

Matching

Statistical Modeling & Causal Inference

Agenda

- Matching logic in experimental context
- Matching (exact & propensity score)
- Common support

- (Balance) Tables in R
- Matching in R

Experiment analogy

- **Conditional randomization:** choose relevant covariates; random treatment assignment within (combinations of) covariate levels (a.k.a. randomized block design)
- **Paired randomization:** as above; only two subjects per (combination of) covariate value (but multiple covariate-identical pairs allowed), one of which is randomly assigned to the treatment.

Randomized block design

- Subjects are assigned to blocks, based on gender
- Within each block, subjects are randomly assigned to treatments (placebo or vaccine)
- It is thought that men and women may react differently to this medication
- This design ensures that each treatment condition has an equal proportion of men and women
- As a result, differences between treatment conditions cannot be attributed to gender

Gender	Treatment	
	Placebo	Vaccine
Male	250	250
Female	250	250

Paired randomization design

- Subjects are grouped into pairs based on some blocking variable(s)
- Within each pair, subjects are randomly assigned to different treatments
- Below, 1000 subjects are grouped into 500 matched pairs
- Each pair is matched on gender and age

Pair	Treatment	
	<i>Placebo</i>	<i>Vaccine</i>
1	1	1
2	1	1
...
500	1	1

Matching

- Through matching techniques, we seek to **explicitly balance** the distribution of covariates between treatment and control groups.
- To overcome the lack of ‘twins’ to compare treated and controlled units, we can match observations to the **most plausible** counterfactual available.
- There are multiple ways to define what “most plausible” means. We must choose a technique for that purpose:
 - Mahalanobis/nearest neighbor covariate matching
 - Propensity score matching
 - Coarsened exact matching

Exact Matching

1. Use theoretical and empirical knowledge to identify **relevant confounder(s) (X)**
2. Starting from treated subjects, **select at least one match** from the control group with the same value(s) on X
3. **Drop subjects** off “common support” (unmatched subjects)
4. Estimate causal effect as the average difference in Y **across pairs of matched subjects.**

$$\hat{\delta}_{ATT} = \frac{1}{N_T} \sum_{D_i=1} (Y_i - Y_{j(i)}) = \frac{1}{N_T} \sum_{D_i=1} (Y_i - \underbrace{\left[\frac{1}{M} \sum_{D_i=1} Y_{j_m(1)} \right]}_{\text{average outcome of matches}})$$

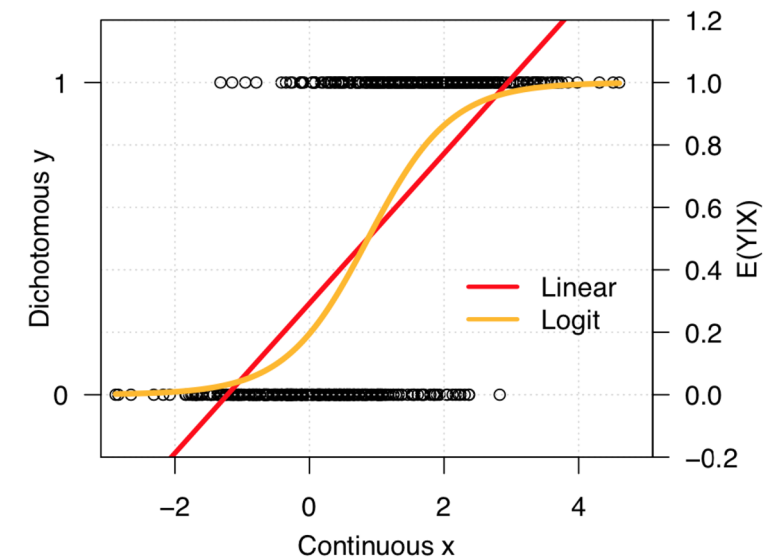
If there is more than one match, you can use their average outcome as the counterfactual.

Propensity score matching

- A propensity score is a measure of the **predicted probability of being in the treatment group**, given the relevant covariates (W).
- We can use propensity scores in order to match treated units with control observations that look **as if** they were treated.
- This is usually modeled with logit/probit regression by which all the potential confounders are used to estimate the single value (PS).

Logit/probit regression

- Similar to linear regression, except we're working with a binary categorical outcome variable.
- Instead of fitting a line to the data, it fits an S-shaped curve that goes from 0 to 1. It tells you the probability of outcome based on the covariates –this is our propensity score!



