P}}(\mathbf{X}_j = \mathbf{x} \mid D_j = 1) &= \widehat{\mathbb{P}}(\mathbf{X}_j = \mathbf{x} \mid D_j = 0, j \text{ is matched}) \\ \rightsquigarrow \mathbb{E}[Y_i(0) \mid D_i = 1] &= \mathbb{E}[Y_j(0) \mid D_j = 0, j \text{ is matched}]\end{aligned}$$

- Issue: not feasible with high-dimensional or continuous \mathbf{X}_i
- **Coarsened exact matching** (Iacus et al, 2011)
 - Discretize and group covariates into substantively meaningful bins
 - Exact match on these bins \rightsquigarrow accounts for interactions
 - Have to drop treated units in bins with no controls \rightsquigarrow changes estimand
 - Allows you to control bias/variance tradeoff through coarsening

Continuous X: Exact Matching Fails

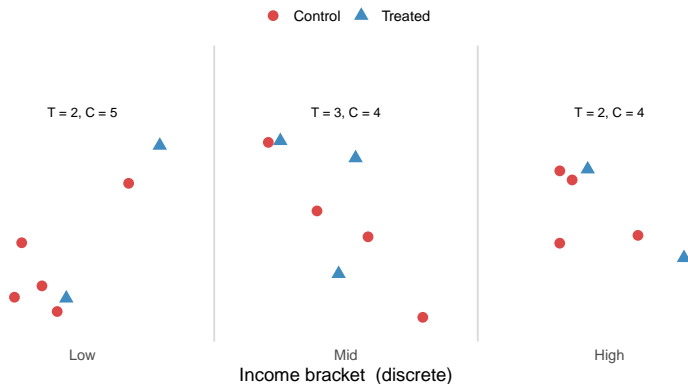
- Say \mathbf{X} = income (continuous). Then $X_j \neq X_i$ “almost surely” (a.s.), i.e., $\mathbb{P}(X_j = X_i) = 0$:



- Fix: **coarsen X** into bins (CEM) and match within each bin.

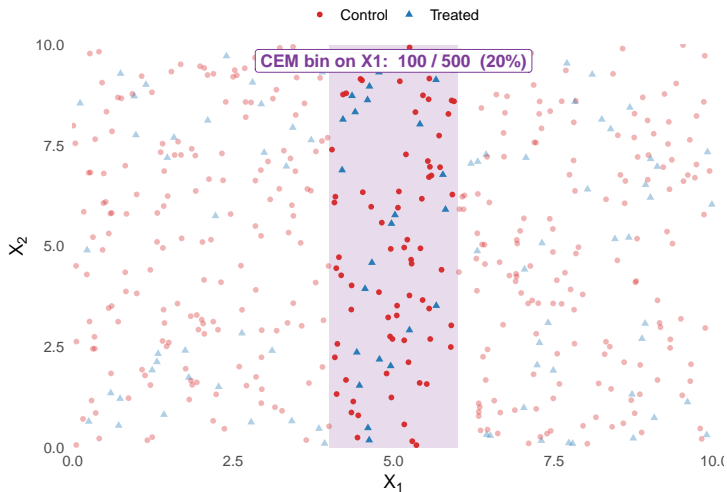
CEM: Bin Continuous X into Brackets

- Income brackets {Low, Mid, High} are a CEM coarsening. Pair each treated with a control in the **same bracket**:



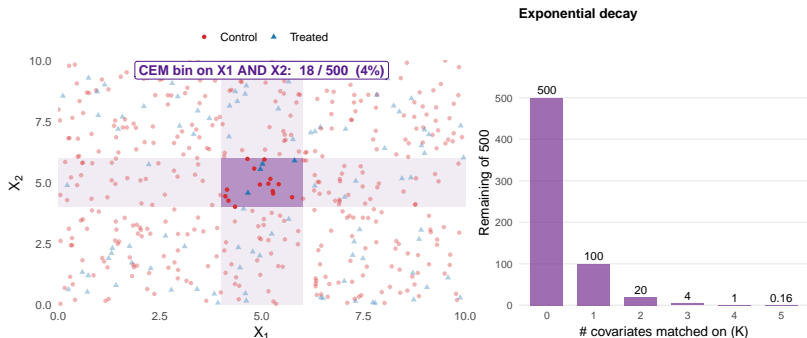
Curse of Dimensionality: CEM on One Covariate

- Coarsen X_1 into bins; keep observations in one bin:



Curse of Dimensionality: CEM on Multiple Covariates

- Now require closeness on *both* X_1 and X_2 (intersection of two bins):



- Matching breaks when a treated unit's bin has no controls. As K grows, this becomes increasingly common.

