

Generalized Random Forest

Seminar Data Science for Economics

Madina Kurmangaliyeva

m.kurmangaliyeva@uvt.nl

Spring 2020

Tilburg University

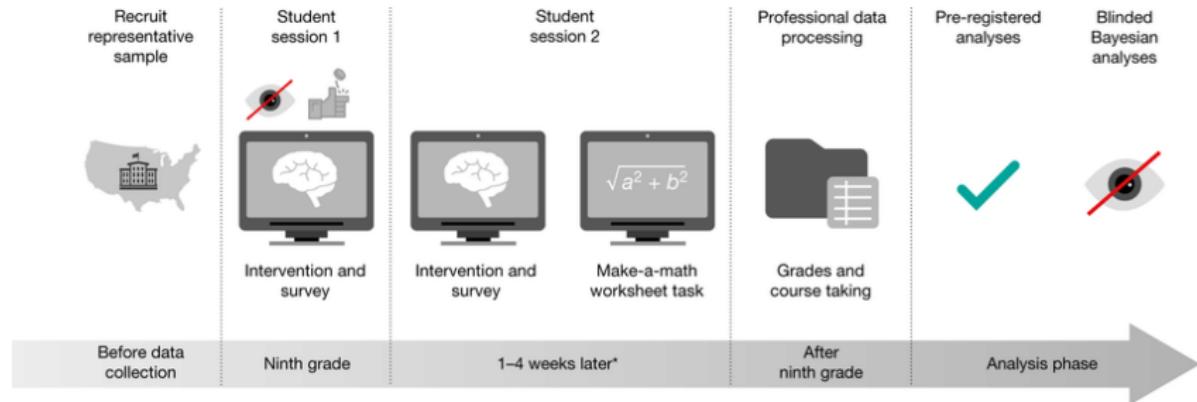
GRF application to the National Study of Learning Mindsets (Athey and Wager, 2019)

National Study of Learning Mindsets

- Large RCT: a short online program designed to foster a growth mindset during the transition to high school
- a low-cost “nudge-like” intervention to change student behavior: “hard work can make you more intelligent” (**learning mindset**)
- as opposed to a **fixed mindset**: “abilities are fixed from birth”
- Primary purpose: seeks to understand heterogeneity in the intervention.

A national experiment reveals where a growth mindset improves achievement

David S. Yeager^{1*}, Paul Hanselman^{2#}, Gregory M. Walton³, Jared S. Murray¹, Robert Crosnoe¹, Chandra Muller¹, Elizabeth Tipton⁴, Barbara Schneider⁵, Chris S. Hulleman⁶, Cintia P. Hinojosa⁷, David Paunesku⁸, Carissa Romero⁹, Kate Flint¹⁰, Alice Roberts¹⁰, Jill Tirott¹⁰, Ronaldo Iachan¹⁰, Jenny Buontempo¹, Sophia Man Yang¹, Carlos M. Carvalho¹, P. Richard Hahn¹¹, Maithreyi Gopalan¹², Pratik Mhatre¹, Ronald Ferguson¹³, Angela L. Duckworth¹⁴ & Carol S. Dweck³



Variables

- Y** a continuous measure of achievement
- W** a binary treatment
- S3** Student: self-reported expectations for success in the future
- C1** Student: race/ethnicity
- C2** Student: gender
- C3** Student: first-generation status
- XC** School: urbanicity (i.e. rural, suburban, etc.)
- X1** School: mean of students' fixed mindsets
- X2** School: achievement level, as measured by test scores and college preparation for the previous 4 cohorts of students
- X3** School: racial/ethnic minority composition
- X4** School: poverty concentration
- X5** School: size - Total # of students

Questions

The main research questions to be addressed include:

1. Was the mindset intervention effective in improving student achievement?
2. Test hypothesis: the effect is moderated by school-level pre-existing mindset norms (X1) and achievement (X2):
 - Effect largest in middle-achieving schools or is decreasing in school-level achievement
3. Exploring possible role of other variables in moderating treatment effects.

Challenges and the GRF

- Multiple hypothesis testing outside a pre-analysis plan (General)
- Observational study: Originally an RCT, but synthetic data led to $\text{corr}(W, S_3) \neq 0$ (General)
- Clustering at school level (GRF-specific)

Generalized Random Forests:

- uses Random Forests machinery
- non-parametrically estimates treatment effects
- captures heterogeneity in treatment effects in a data-driven manner

Results using GRF

Using Generalized Random Forest (Athey, Tibshirani, and Wager, 2019):

- Treatment had a large positive effect on average
- **No** evidence of strong heterogeneity in treatment effects
- **Some** evidence that X_1 – school-level mean of student's fixed mindsets
- **No** evidence that X_2 – school achievement level – mediates the effect

All achieved with just a few lines of code in R using *grf* package.

Plan for today

1. Brief recap of random forests
2. random forests → causal trees → generalized random forests
3. grf package commands in R and the mindset program evaluation

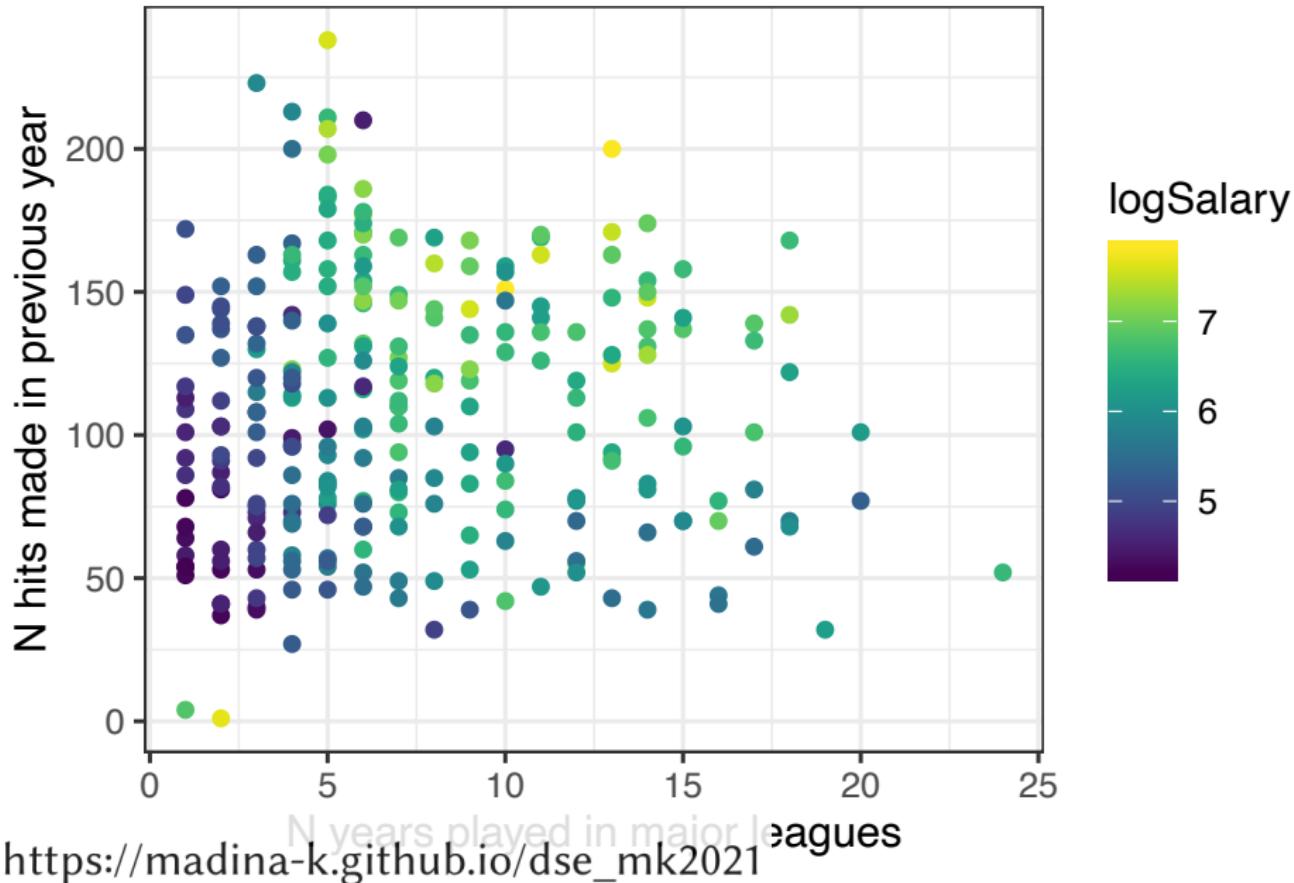
Decision trees (briefly)

