

# Adaptive Hyper-box Matching for Interpretable Individualized Treatment Effect Estimation

Marco Morucci\*, Vittorio Orlandi\*,  
Sudeepa Roy, Cynthia Rudin, Alexander Volfovsky

Duke | ALMOST MATCHING  
EXACTLY LAB



# Adaptive Hyper-Box (AHB) Matching

Causal inference:  $Y_i(1) - Y_i(0)$

Confounding in observational studies can bias results

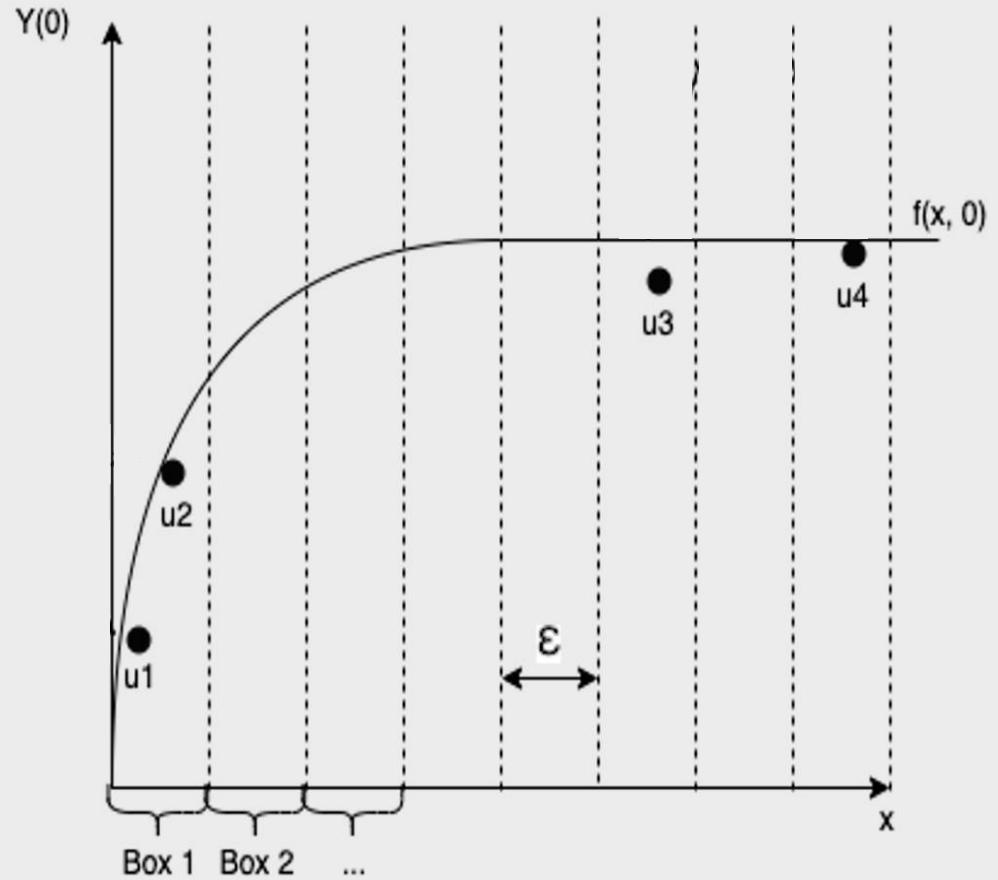
