

Lecture 24: Causal Trees (Cont.)

Big Data and Machine Learning for Applied Economics
Econ 4676

Ignacio Sarmiento-Barbieri

Universidad de los Andes

November 5, 2020

Agenda

- 1 Recap: Causal Trees
 - Causality Review: ATE, CATE, HTE
- 2 Causal Tree: Theory
 - Honest Inference for Treatment Effects
 - Observational Studies with Unconfoundedness
- 3 Causal Forests
- 4 Review & Next Steps
- 5 Further Readings

Treatment Effects

- ▶ We observe a sequence of triples $\{(W_i, Y_i, X_i)\}_i^N$, where
 - ▶ $W_i \in \{0, 1\}$: is a binary variable indicating whether the individual was treated (1) or not (0)
 - ▶ $Y_i^{obs} \in \mathbb{R}$: a real variable indicating the observed outcome for that individual
 - ▶ X_i : is a p -dimensional vector of observable pre-treatment characteristics
- ▶ Moreover, in the Neyman-Rubin potential-outcomes framework, we will denote by
 - ▶ $Y_i(1)$: the outcome unit i would attain if they received the treatment
 - ▶ $Y_i(0)$: the outcome unit i would attain if they were part of the control group

Treatment Effects

The individual treatment effect for subject i can then be written as

$$Y_i(1) - Y_i(0)$$

Unfortunately, in our data we can only observe one of these two potential outcomes.

Education (X_i)	Treated W_i	No Subsidy $Y_i(0)$	Subsidy $Y_i(1)$	Treatment effect $\tau_i = Y_i(1) - Y_i(0)$
<i>High</i>	1	?	$Y_1(1)$?
<i>High</i>	0	$Y_2(0)$?	?
<i>Low</i>	0	$Y_3(0)$?	?
<i>Low</i>	1	?	$Y_4(1)$?

Using the potential outcome notation above, the observed outcome can also be written as

$$Y_i = W_i Y_i(1) + (1 - W_i) Y_i(0)$$

Average Treatment Effects

- ▶ Computing the difference for each individual is impossible.
- ▶ But we can get the Average Treatment Effect (ATE):

$$\tau := E[Y_i(1) - Y_i(0)] \quad (1)$$

- ▶ Heterogeneous Treatment Effects: Same treatment may affect different individuals differently
- ▶ Conditional Average Treatment Effect (CATE)

$$\tau(x) := E[Y_i(1) - Y_i(0) | X_i = x] \quad (2)$$

Heterogeneous Treatment Effects

Concerns

- ▶ Issues:
 - ▶ Ad hoc searches for particularly responsive subgroups may mistake noise for a true treatment effect.
 - ▶ Concerns about ex-post “data-mining” or p-hacking
 - ▶ preregistered analysis plan can protect against claims of data mining
 - ▶ But may also prevent researchers from discovering unanticipated results and developing new hypotheses
- ▶ But how is researcher to predict all forms of heterogeneity in an environment with many covariates?
- ▶ Athey and Imbens to the rescue
 - ▶ Allow researcher to specify set of potential covariates
 - ▶ Data-driven search for heterogeneity in causal effects with valid standard errors

Heterogeneous Treatment Effects

- ▶ Before proceeding we need to make a couple of assumptions
- ▶ Assumption 1: Unconfoundedness

$$Y_i(1), Y_i(0) \perp W_i \mid X_i \quad (3)$$

- ▶ The *unconfoundedness* assumption states that, once we condition on observable characteristics, the treatment assignment is independent to how each person would respond to the treatment.
- ▶ i.e., the rule that determines whether or not a person is treated is determined completely by their observable characteristics.
- ▶ This allows, for example, for experiments where people from different genders get treated with different probabilities,
- ▶ **rules out** experiments where people self-select into treatment due to some characteristic that is not observed in our data.