Decision trees

Decision trees:

- are very popular for prediction tasks
- help to **partition/segment** the predictor space
- ... for both **regression** (i.e., continuous target variables) and **classification** (categorical target variables)

Baseball players' salary plot: Color indicates log(Salary)

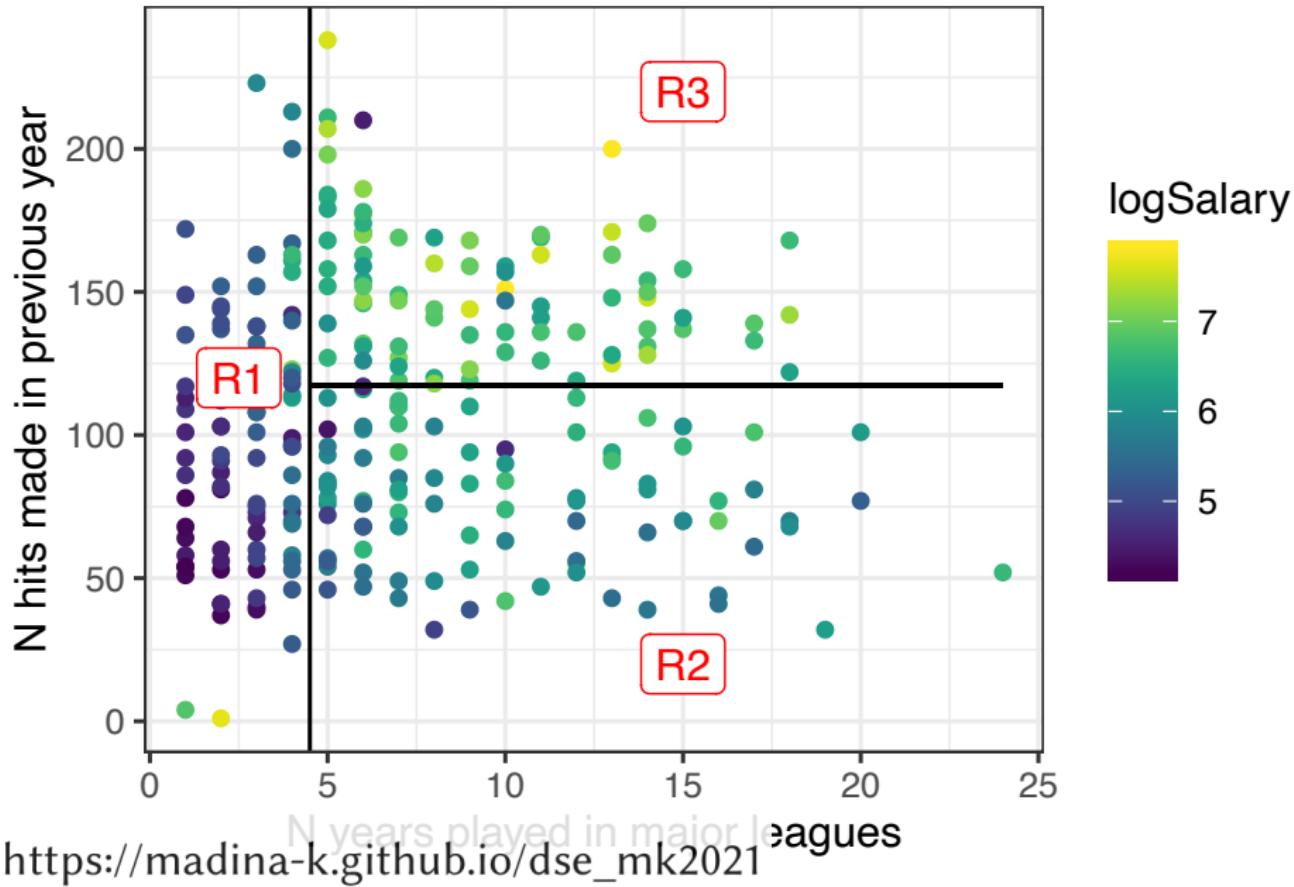


Baseball players' salary: Decision tree

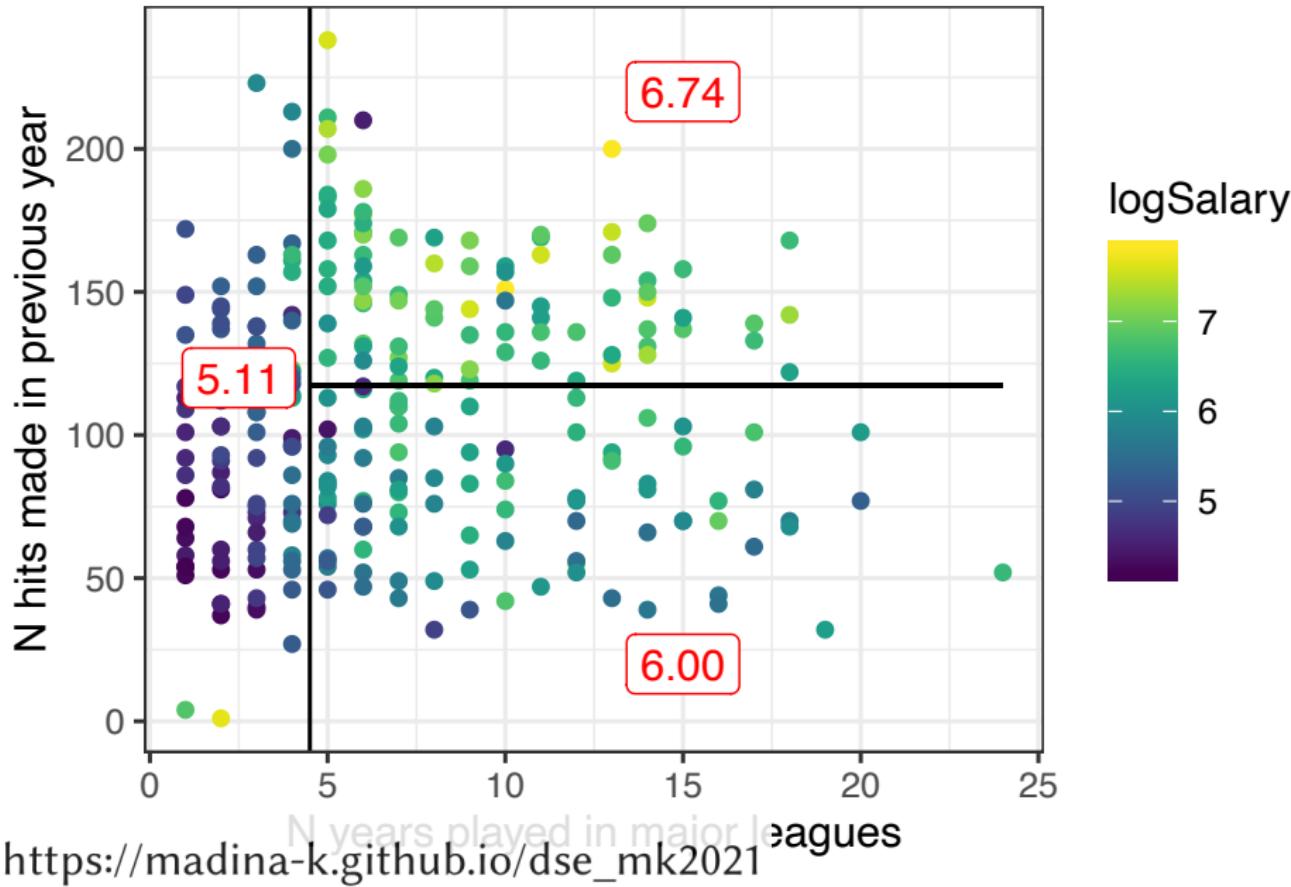
If we train a decision tree to predict the log(Salary) based on the player's characteristics, this is an example of tree that we may get:



Baseball players' salary: regions $R_j, j \in \{1, 2, 3\}$



Baseball players' salary: predictions $\hat{y}_{R_j}, j \in \{1, 2, 3\}$



Prediction

Hence, for any new player we predict his log-salary depending on which region (R_1 , R_2 , or R_3) he belongs to, i.e.:

$$\hat{y}_i(X_i) = \hat{y}_{R_j} \text{ such that } X_i \in R_j, j \in \{1, 2, 3\} \quad (1)$$

How does the tree decide to split data into regions?