Matching similar units reduces bias

We match units to others falling within an axis-aligned hyper-box

The resulting matches are **case-based** and **interpretable**

# Why Should Boxes Be Adaptive?

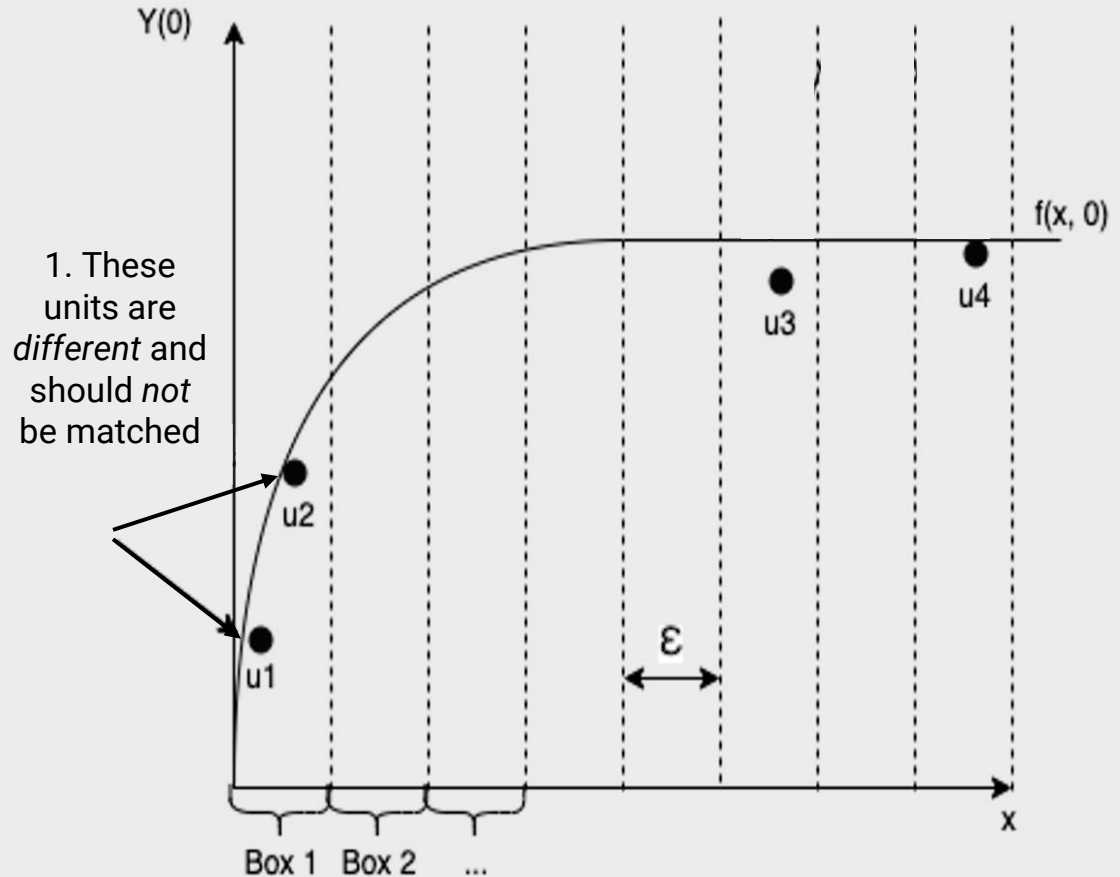
Two examples:



# Why Should Boxes Be Adaptive?

Two examples:

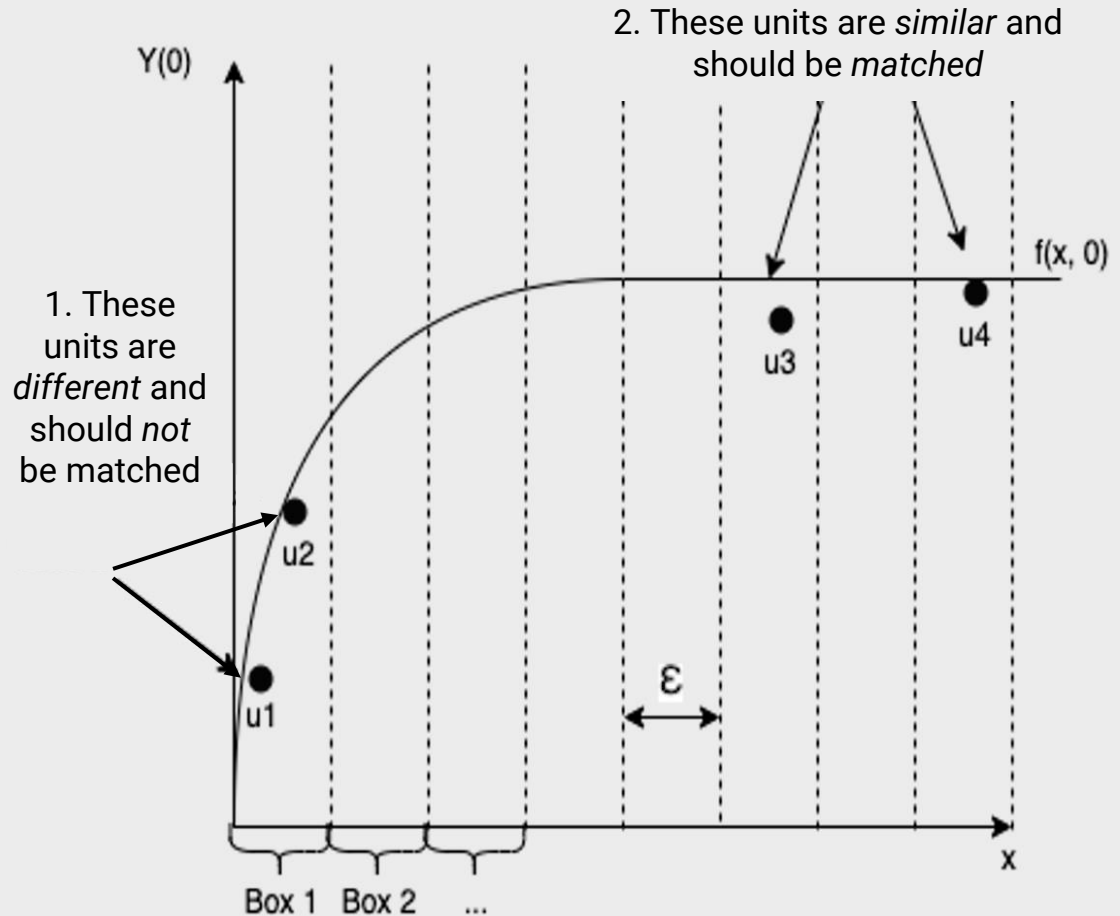
1. Should *not* be matched, but are → **biased** estimates.