Matching in High Dimensions

- Even CEM can break down with high-dimensional \mathbf{X}_j .
- We can define closeness using lower dimensional **distance metrics**
 - Reducing dimensionality: mapping two vectors to a single number.

- **Mahalanobis distance:**

$$D(\mathbf{x}_i, \mathbf{x}_j) = \sqrt{(\mathbf{x}_i - \mathbf{x}_j)' \widehat{\Sigma}^{-1} (\mathbf{x}_i - \mathbf{x}_j)}$$

- $\widehat{\Sigma}$ is the estimated variance-covariance matrix of the observations:

$$\widehat{\Sigma} = \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})^T$$

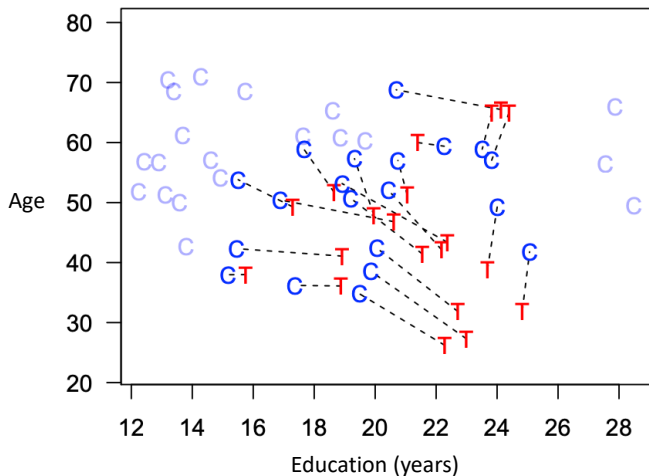
- **Estimated propensity score** (Rosenbaum and Rubin, 1983):

$$D(\mathbf{x}_i, \mathbf{x}_j) = |\widehat{\pi}(\mathbf{x}_i) - \widehat{\pi}(\mathbf{x}_j)| = |\widehat{\mathbb{P}}(D_i = 1 | \mathbf{x}_i) - \widehat{\mathbb{P}}(D_j = 1 | \mathbf{x}_j)|$$

- Some use linear predictor: $\text{Dist}_{ij} = |\text{logit}(\widehat{\pi}(\mathbf{x}_i)) - \text{logit}(\widehat{\pi}(\mathbf{x}_j))|$

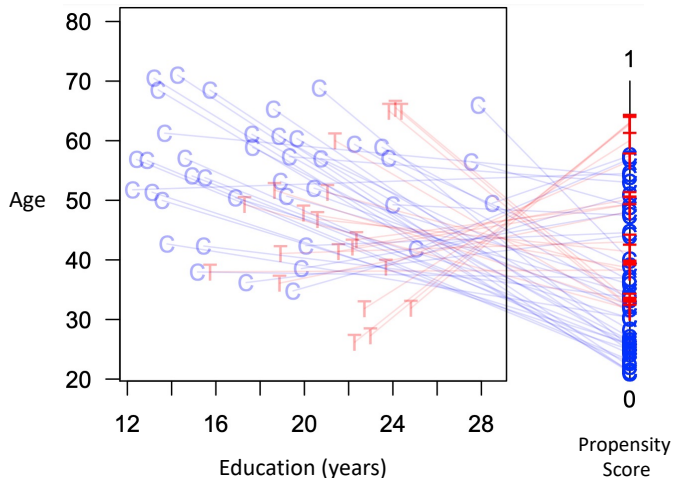
Mahalanobis Distance Matching

- **Mahalanobis Distance Matching:**
 - Prune observations where $\text{Dist}_{ij} > \text{caliper}$ (\rightsquigarrow caliper? = cutoff point)