The goal is to find the best possible partition of X -space into R_1, \dots, R_J that minimizes the RSS:

$$\sum_{j=1}^J \sum_{i \in R_j} (y_i - \underbrace{\hat{y}_{R_j}}_{\text{prediction for } R_j})^2 \quad (2)$$

We search for the best partition using **recursive binary splitting**.

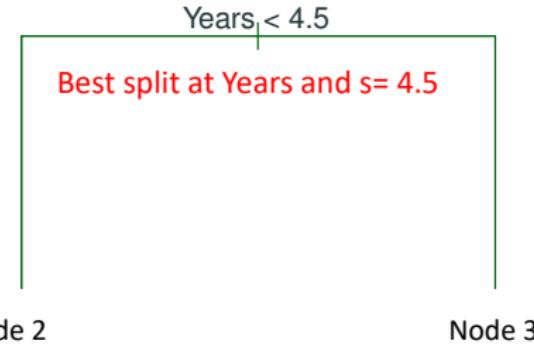
STEP 1

Node 1 (all observations)

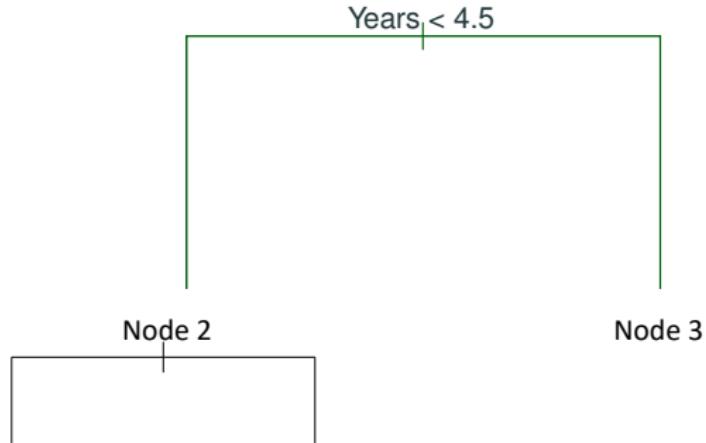


Best split at Years or at Hits? At which cutpoint?

STEP 1

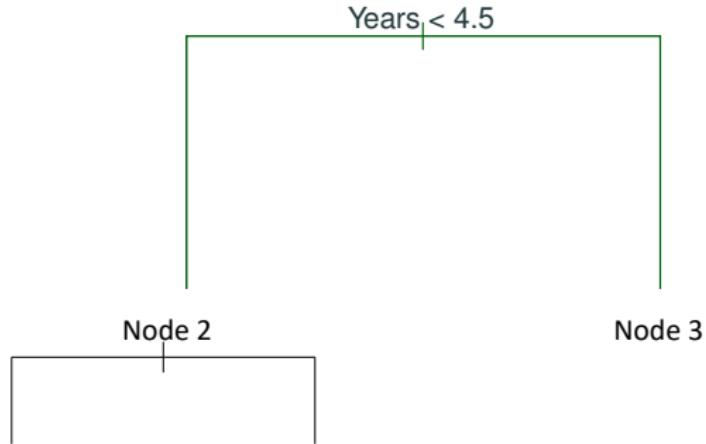


STEP 2



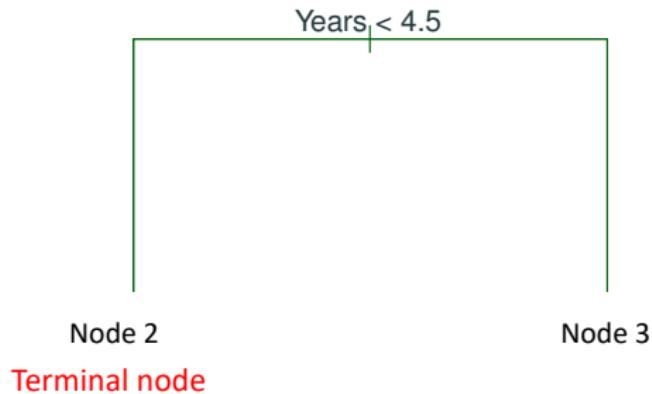
Best split at Years or at Hits? At which cutpoint?

STEP 2

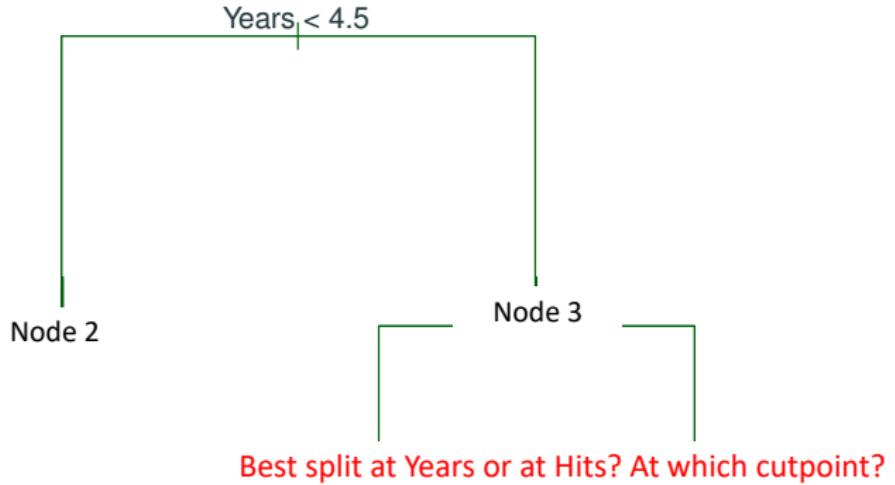


The best split does not pass through the stopping criteria, Node 2 becomes a leaf

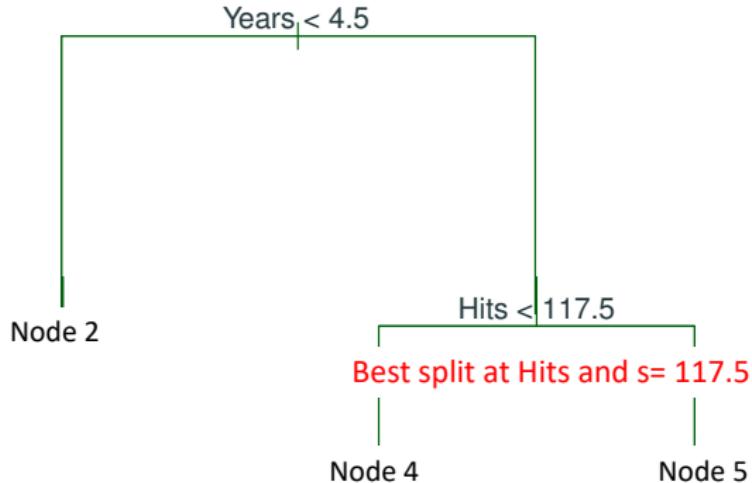
STEP 2

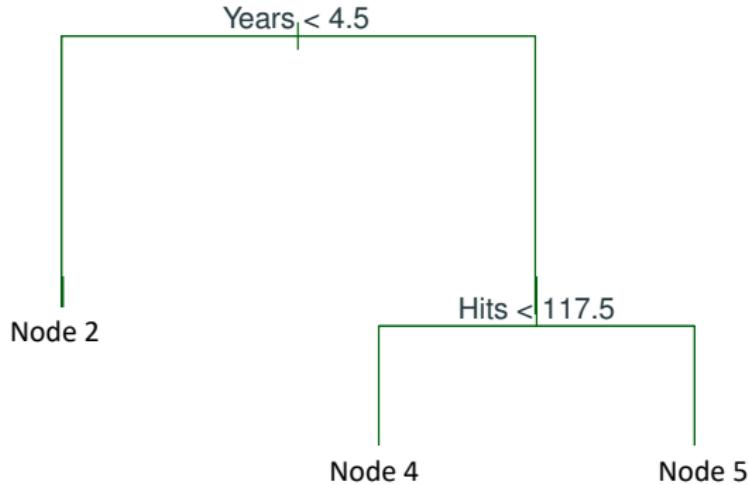


STEP 3



STEP 3





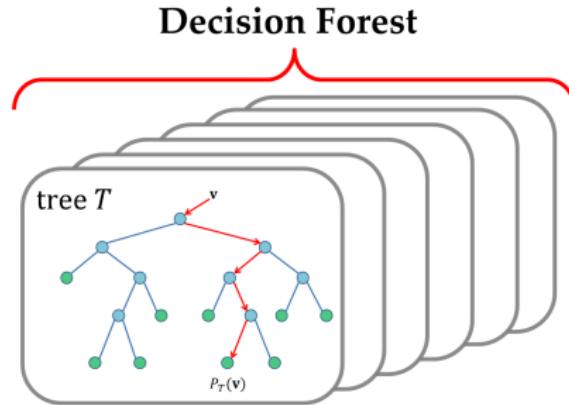
All end nodes are now terminal,
algorithm stops

Calculate predictions

Finally, for each terminal node calculate the mean response value (target variable) for the observations in the training sample which fall into that terminal node.