Propensity score matching

- If the model for estimating the propensity score is well specified (ie. if we chose the right covariates to fulfill back-door criterion), we can control for (match on) the propensity scores and achieve conditional independence.

$$Y_0, Y_1 \perp\!\!\!\perp D \mid pr(W)$$

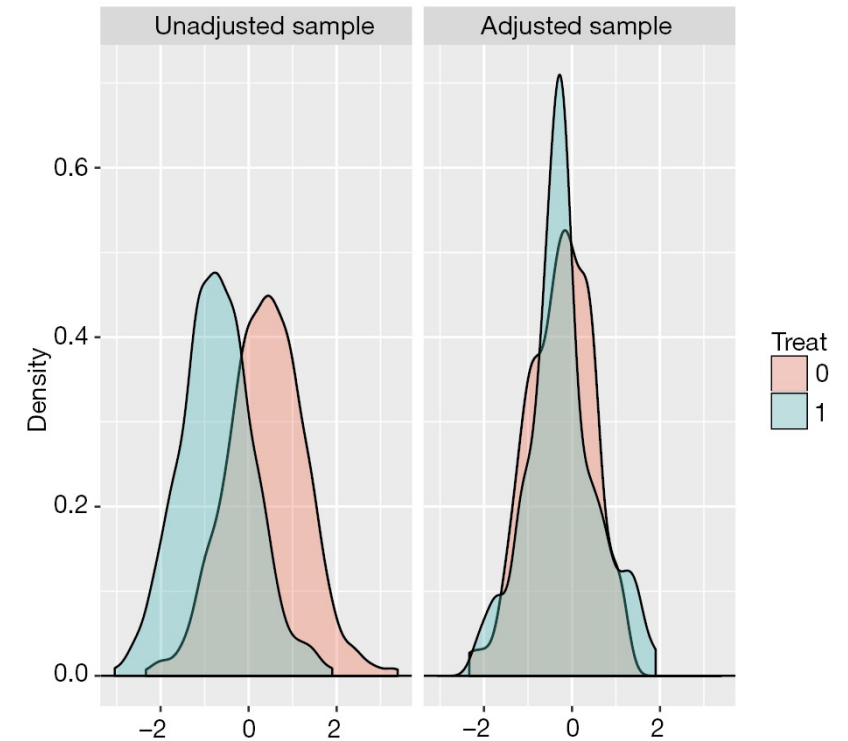
- When there are no exact matches on PS, we can define an algorithm to find the most plausible counterfactual based on PS
→ implies defining issues like replacement, caliper/trimming.

Steps for matching using propensity scores

- Define the set of **potential confounders** (W) by laying out the causal graph.
- Model the probability of $P(D = 1|X)$, using a **logit/probit regression model**.
- Use predicted treatment probabilities as an estimate of **propensity scores**.
- Inspect PS distribution to define whether to **trim** or not. (discard observations unlikely to have a plausible match. Renounce ATE).
- **Match** subjects from treatment and control group applying an algorithm of your choice.
- Check whether your treatment ($D=1$) and control ($D=0$) groups are balanced in terms of the covariates you defined (t-tests). If not balanced, repeat.
- Only then estimate treatment effect (the matching method in itself does not estimate the effect).

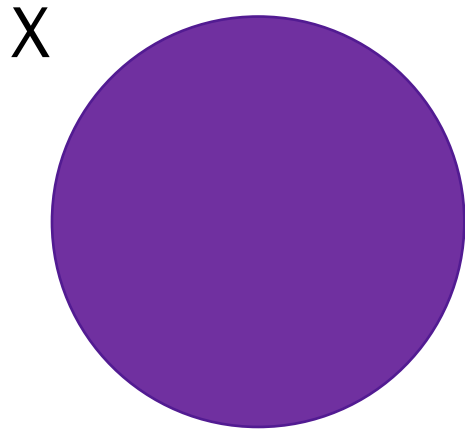
Common support (the curse of dimensionality)

- The more confounders we consider, the less likely it is that we find units with otherwise identical characteristics in the treatment and control groups.
- We **cannot compare all units to a ‘twin’**: they lack common support.
- Without common support for all units, **we cannot estimate the ATE**.
- Knowing **which information** is missing is important. Depending on where the gaps are, we can estimate other effects.



Zhang, Z., Kim, H. J., Lonjon, G., & Zhu, Y. (2019). Balance diagnostics after propensity score matching. *Annals of translational medicine*, 7(1).