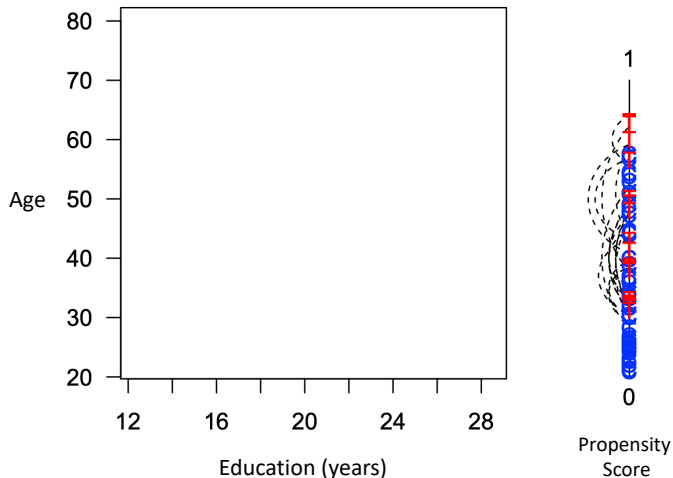
Credit: Figure from Gary King's talk on "The Balance-Sample Size Frontier in Matching Methods for Causal Inference," at University of Michigan, January 24, 2014.

Propensity Score Matching Illustrated



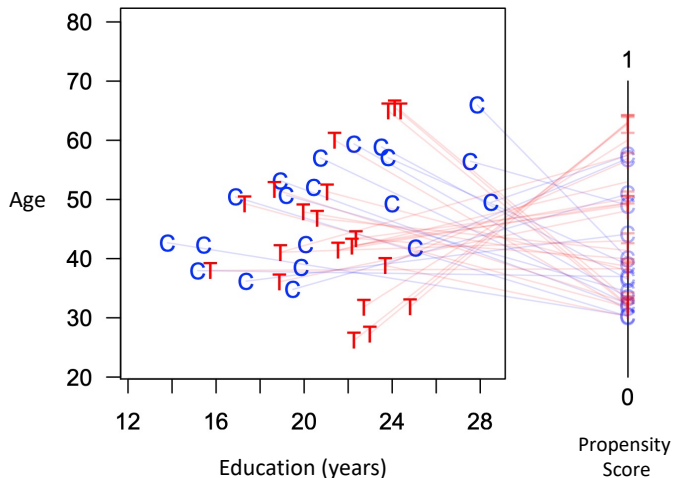
Credit: Figure from Gary King's talk on "The Balance-Sample Size Frontier in Matching Methods for Causal Inference," at University of Michigan, January 24, 2014.

Propensity Score Matching Illustrated



Credit: Figure from Gary King's talk on "The Balance-Sample Size Frontier in Matching Methods for Causal Inference," at University of Michigan, January 24, 2014.

Propensity Score Matching Illustrated



Credit: Figure from Gary King's talk on "The Balance-Sample Size Frontier in Matching Methods for Causal Inference," at University of Michigan, January 24, 2014.

Other Matching Choices

- **Matching ratio:** how many control units per treated unit?
 - Lower reduces bias (only use the closest matches)
 - Lower increases variance
- **With or without replacement:** can the same control be matched to multiple treated?
 - **With** replacement: each treated picks its closest control independently \rightsquigarrow better matches, order doesn't matter. But the same control reappears in many pairs \rightsquigarrow pairs are not independent \rightsquigarrow SE must **cluster on the matched control** (Abadie & Imbens, 2006).
 - **Without** replacement: each control used at most once \rightsquigarrow no within-control duplication \rightsquigarrow standard (non-clustered) SE formulas apply.
- **Caliper:** drop poor matches?
 - Only keep matches below a distance threshold, $D(\mathbf{X}_i, \mathbf{X}_j) \leq c$
 - Rosenbaum and Rubin (1985): use c size equiv. to $0.25 \times \text{sd}$ of the PS.
 - Reduces imbalance, but if you drop treated units, estimand changes.
 - \rightsquigarrow If we drop treated units, what are we estimating other than the ATT?

More on Propensity Scores

- **Balancing theorem (Rosenbaum and Rubin, 1983)**. Covariates are balanced conditional on the true propensity score:

$$D_i \perp\!\!\!\perp \mathbf{X}_i \mid \pi(\mathbf{X}_i)$$

- Together with no unmeasured confounders, this implies we only need to match or balance on $\pi(\mathbf{x})$:

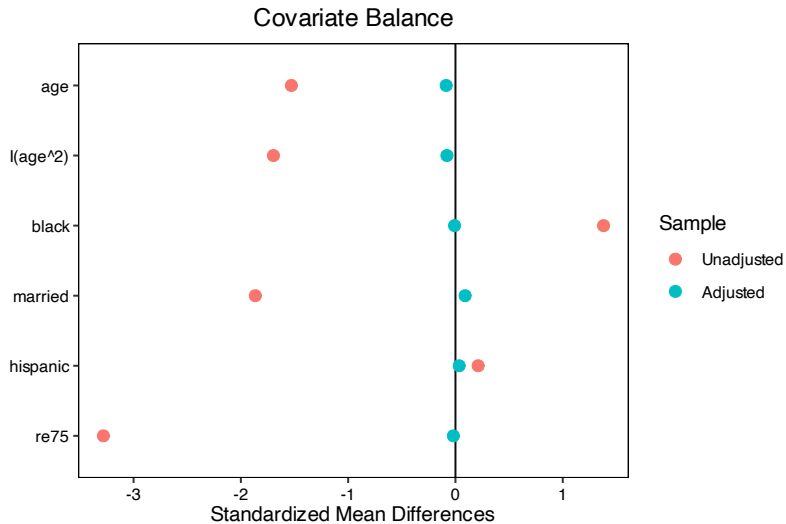
$$\underbrace{(Y_i(0), Y_i(1)) \perp\!\!\!\perp D_i \mid \mathbf{X}_i}_{\text{conditional unconfoundedness}} \iff (Y_i(0), Y_i(1)) \perp\!\!\!\perp D_i \mid \pi(\mathbf{X}_i)$$

- In observational data, we never know the true $\pi(\mathbf{x}) \rightsquigarrow$ estimate $\hat{\pi}(\mathbf{x})$.
- Is balancing on $\hat{\pi}(\mathbf{x})$ sufficient? **No idea!**
 - Have to check whether \mathbf{X}_i is actually balanced.
 - Somewhat deflates the benefits of PS balancing/matching.
- \rightsquigarrow “propensity score tautology”.

Assessing Balance

- Goal of matching is to maximize balance: $\widehat{F}_1(\mathbf{x}) \approx \widehat{F}_{0,\mathcal{J}}(\mathbf{x})$
 - Joint distribution of \mathbf{X}_j is similar between treated and matched controls.
 - Difficult to assess balance across many dimensions \rightsquigarrow summaries.
- Options:
 - Differences-in-means/medians, standardized.
 - Visualize balance using Love plot.
 - QQ plots/KS statistics for comparing the entire distribution of \mathbf{X}_j .
- Hypothesis tests for balance are problematic:
 - Dropping units can lower power (\uparrow p-values) without a change in balance.

Balance Diagnostics: Love Plot



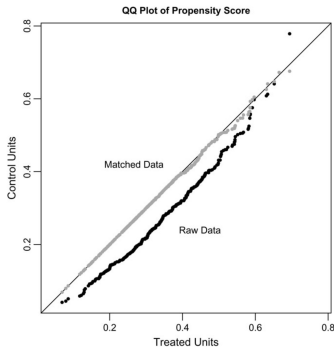
Balance Diagnostics: Table and Density

| | Nonadopters | | Adopters | | Differences in Means |
|---|-------------|---------|----------|---------|----------------------|
| | Mean | Std dev | Mean | Std dev | |
| Customer spending: | | | | | |
| Purchase amount (monetary) | 21.133 | 14.251 | 21.207 | 15.407 | -0.074 |
| Number of transactions (frequency) | 1.318 | 0.578 | 1.321 | 0.536 | -0.003 |
| Number of transactions per trip | 1.161 | 0.450 | 1.149 | 0.422 | 0.012 |
| Days since last purchase (recency) | 118.011 | 83.374 | 117.085 | 85.634 | 0.926 |
| Number of books purchased (quantity) | 1.819 | 1.173 | 1.848 | 1.256 | -0.029 |
| Maximum price of the books purchased | 12.915 | 5.782 | 12.957 | 6.311 | -0.042 |
| Consumption variety: | | | | | |
| Number of unique books | 1.733 | 1.101 | 1.766 | 1.149 | -0.033 |
| Number of unique genres | 1.272 | 0.511 | 1.286 | 0.512 | -0.014 |
| Number of unique authors | 1.975 | 1.448 | 2.019 | 1.473 | -0.044 |
| Number of unique publishers | 1.685 | 0.988 | 1.696 | 0.970 | -0.011 |
| Concentration on personal favorites: | | | | | |
| Genre concentration (norm. HHI) | 0.907 | 0.161 | 0.903 | 0.156 | 0.004 |
| Author concentration (norm. HHI) | 0.753 | 0.266 | 0.750 | 0.256 | 0.003 |
| Publisher concentration (norm. HHI) | 0.788 | 0.239 | 0.791 | 0.228 | -0.003 |
| Best-seller purchases: | | | | | |
| Shares of top 10 best sellers | 0.076 | 0.219 | 0.074 | 0.201 | 0.002 |
| Shares of top 50 best sellers | 0.147 | 0.298 | 0.142 | 0.272 | 0.005 |
| Usage of Sales Channels: | | | | | |
| Share of in-store mobile sales | 0.147 | 0.328 | 0.190 | 0.362 | -0.043 |
| Share of m-commerce sales | 0.071 | 0.237 | 0.096 | 0.268 | -0.025 |
| Share of e-commerce sales | 0.086 | 0.261 | 0.094 | 0.267 | -0.008 |
| Share of cash register sales | 0.103 | 0.279 | 0.107 | 0.278 | -0.004 |
| Customer demographics: | | | | | |
| Age | 31.193 | 7.310 | 31.151 | 7.394 | 0.042 |
| Gender (male = 1) | 0.333 | 0.471 | 0.341 | 0.474 | -0.008 |