Random Forest

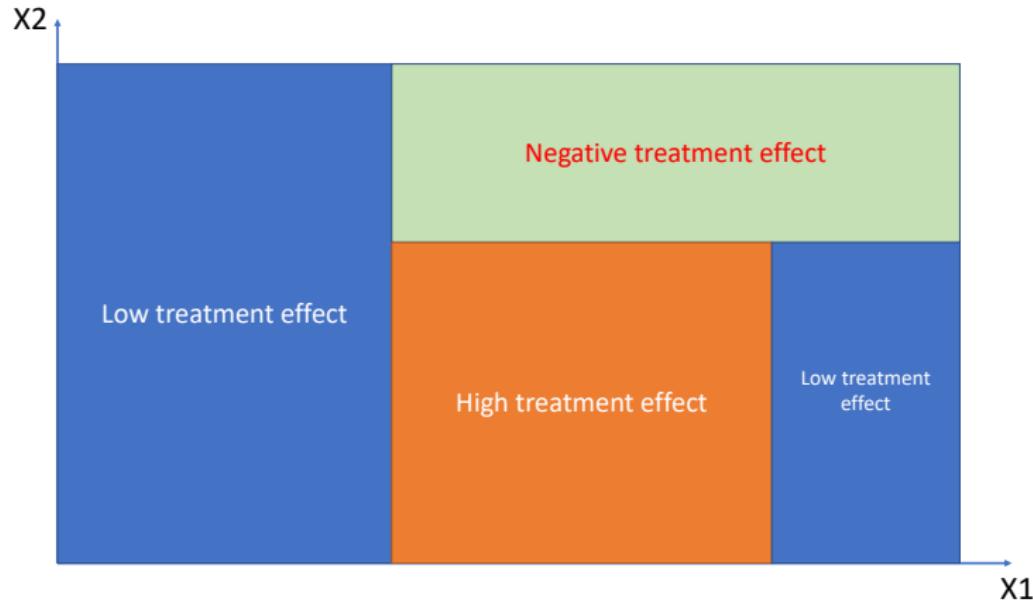


$$\hat{y}^{(RF)} = \underbrace{\sum_{t=1}^T \hat{y}_t}_{\text{average of decision trees}}$$

- Bootstrap: Draw T **random subsamples** of the training data $\Rightarrow T$ trees
- At each split, consider only **some randomly-chosen predictors** as splitting candidates

A random forest for causal inference?

Random forest and Causal inference



Suppose we want to partition covariate space based on the **heterogeneity in treatment effects/CATE** given X. Can we use a random forest?

Cannot use Random Forest for treatment effects

If we had been able to observe the **counterfactuals** directly, ...

| | $Y^{W=1}$ | $Y^{W=0}$ | τ_i | W_i | X_i |
|-----|-------------|-------------|----------|-------|-------|
| 1 | $y_1^{W=1}$ | $y_1^{W=0}$ | τ_1 | 0 | x_1 |
| 2 | $y_2^{W=1}$ | $y_2^{W=0}$ | τ_2 | 1 | x_2 |
| 3 | $y_3^{W=1}$ | $y_3^{W=0}$ | τ_3 | 0 | x_3 |
| ... | ... | ... | ... | ... | ... |
| n | $y_n^{W=1}$ | $y_n^{W=0}$ | τ_n | 1 | x_n |

...then we would have been able to simply plug τ_i as the target variable in a Random forest.

Reality

In reality, the **counterfactuals are missing**

| | $Y^{W=1}$ | $Y^{W=0}$ | τ_i | W_i | X_i |
|-----|-------------|-------------|----------|-------|-------|
| 1 | ? | $y_1^{W=0}$ | ? | 0 | x_1 |
| 2 | $y_2^{W=1}$ | ? | ? | 1 | x_2 |
| 3 | ? | $y_3^{W=0}$ | ? | 0 | x_3 |
| ... | ... | ... | ... | ... | ... |
| n | $y_n^{W=1}$ | ? | ? | 1 | x_n |

So we cannot use RF, we do not observe τ_i !

Naive and wrong approach

Some people wrongly use a RF to simply learn two predictions and take the difference

$$\hat{\tau}^{(RF)}(x) \equiv \underbrace{\hat{y}^{(RF)}(W = 1, X = x) - \hat{y}^{(RF)}(W = 0, X = x)}_{\text{WRONG}}$$

Even under an RCT, $\hat{\tau}^{(RF)}(x)$ estimator will be **biased**

Why? Because **we used the same data** to decide how **to partition** into leaves and **to estimate** mean values within each leaf.

In general

Prediction ≠ Causal inference

What do we want

We want to have a partition of X-space such that:

- captures **heterogeneity** in treatment effects
- **unbiased** estimate
- accounts for the need to estimate **standard errors** (not too narrow partitions)