# Why Should Boxes Be Adaptive?

Two examples:

1. Should *not* be matched, but are → **biased** estimates.
2. Should be matched, but aren't → **variable** estimates.



# Finding Boxes

Good boxes should have:

1. Low variability



Treated unit to  
match



Smaller potential  
outcome

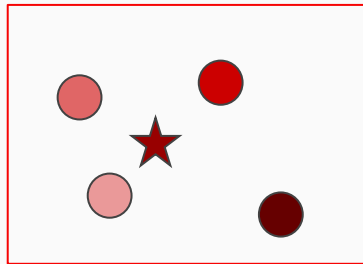


Larger potential  
outcome

# Finding Boxes

Good boxes should have:

1. Low variability



Bad Box



Treated unit to  
match



Smaller potential  
outcome

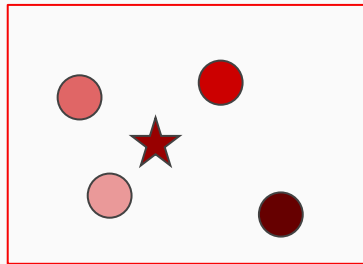


Larger potential  
outcome

# Finding Boxes

Good boxes should have:

1. Low variability
2. Low error



Bad Box



Treated unit to  
match



Smaller potential  
outcome



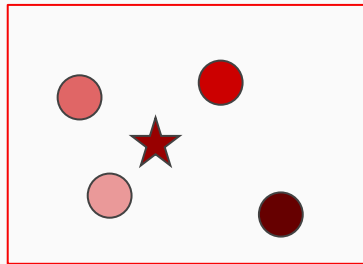
Larger potential  
outcome



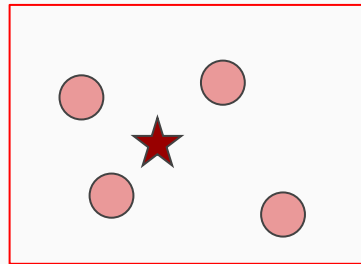
# Finding Boxes

Good boxes should have:

1. Low variability
2. Low error



Bad Box



Bad Box



Treated unit to  
match



Smaller potential  
outcome

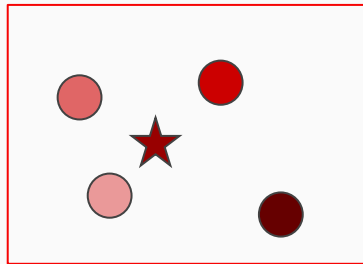


Larger potential  
outcome

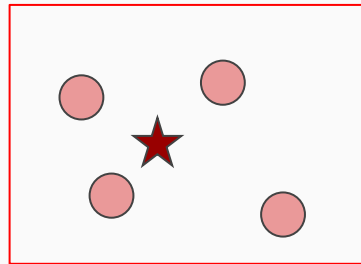
# Finding Boxes

Good boxes should have:

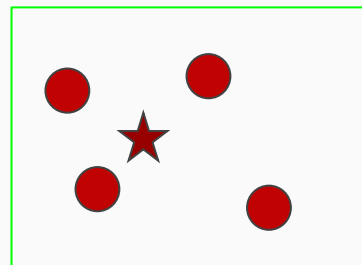
1. Low variability
2. Low error



Bad Box



Bad Box



Good Box



Treated unit to  
match



Smaller potential  
outcome



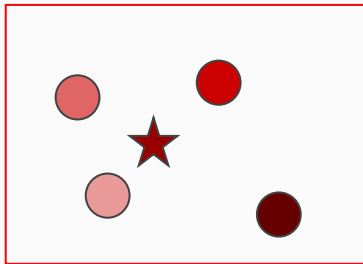
Larger potential  
outcome

# Finding Boxes

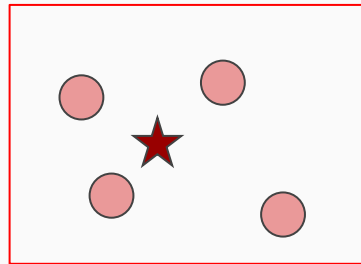
Good boxes should have:

1. Low variability
2. Low error

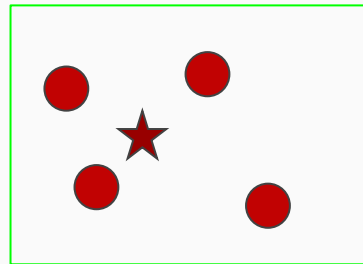
Black Box + Separate Training Set →  
**Honest** learning about outcome  
variability



Bad Box



Bad Box



Good Box



Treated unit to  
match



Smaller potential  
outcome



Larger potential  
outcome

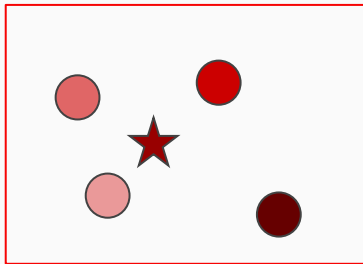
# Finding Boxes

Good boxes should have:

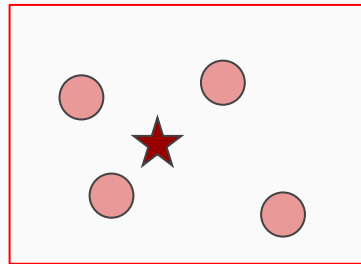
1. Low variability
2. Low error

Black Box + Separate Training Set →  
**Honest** learning about outcome  
variability

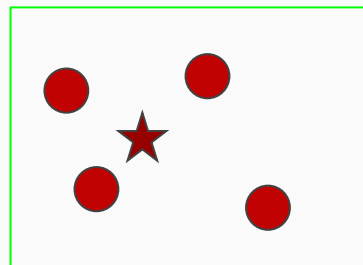
Optimization solved by MIP solver or  
fast approximation



Bad Box



Bad Box



Good Box



Treated unit to  
match



Smaller potential  
outcome



Larger potential  
outcome

# Simulations

$Y = \text{confounder} + (\text{TE modifier}) * \text{treatment} + \text{error}$

Various forms of **confounder**, **treatment effect modifier**:

- Constant
- Linear
- Quadratic
- Constant within Box

Goal: estimate ITE of treated units

# Simulations

$Y = \text{confounder} + (\text{TE modifier}) * \text{treatment} + \text{error}$

Various forms of **confounder**, **treatment effect modifier**

- Constant
- Linear
- Quadratic
- Constant within Box

Goal: estimate ITE of treated units

Compare to:

- BART
- Propensity Score Matching
- Prognostic Score Matching
- Genetic Matching
- Coarsened Exact Matching
- Full Matching
- Prognostic Score Matching w/ True Counterfactuals
- Mahalanobis Distance Matching

# Results

Estimation of ITE of treated units: mean absolute error normalized by ATT

$p$	Method $g / h$	AHB		Black Box	Benchmark	Matching					
		MIP	Fast	BART	Best CF	CEM	Full Matching	GenMatch	Mahal	Nearest Neighbor	Prognostic
(0, 0, 2)	None / Const	0.09	0.05	<b>0.04</b>	0.25	1.01	0.32	0.36	0.34	0.37	0.25
(2, 0, 0)	Box / Const	<b>0.11</b>	0.16	0.24	0.24	0.24	3.03	0.66	0.62	3.05	0.29
(2, 0, 0)	Linear / Const	0.17	0.22	<b>0.14</b>	0.26	0.23	0.82	0.38	0.36	0.91	0.28
(2, 0, 0)	Quad / Const	0.10	<b>0.04</b>	0.08	0.25	0.22	0.42	0.38	0.37	0.45	0.27
(2, 0, 4)	Quad / Const	<b>0.02</b>	<b>0.02</b>	<b>0.02</b>	0.16	NA	0.21	0.12	0.11	0.24	0.04
(1, 1, 0)	Box / Box	<b>0.30</b>	0.45	0.65	0.73	0.58	2.59	2.37	1.02	2.30	0.94
(1, 0, 1)	Binary / Const	<b>0.02</b>	<b>0.02</b>	<b>0.02</b>	0.09	<b>0.02</b>	0.12	0.49	0.10	0.10	0.09
(1, 1, 6)	Binary / Binary	<b>0.06</b>	<b>0.06</b>	0.09	0.17	0.20	0.71	0.97	0.27	0.61	0.18
(2, 0, 0)	Mixed / Const	0.07	0.12	<b>0.06</b>	0.09	0.12	0.48	0.15	0.15	0.55	0.10

# Summary

AHB constructs boxes to match units that are:

1. *Large* where the treatment effect is **near-constant**
2. *Small* where the treatment effect is **highly variable**

The resulting matches are **case-based** and **interpretable**



# Duke

ALMOST MATCHING  
EXACTLY LAB

Thank You!



`almostmatchingexactly.github.io`  
Marco, Vittorio, Sudeepa, Cynthia, Alex