Note. Means and standard deviations of the observed variables, calculated in the 11-month pre-introduction period before February 1, 2015. The unit of analysis is an individual customer.

* Significant at 10%; ** significant at 5%; *** significant at 1%.

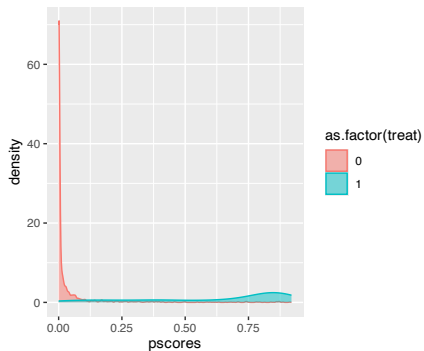
Source: Table 3 of Yoo et al. "Mobile Payment and In-Store Mobile Purchase Behavior" KAIST Working Paper Series



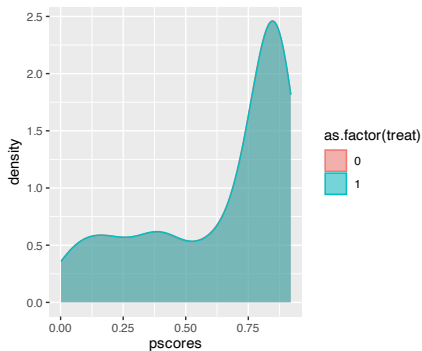
Source: Figure 3 in Ho, Daniel E., et al. "Matching as nonparametric preprossing for reducing model dependence in parametric causal inference." Political analysis 15.3 (2007): 199-236.

Balance Diagnostics: PS Overlap

Propensity Score Distribution Before Matching



Propensity Score Distribution After Matching



Example: LaLonde Dataset

- The effectiveness of a job training program (National Supported Work Demonstration; NSW) on wage increases.
 - The federal government instituted a randomized evaluation of this program.
 - How well the result may be recovered when the experimental controls are replaced with a set of observational controls (Population Survey of Income Dynamics; PSID)?
 - Data publicly available at [NBER data archive](#).
- **Problem:** Imbalances between the experimental and observational data \rightsquigarrow use matching.

Example: LaLonde dataset

- **Data:**
 - Treated: 297 units from NSW
 - Control: 2490 units from PSID
 - Treatment: Participation in the job training program (`treat`)
 - Outcome: 1978 earnings (in dollars; `re78`)
 - Pre-treatment covariates: age, race, marriage, past earnings, past employment

Example R Codes: Data Import

- Import and process data:

```
1 > pacman::p_load(tidyverse, broom, cobalt, Matching, MatchIt)
2
3 > lalonde_nsw <- haven::read_dta(url("http://www.nber.org/~rdehejia/data/nsw.dta"))
4 > PSID_obs <- haven::read_dta(url("http://www.nber.org/~rdehejia/data/psid_controls.dta"))
5
6 > lalonde_ECI <- full_join(lalonde_nsw |>
7   filter(treat == 1),
8   PSID_obs); lalonde_ECI
9
10 # A tibble: 2,787 x 11
11   data_id   treat   age education black hispanic married nodegree re75 re78
12   <chr>     <dbl> <dbl>   <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl>
13 1 Lalonde S...     1    37      11     1     0     1     1     0 9930.
14 2 Lalonde S...     1    22     9     0     1     0     1     0 3596.
15 3 Lalonde S...     1    30     12     1     0     0     0     0 24909.
16 # ... 2,784 more rows
```