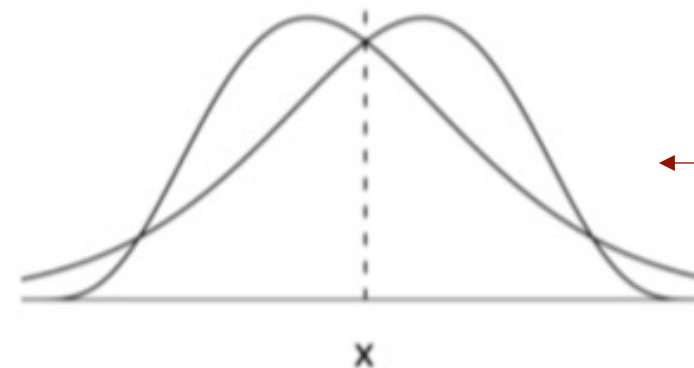
Common support - ATE



 $D = 1$
 $D = 0$

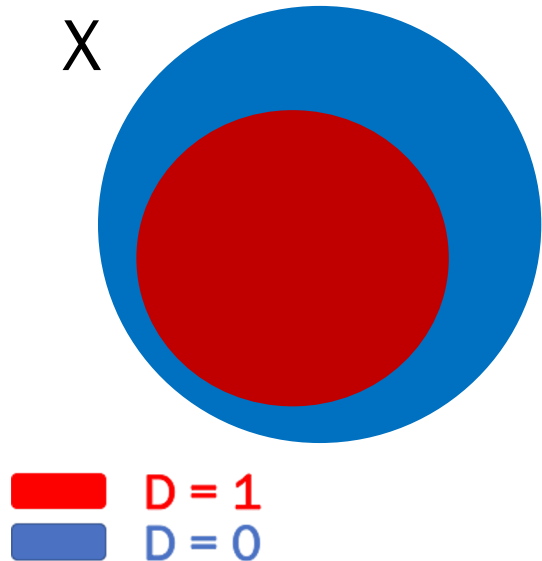
Name	Y	D	X
Jake	10	1	3
Gina	8	1	2
Terry	6	1	1
Rosa	8	0	3
Charles	6	0	2
Ray	4	0	1

In this case, we have full common support, meaning that the distributions of X under both treatment and control are equal. We could gather the ATE.



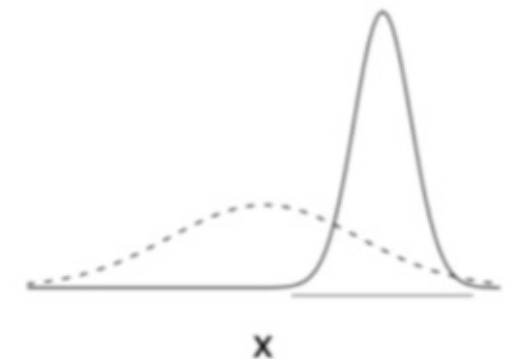
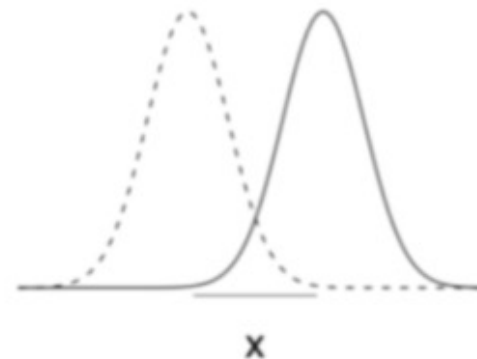
← Not equal but large overlap (with imbalance)

Partial common support - ATT



Name	Y	D	X
Jake	10	1	3
Gina	12	1	3
Terry	8	1	2
Rosa	6	0	3
Charles	3	0	2
Ray	1	0	1

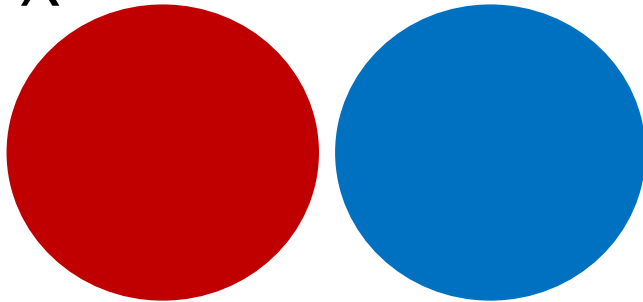
In this case, all our treated units have common support in the control group. But not our controls have a “twin” in the treatment group. We can gather the Average Treatment Effect for the Treated.



No common support



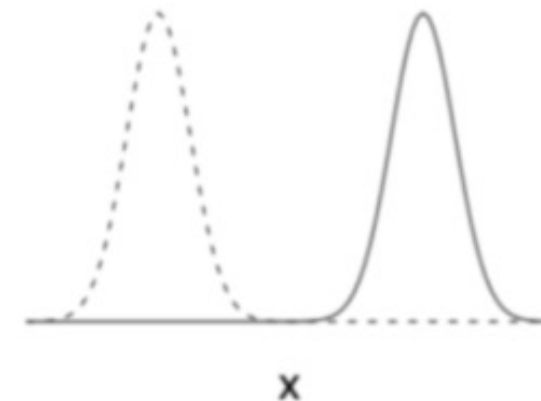
X



■ D = 1
■ D = 0

Name	Y	D	X
Jake	15	1	6
Gina	10	1	5
Terry	5	1	4
Rosa	10	0	1
Charles	6	0	2
Ray	4	0	3

In this case, none of our control and treated units have common support. Our units are non-comparable in their levels of X.



Further Resources

For any coding issues – [Stackoverflow](#)

Hertie's Data Science Lab – [Research Consulting](#)