x2

Captures heterogeneity in TE

$\hat{\tau}_1$

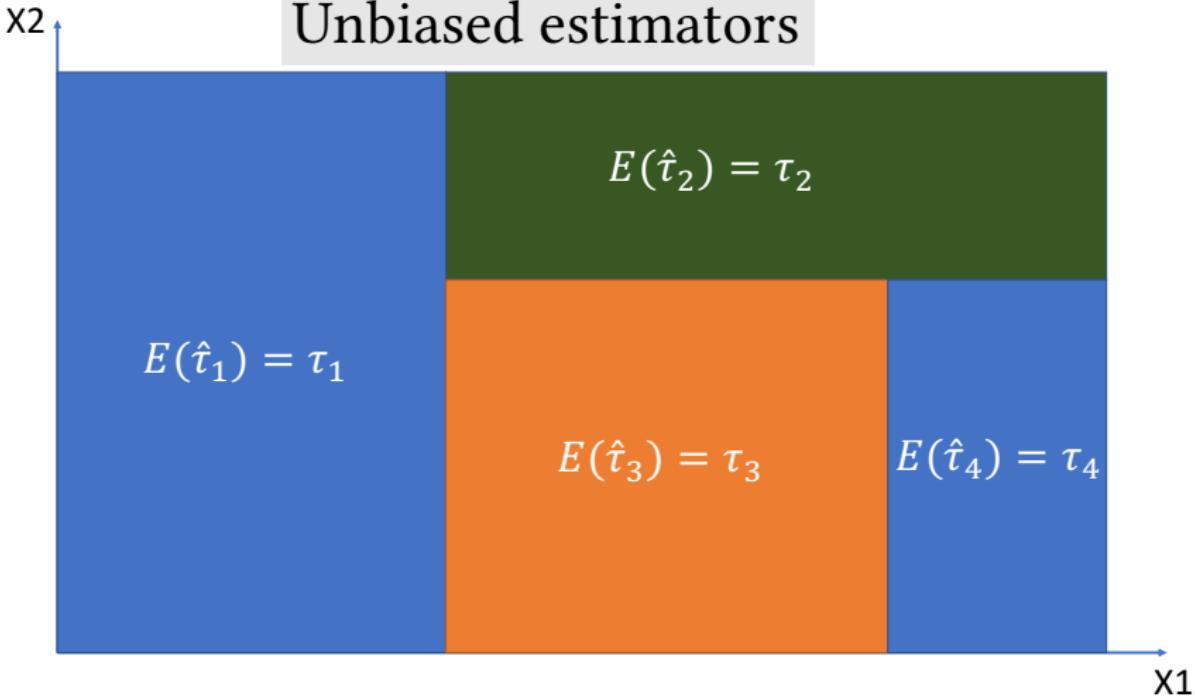
$\hat{\tau}_2$

$\hat{\tau}_3$

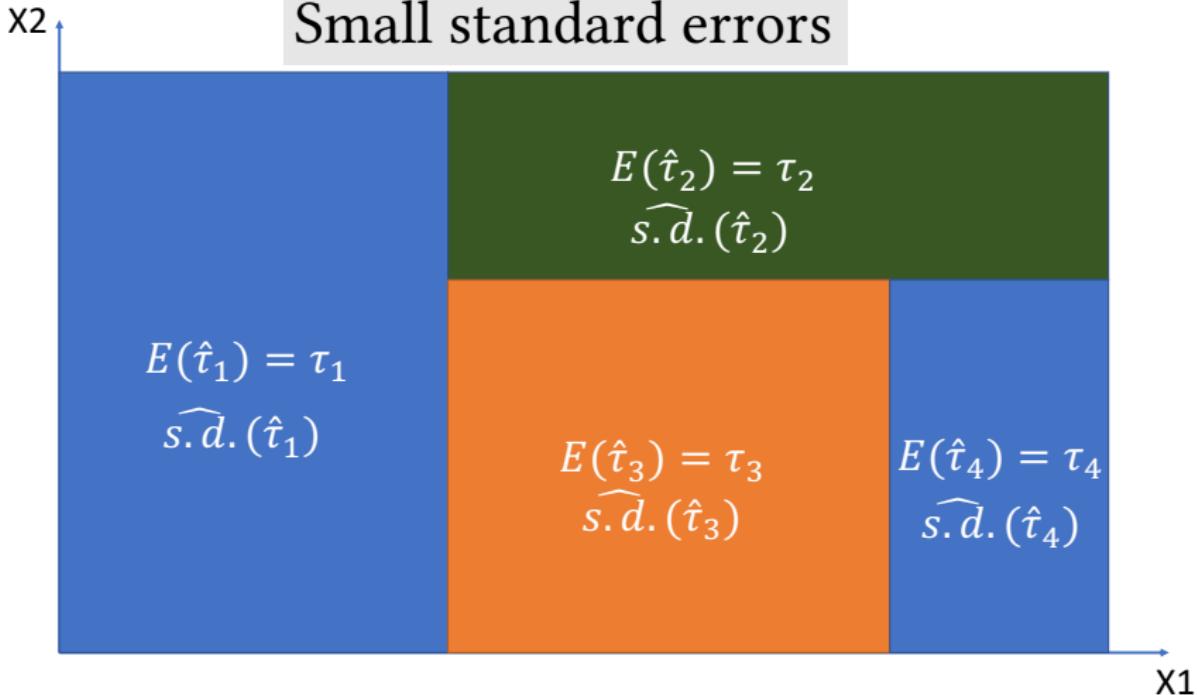
$\hat{\tau}_4$

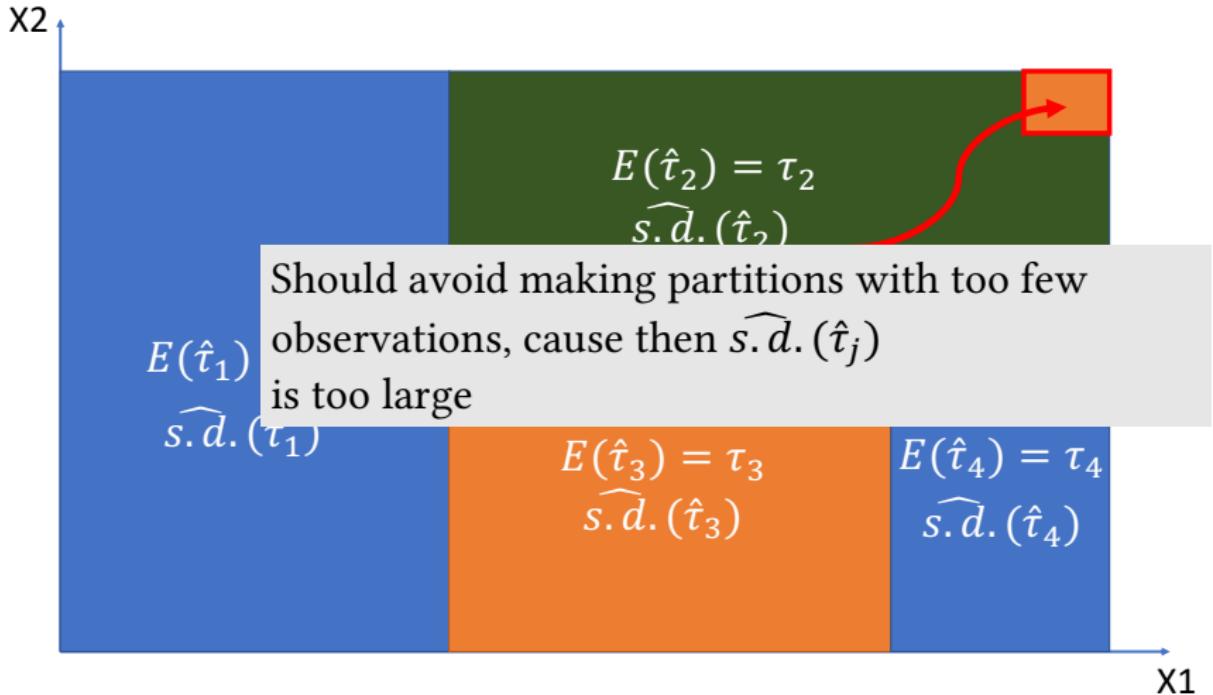
x1

Unbiased estimators



Small standard errors





Causal Trees: RCT, binary treatment

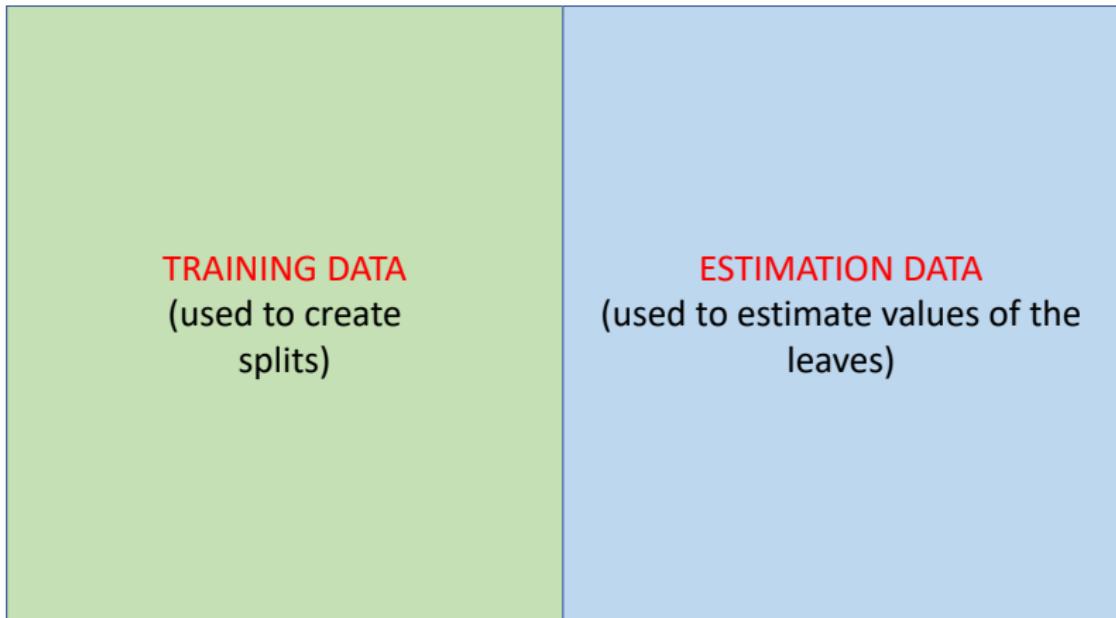
(Athey and Imbens, PNAS 2016)

Preview

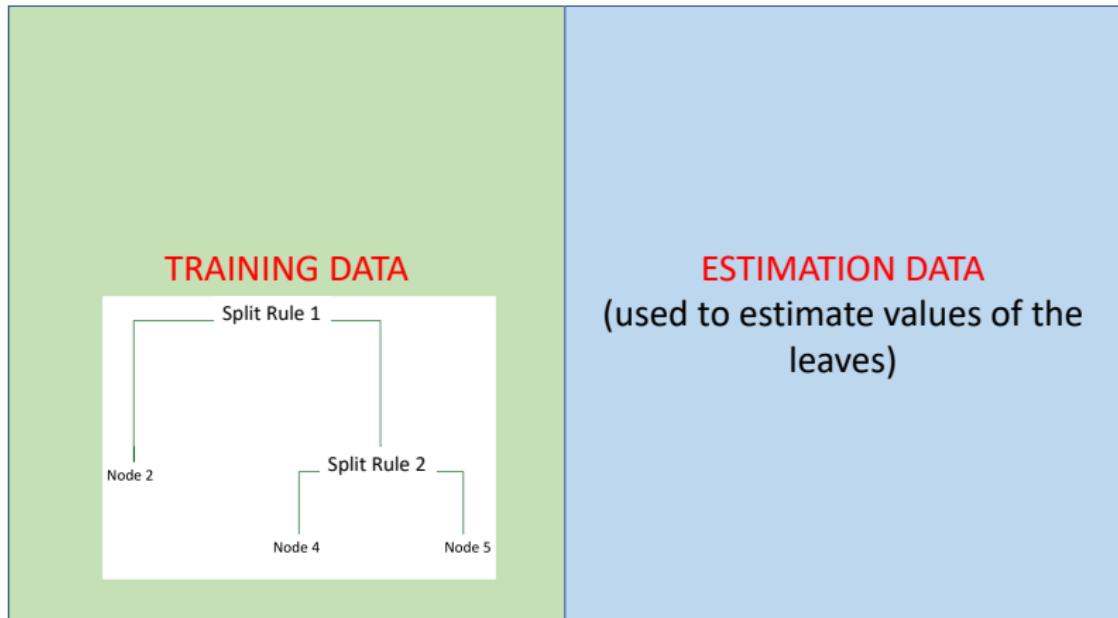
Causal trees solve the bias problem by adopting Honest Splitting:

- uses an **independent sample to estimate** leaf means
- modifies the splitting functions to generate **unbiased estimates** of treatment effects, but
- accounting for the fact that finer leaves:
 - more unbiased treatment effects
 - higher standard errors

Building block 1. Honest splitting



Building block 1. Honest splitting



Building block 1. Honest splitting



Normal decision trees' splitting criteria

$$\min \widehat{MSE}_\mu(\mathcal{S}^{te}, \mathcal{S}^{tr}, \Pi^{tr}) \equiv - \frac{1}{N^{tr}} \sum_{i \in \mathcal{S}^{tr}} \hat{\mu}^2(X_i; \mathcal{S}^{tr}, \Pi^{tr})$$

Normal decision trees' splitting criteria

$$\min \widehat{MSE}_\mu(\mathcal{S}^{te}, \mathcal{S}^{tr}, \Pi^{tr}) \equiv - \frac{1}{N^{tr}} \sum_{i \in \mathcal{S}^{tr}} \hat{\mu}^2(X_i; \mathcal{S}^{tr}, \Pi^{tr})$$


test set

Normal decision trees' splitting criteria

$$\min \widehat{MSE}_\mu(\mathcal{S}^{te}, \mathcal{S}^{tr}, \Pi^{tr}) \equiv - \frac{1}{N^{tr}} \sum_{i \in \mathcal{S}^{tr}} \hat{\mu}^2(X_i; \mathcal{S}^{tr}, \Pi^{tr})$$

The diagram shows two gray rectangular boxes at the bottom. The left box contains the text "test set" and the right box contains the text "training set". Two black arrows originate from these boxes and point upwards towards the \mathcal{S}^{te} and \mathcal{S}^{tr} terms in the mathematical equation above.

Normal decision trees' splitting criteria

$$\min \widehat{MSE}_\mu(\mathcal{S}^{te}, \mathcal{S}^{tr}, \Pi^{tr}) \equiv - \frac{1}{N^{tr}} \sum_{i \in \mathcal{S}^{tr}} \hat{\mu}^2(X_i; \mathcal{S}^{tr}, \Pi^{tr})$$

test set training set tree partitions based on training set

Normal decision trees' splitting criteria

equivalent to maximizing the sum of the final nodes' squared predictions

$$\min \widehat{MSE}_\mu(\mathcal{S}^{te}, \mathcal{S}^{tr}, \Pi^{tr}) \equiv - \frac{1}{N^{tr}} \sum_{i \in \mathcal{S}^{tr}} \hat{\mu}^2(X_i; \mathcal{S}^{tr}, \Pi^{tr})$$

test set

training set

tree partitions based on
training set

Normal decision trees' splitting criteria

equivalent to maximizing the sum of the final nodes' squared predictions

$$\min \widehat{MSE}_\mu(\mathcal{S}^{te}, \mathcal{S}^{tr}, \Pi^{tr}) \equiv - \frac{1}{N^{tr}} \sum_{i \in \mathcal{S}^{tr}} \hat{\mu}^2(X_i; \mathcal{S}^{tr}, \Pi^{tr})$$

test set

training set

tree partitions based on
training set

It needs only the training sample and the number of observations in the training sample.

Building block 2. New objective function

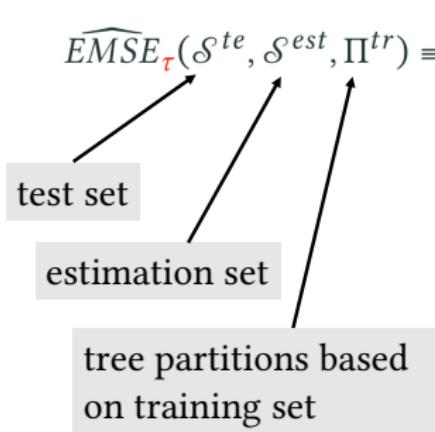
Causal Trees:

$$\widehat{EMSE}_{\tau}(\mathcal{S}^{te}, \mathcal{S}^{est}, \Pi^{tr}) \equiv - \underbrace{\frac{1}{N^{tr}} \sum_{i \in \mathcal{S}^{tr}} \hat{\tau}^2(X_i; \mathcal{S}^{tr}, \Pi^{tr})}_{\text{rewards high heterogeneity}} + \underbrace{\left(\frac{1}{N^{tr}} + \frac{1}{N^{est}} \right) \sum_{l \in \Pi^{tr}} \left(\frac{S_{\mathcal{S}^{tr} \text{treat}}^2(l)}{p} + \frac{S_{\mathcal{S}^{tr} \text{control}}^2(l)}{1-p} \right)}_{\text{Penalizes splits leading to small leafs}}$$

Building block 2. New objective function

Causal Trees:

$$\widehat{EMSE}_{\tau}(\mathcal{S}^{te}, \mathcal{S}^{est}, \Pi^{tr}) \equiv - \frac{1}{N^{tr}} \sum_{i \in \mathcal{S}^{tr}} \hat{\tau}^2(X_i; \mathcal{S}^{tr}, \Pi^{tr})$$

A diagram illustrating the components of the objective function. Three arrows point from labels on the left to specific terms in the equation:

- An arrow from "test set" points to the first term, \widehat{EMSE}_{τ} .
- An arrow from "estimation set" points to the second term, $\sum_{i \in \mathcal{S}^{tr}}$.
- An arrow from "tree partitions based on training set" points to the third term, $\left(\frac{1}{N^{tr}} + \frac{1}{N^{est}} \right) \sum_{l \in \Pi^{tr}}$.

rewards high heterogeneity

$$+ \left(\frac{1}{N^{tr}} + \frac{1}{N^{est}} \right) \sum_{l \in \Pi^{tr}} \left(\frac{S_{\mathcal{S}^{tr}}^2(l)}{p} + \frac{S_{\mathcal{S}^{tr}}^2(l)}{1-p} \right)$$

Penalizes splits leading to small leafs

Building block 2. New objective function

Causal Trees:

$$\widehat{EMSE}_{\tau}(\mathcal{S}^{te}, \mathcal{S}^{est}, \Pi^{tr}) \equiv - \frac{1}{N^{tr}} \sum_{i \in \mathcal{S}^{tr}} \hat{\tau}^2(X_i; \mathcal{S}^{tr}, \Pi^{tr})$$

rewards high heterogeneity

$$+ \left(\frac{1}{N^{tr}} + \frac{1}{N^{est}} \right) \sum_{l \in \Pi^{tr}} \left(\frac{S_{\mathcal{S}^{tr} \text{ treat}}^2(l)}{p} + \frac{S_{\mathcal{S}^{tr} \text{ control}}^2(l)}{1-p} \right)$$

Penalizes splits leading to small leafs

test set estimation set tree partitions based on training set

Building block 2. New objective function

Causal Trees:

$$\widehat{EMSE}_{\tau}(\mathcal{S}^{te}, \mathcal{S}^{est}, \Pi^{tr}) \equiv - \frac{1}{N^{tr}} \sum_{i \in \mathcal{S}^{tr}} \hat{\tau}^2(X_i | \mathcal{S}^{tr}, \Pi^{tr})$$

rewards high heterogeneity

$$+ \left(\frac{1}{N^{tr}} + \frac{1}{N^{est}} \right) \sum_{l \in \Pi^{tr}} \left(\frac{S_{\mathcal{S}^{tr} \text{ treat}}^2(l)}{p} + \frac{S_{\mathcal{S}^{tr} \text{ control}}^2(l)}{1-p} \right)$$