Example R Codes: Balance Before Matching

- Assessing balance before matching:

```
1 > bal.tab(x = lalonde_ECI |> dplyr::select(age:re78),
2         treat = lalonde_ECI$treat, continuous = "std", binary = "std")
3
4 Note: `s.d.denom` not specified; assuming "pooled".
5 Balance Measures
6           Type Diff.Un
7 age      Contin. -1.1662
8 education Contin. -0.6862
9 black    Binary  1.3222
10 hispanic Binary  0.2554
11 married  Binary -1.9513
12 nodegree Binary  0.9409
13 re75     Contin. -1.5662
14 re78     Contin. -1.2939
15
16 Sample sizes
17   Control Treated
18 All    2490    297
```

Example R Codes: Propensity Score Matching

- Estimate propensity score (logit) and match 1-to-1 NN:

```
1 # Propensity score
2 > pcores <- glm(treat ~ age + I(age^2) + black + married + hispanic + re75,
3               family = binomial(), data = lalonde_ECI)$fitted.values
4
5 # Conduct one-to-one nearest neighbor propensity score matching
6 > require(Matching)
7 > match_ps <- Match(Y=lalonde_ECI$re78, Tr=lalonde_ECI$treat,
8                   X=pcores, M=1, replace = TRUE)
9 > c(est = as.numeric(match_ps$est), se = match_ps$se)
10      est      se
11 -2030.895  1300.975
```

LaLonde: ATT Estimates Compared

- Three views of the same NSW-vs-PSID effect on 1978 earnings:

| Method | ATT (\$) | SE (\$) |
|---|----------|---------|
| Experimental benchmark (NSW only, $n_T = 297$, $n_C = 425$) | +886 | 488 |
| Naive observational (PSID controls, no matching) | -15,578 | 509 |
| PSM, 1-to-1 NN with replacement | -2,031 | 1,301 |

- Naive obs flips the *sign*: -\$16K vs. +\$886 in the RCT.
- PSM closes $\approx 87\%$ of the bias gap, but not all of it.
- Residual gap: selection on *unmeasured* confounders (motivation, networks) + finite-sample matching is imperfect on observables.



On to the Presentations & Discussions!

Contact Information:

jaewon.yoo@iss.nthu.edu.tw

<https://j1yoo.github.io/>



Appendix

Bias of Inexact Matching

- To show the bias on matching, focus on finding a single control match.
- Let $j(i)$ be the matched control for unit i , the bias is:

$$\mathbb{E}[Y_j | D_i = 1, \mathbf{X}_i, \mathbf{X}_j] - \mathbb{E}[Y_i(0) | D_i = 1, \mathbf{X}_i] = \underbrace{(\mu_0(\mathbf{X}_i) - \mu_0(\mathbf{X}_{j(i)}))}_{\text{unit-level bias}}$$

- Bias is 0 if matching is exact since $\mathbf{X}_i = \mathbf{X}_{j(i)}$
- Bias grows with **matching discrepancy**/imbalance.
- **Bias correction:** estimate $\widehat{\mu}_0(\mathbf{x})$ with regression and estimate bias.

$$\widehat{Y}_i(0) = Y_{j(i)} - (\widehat{\mu}_0(\mathbf{X}_i) - \widehat{\mu}_0(\mathbf{X}_{j(i)}))$$

- Imputation of missing potential outcome now matching + regression.
- Generalizes easily to any number of matches.

Sampling Variance of Matching Estimators

- Matching with replacement: cluster on the match.
 - Can either use clustered SEs or cluster bootstrap.
 - Valid for post-matching regression (Abadie and Spiess, 2021).
- Matching without replacement: more complicated.
 - Same control unit matched to multiple treated: no easy clustering.
 - $K_M(i)$ is the number of times a unit is used as a match.
- Assuming units are well-matched so bias can be ignored,

$$\mathbb{V}(\widehat{\tau}_m) = \frac{1}{n_1} \left(\underbrace{\mathbb{E} [(\tau(\mathbf{X}_i) - \tau_{\text{ATT}})^2 \mid D_i = 1]}_{\text{variance of CATE on treated}} + \underbrace{\mathbb{V}[\widehat{\tau}_m \mid \mathbb{X}, \mathbf{D}]}_{\text{conditional variance}} \right)$$

- Abadie and Imbens (2006) provide matching-based variance estimators.