Penalizes splits leading to small leafs

treatment effects

only using training set

test set

estimation set

tree partitions based on training set

Building block 2. New objective function

Causal Trees:

$$\widehat{EMSE}_{\tau}(\mathcal{S}^{te}, \mathcal{S}^{est}, \Pi^{tr}) \equiv - \frac{1}{N^{tr}} \sum_{i \in \mathcal{S}^{tr}} \hat{\tau}^2(X_i | \mathcal{S}^{tr}, \Pi^{tr})$$

rewards high heterogeneity

$$+ \left(\frac{1}{N^{tr}} + \frac{1}{N^{est}} \right) \sum_{l \in \Pi^{tr}} \left(\frac{S_{\mathcal{S}^{tr} \text{ treat}}^2(l)}{p} + \frac{S_{\mathcal{S}^{tr} \text{ control}}^2(l)}{1-p} \right)$$

Penalizes splits leading to small leafs

treatment effects

only using training set

variance ...

test set

estimation set

tree partitions based on training set

Building block 2. New objective function

Causal Trees:

$$\widehat{EMSE}_{\tau}(\mathcal{S}^{te}, \mathcal{S}^{est}, \Pi^{tr}) \equiv - \frac{1}{N^{tr}} \sum_{i \in \mathcal{S}^{tr}} \hat{\tau}^2(X_i | \mathcal{S}^{tr}, \Pi^{tr})$$

rewards high heterogeneity

$$+ \left(\frac{1}{N^{tr}} + \frac{1}{N^{est}} \right) \sum_{l \in \Pi^{tr}} \left(\frac{S_{\mathcal{S}^{tr} \text{ treat}}^2(l)}{p} + \frac{S_{\mathcal{S}^{tr} \text{ control}}^2(l)}{1-p} \right)$$

Penalizes splits leading to small leafs

in the sample of treated observations in leaf (l)

The diagram illustrates the components of the \widehat{EMSE}_{τ} objective function. It starts with the formula:
$$\widehat{EMSE}_{\tau}(\mathcal{S}^{te}, \mathcal{S}^{est}, \Pi^{tr}) \equiv - \frac{1}{N^{tr}} \sum_{i \in \mathcal{S}^{tr}} \hat{\tau}^2(X_i | \mathcal{S}^{tr}, \Pi^{tr})$$
. A red box highlights the term $\hat{\tau}^2(X_i | \mathcal{S}^{tr}, \Pi^{tr})$, which is labeled 'treatment effects'. An arrow from this term points to a box labeled 'only using training set'. Another arrow from the same term points to a box labeled 'variance ...'. Below the first term, a bracket indicates it 'rewards high heterogeneity'. The second part of the formula is:
$$+ \left(\frac{1}{N^{tr}} + \frac{1}{N^{est}} \right) \sum_{l \in \Pi^{tr}} \left(\frac{S_{\mathcal{S}^{tr} \text{ treat}}^2(l)}{p} + \frac{S_{\mathcal{S}^{tr} \text{ control}}^2(l)}{1-p} \right)$$
. This part is labeled 'Penalizes splits leading to small leafs'. An arrow from the second term points to a box labeled 'in the sample of treated observations in leaf (l)'. To the left of the formula, three boxes are connected by arrows pointing to the formula: 'test set' (top), 'estimation set' (middle), and 'tree partitions based on training set' (bottom). The 'test set' and 'estimation set' boxes both have arrows pointing to the first term of the formula.

Building block 2. New objective function

Causal Trees:

$$\widehat{EMSE}_{\tau}(\mathcal{S}^{te}, \mathcal{S}^{est}, \Pi^{tr}) \equiv - \frac{1}{N^{tr}} \sum_{i \in \mathcal{S}^{tr}} \hat{\tau}^2(X_i | \mathcal{S}^{tr}, \Pi^{tr})$$

rewards high heterogeneity

$$+ \left(\frac{1}{N^{tr}} + \frac{1}{N^{est}} \right) \sum_{l \in \Pi^{tr}} \left(\frac{S_{\mathcal{S}^{tr} \text{ treat}}^2(l)}{p} + \frac{S_{\mathcal{S}^{tr} \text{ control}}^2(l)}{1-p} \right)$$

Penalizes splits leading to small leafs

test set

estimation set

tree partitions based
on training set

only using training set

variance ...

in the sample of control observations in leaf (l)

Building block 2. New objective function

Causal Trees:

$$\widehat{EMSE}_{\tau}(\mathcal{S}^{te}, \mathcal{S}^{est}, \Pi^{tr}) \equiv - \frac{1}{N^{tr}} \sum_{i \in \mathcal{S}^{tr}} \hat{\tau}^2(X_i | \mathcal{S}^{tr}, \Pi^{tr})$$

treatment effects

only using training set

rewards high heterogeneity

$$+ \left(\frac{1}{N^{tr}} + \frac{1}{N^{est}} \right) \sum_{l \in \Pi^{tr}} \left(\frac{S_{\mathcal{S}^{tr} \text{ treat}}^2(l)}{p} + \frac{S_{\mathcal{S}^{tr} \text{ control}}^2(l)}{1-p} \right)$$

Penalizes splits leading to small leafs

test set

estimation set

tree partitions based on training set

probability of treatment

Building block 2. New objective function

Causal Trees:

$$\widehat{EMSE}_{\tau}(\mathcal{S}^{te}, \mathcal{S}^{est}, \Pi^{tr}) \equiv - \frac{1}{N^{tr}} \sum_{i \in \mathcal{S}^{tr}} \hat{\tau}^2(X_i | \mathcal{S}^{tr}, \Pi^{tr})$$

underbrace
rewards high heterogeneity

$$+ \left(\frac{1}{N^{tr}} + \frac{1}{N^{est}} \right) \sum_{l \in \Pi^{tr}} \left(\frac{S_{\mathcal{S}^{tr} \text{ treat}}^2(l)}{p} + \frac{S_{\mathcal{S}^{tr} \text{ control}}^2(l)}{1-p} \right)$$

underbrace
Penalizes splits leading to small leafs

treatment effects

only using training set

test set

estimation set

tree partitions based on training set

Number of obs

Building block 2. New objective function

Causal Trees:

$$\widehat{EMSE}_\tau(\mathcal{S}^{te}, \mathcal{S}^{est}, \Pi^{tr}) \equiv - \frac{1}{N^{tr}} \sum_{i \in \mathcal{S}^{tr}} \hat{\tau}^2(X_i | \mathcal{S}^{tr}, \Pi^{tr})$$

treatment effects
rewards high heterogeneity

$$+ \left(\frac{1}{N^{tr}} + \frac{1}{N^{est}} \right) \sum_{l \in \Pi^{tr}} \left(\frac{S_{\mathcal{S}^{tr} \text{ treat}}^2(l)}{p} + \frac{S_{\mathcal{S}^{tr} \text{ control}}^2(l)}{1-p} \right)$$

Penalizes splits leading to small leafs

test set
estimation set
tree partitions based on training set

It needs only the training sample and the number of observations in the training and the estimation sample.

Building block 2. New objective function

Causal Trees:

$$\widehat{EMSE}_{\hat{\tau}}(\mathcal{S}^{te}, \mathcal{S}^{est}, \Pi^{tr}) \equiv - \frac{1}{N^{tr}} \sum_{i \in \mathcal{S}^{tr}} \hat{\tau}^2(X_i; \mathcal{S}^{tr}, \Pi^{tr})$$

$\underbrace{\quad\quad\quad}_{\text{rewards high heterogeneity}}$

$$+ \left(\frac{1}{N^{tr}} + \frac{1}{N^{est}} \right) \sum_{l \in \Pi^{tr}} \left(\frac{S_{\mathcal{S}^{tr} \text{treat}}^2(l)}{p} + \frac{S_{\mathcal{S}^{tr} \text{control}}^2(l)}{1-p} \right)$$

$\underbrace{\quad\quad\quad}_{\text{Penalizes splits leading to small leafs}}$

Estimated treatment effect is the difference in the average y between treated and control observations:

$$\hat{\tau}(x; \Pi^{tr}) \equiv \hat{\mu}(W = 1, x; \mathcal{S}^{tr}, \Pi^{tr}) - \hat{\mu}(W = 0, x; \mathcal{S}^{tr}, \Pi^{tr})$$

Summary: Decision trees vs Causal Trees

| | Regression Tree | Causal Tree |
|------------------------------------|-----------------|----------------------|
| Predictions based on | training sample | estimation sample |
| Splitting rule minimizes in-sample | RSS | Honest Target |
| Segments X for heterogeneity | in outcomes | in treatment effects |

In DML

The last step:

$$Y_i = \beta W_i + \epsilon_i \quad (3)$$

Generalized Random Forest

Random Effects Model:

$$Y_i = b_i W_i + \epsilon_i; \quad \beta(x) = E(b_i | X_i = x); \quad (4)$$

(Note: under non-RCT, we first partial out X_i from Y_i and W_i : 3-steps procedure)

Generalized Random Forest

Random Effects Model:

$$Y_i = b_i W_i + \epsilon_i; \quad \beta(x) = E(b_i | X_i = x); \quad (4)$$

(Note: under non-RCT, we first partial out X_i from Y_i and W_i : 3-steps procedure)

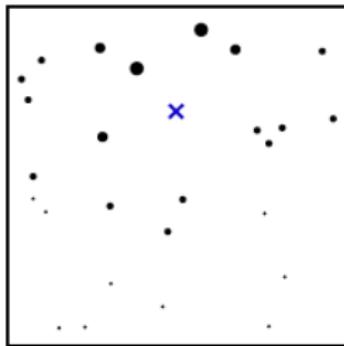
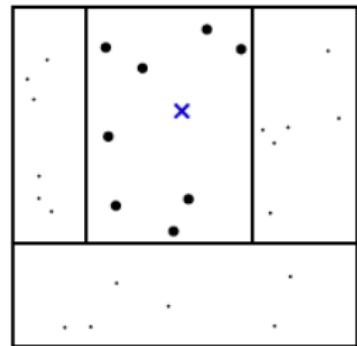
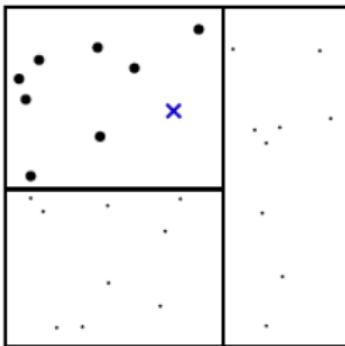
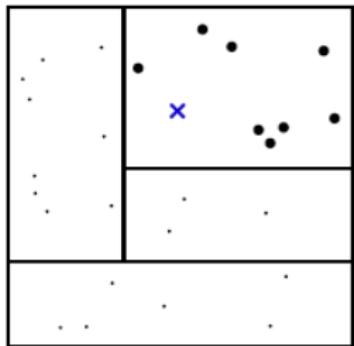
In particular, the generalized random forest estimates $\beta(x)$ non-parametrically:

$$\hat{\beta}(x) = \frac{\sum_{i=1}^n \alpha_i(x)(W_i - \bar{W}_\alpha)(Y_i - \bar{Y}_\alpha)}{\sum_{i=1}^n \alpha_i(x)(W_i - \bar{W}_\alpha)^2} \quad (5)$$

where

- $\alpha_i(x)$ is a weight determined by Causal Forest
- $\bar{W}_\alpha = \sum \alpha_i(x) W_i$ is a weighted average treatment
- $\bar{Y}_\alpha = \sum \alpha_i(x) Y_i$ is a weighted average outcome

Weights from Causal Forest



Weights from Causal Forest

$$\alpha_{bi}(x) = \frac{\mathbf{1}(\{X_i \in L_b(x)\})}{|L_b(x)|}, \quad \alpha_i(x) = \frac{1}{B} \sum_{b=1}^B \alpha_{bi}(x).$$

Back to the National Study of Learning Mindsets (Athey and Wager, 2019)

grf algorithm

```
cf <- causal_forest(X, Y, W, clusters = school.id,  
                      equalize.cluster.weights = TRUE)
```

This code performs the following steps:

1. Predict Y from X using Random Forest, collect out-of-sample \hat{Y}
2. Predict W from X using Random Forest, collect out-of-sample \hat{W}
3. Run a causal forest of:
 - $Y - \hat{Y}$ on $W - \hat{W}$
 - using X as potential mediators
 - and school ids as clusters
 - drawing the same # of obs per cluster

Cluster-robust random/causal forests

To account for clusters, random/causal forests:

- At bootstrapping stage:
 1. First draw a subsample of clusters
 2. Then random samples of data within the clusters
- For out-of-bag predictions:
 - consider an observation i to be out-of-bag if its cluster was not drawn at the bootstrapping stage

(Conditional) average treatment effect

```
average_treatment_effect(cf, subset = X1>0)
```

calculates augmented inverse-propensity weighting CATE (Robins, Rotnitzky, and Zhao, 1994) for the subset of observations where X1>0

$$\hat{\tau}_j = \frac{1}{n_j} \sum_{\{i: A_i=j\}} \hat{\Gamma}_i, \quad \hat{\tau} = \frac{1}{J} \sum_{j=1}^J \hat{\tau}_j, \quad \hat{\sigma}^2 = \frac{1}{J(J-1)} \sum_{j=1}^J (\hat{\tau}_j - \hat{\tau})^2,$$
$$\hat{\Gamma}_i = \hat{\tau}^{(-i)}(X_i) + \frac{W_i - \hat{e}^{(-i)}(X_i)}{\hat{e}^{(-i)}(X_i)(1 - \hat{e}^{(-i)}(X_i))} \left(Y_i - \hat{m}^{(-i)}(X_i) - (W_i - \hat{e}^{(-i)}(X_i)) \hat{\tau}^{(-i)}(X_i) \right).$$

In the study:

95% CI for the ATE : 0.247 +/- 0.04

Test for the presence of heterogeneity

`test_calibration(cf)`

tests for the presence of heterogeneity in treatment effects (Chernozhukov et al, 2018):

- Creates two synthetic variables:

$$C_i = \bar{\tau}(W_i - \hat{W}_i) \quad \text{prediction by ATE}$$

$$D_i = (\hat{\tau}^{cf} - \bar{\tau})(W_i - \hat{W}_i) \quad \text{additional pred. by CF}$$

- Runs regression:

$$Y_i - \hat{Y}_i = \gamma C_i + \delta D_i$$

- if $\delta = 1 \Rightarrow$ treatment heterogeneity is well calibrated
- if $\delta = 0 \Rightarrow$ cf captures no heterogeneity

In the study, p-value of $\delta = 0.294 \Rightarrow$ No evidence of heterogeneity in TE

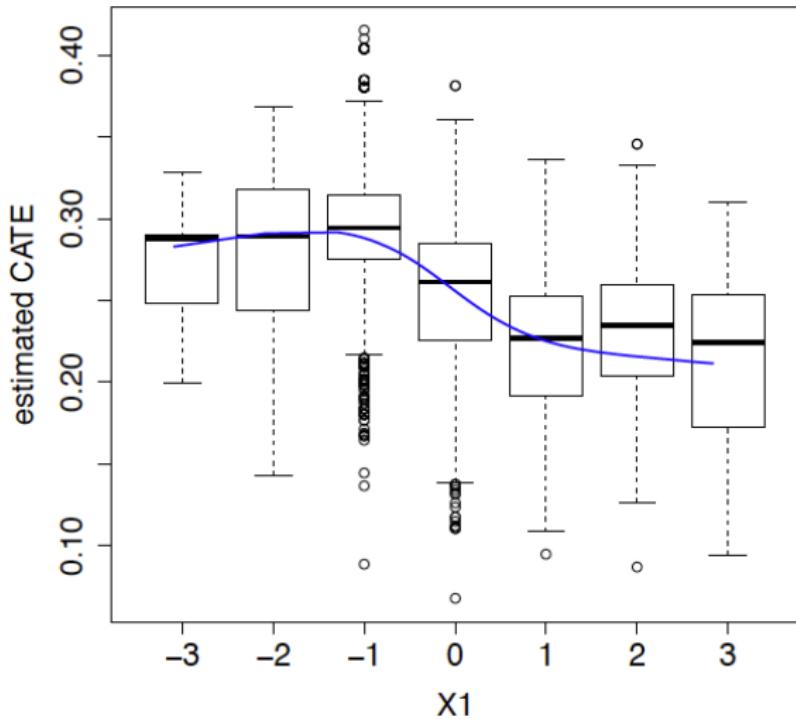
Variable importance

```
variable_importance(cf)
```

calculates a simple weighted sum of how many times feature i was split on at each depth in the forest.

In the study, 24% of splits were on X1

TE predictions by school-level mindset levels



There may be some heterogeneity within X1 levels, after performing

formal tests (See page 9)

https://madina-k.github.io/dse_mk2021

Summary

The GRF provides an elegant solution to:

- estimate ATE non-parametrically
- capture heterogeneity in TE and formally test it
- find which variables mediate the heterogeneous response the most
- all accounting for clustering and potential confounding

Moreover, the GRF can be used for any local moment equations estimation:

- quantile regressions
- IV estimation
- etc.