Heterogeneous Treatment Effects

► Assumption 2: Overlap

$$\forall x \in \text{supp}(X), \quad 0 < P(W = 1 \mid X = x) < 1 \quad (4)$$

- The *overlap* assumption states that at every point of the covariate space we can always find treated and control individuals.
- i.e., in order to estimate the treatment effect for a person with particular characteristics $X_i = x$, we need to ensure that we are able to observe treated and untreated people with those same characteristics so that we can compare their outcomes.

Causal Tree: Theory

- ▶ Work well in RCTs
- ▶ Issue: we do not observe the ground truth
- ▶ Honest estimation (Innovation):
 - ▶ One sample to choose partition
 - ▶ One sample to estimate leaf effects
- ▶ Why is the split critical?
- ▶ Fitting both on the training sample risks overfitting: Estimating many “heterogeneous effects” that are really just noise idiosyncratic to the sample.
- ▶ We want to search for true heterogeneity, not noise

Trees

- ▶ A simple tree

$$MSE_0 = \frac{1}{N} \sum (Y_i - \bar{Y})^2 \quad \text{All observations}$$

$$MSE_1 = \frac{1}{N} \sum (Y_i - \bar{Y}_{j:x_j \in l(x_i|\Pi)})^2 \quad X_i < c_1 \quad X_i \geq c_2$$

- ▶ Partition $\Pi \in P$

$$\{l_1 = \{x_i : x_i < c_1\}, l_2 = \{x_i : x_i \geq c_2\}\} \quad (5)$$

- ▶ Prediction is

$$\hat{\mu}(x) = \bar{Y}_{j:x_j \in l(x_i|\Pi)} \quad (6)$$

The Honest Target: Athey and Imbens Innovation

- ▶ Given a partition Π define

$$MSE_{\mu}(S^{te}, S^{est}, \Pi) = \frac{1}{\#(S^{te})} \sum_{i \in S^{te}} \left\{ (Y_i - \hat{\mu}(X_i, S^{est}, \Pi))^2 - Y_i^2 \right\} \quad (7)$$

- ▶ The expected MSE is the expectation of $MSE_{\mu}(S^{te}, S^{est}, \Pi)$ over estimation and test samples (independent)

$$EMSE_{\mu}(\Pi) = E_{S^{te}, S^{est}} [MSE_{\mu}(S^{te}, S^{est}, \Pi)] \quad (8)$$

The Honest Target: Athey and Imbens Innovation

- The ultimate goal is to construct and assess an algorithm $\pi(\cdot)$ that maximizes the honest criterion

$$\max Q^H(\pi) = -E_{S^{te}, S^{est}, S^{tr}} [MSE_{\mu}(S^{te}, S^{est}, S^{tr}, \pi(S^{tr}))] \quad (9)$$

- In CART the target is different (adaptive target)

$$\max Q^C(\pi) = -E_{S^{te}, S^{tr}} [MSE_{\mu}(S^{te}, S^{tr}, \pi(S^{tr}))] \quad (10)$$

The Honest Criterion

$$\max Q^H(\pi) = -E_{S^{te}, S^{est}, S^{tr}} [MSE_{\mu}(S^{te}, S^{est}, S^{tr}, \pi(S^{tr}))] \quad (11)$$

The Honest Criterion

- Understanding $EMSE_{\mu}(\Pi)$:

$$\begin{aligned} -EMSE_{\mu}(\Pi) &= -E_{S^{te}, S^{est}} \left[(Y_i - \hat{\mu}(X_i, S^{est}, \Pi))^2 - Y_i^2 \right] \\ &= -E_{S^{te}, S^{est}} \left[(Y_i - \mu(X_i, \Pi) + \mu(X_i, \Pi) - \hat{\mu}(X_i, S^{est}, \Pi))^2 - Y_i^2 \right] \\ &= -E_{S^{te}, S^{est}} \left[(Y_i - \mu(X_i, \Pi))^2 - Y_i^2 \right] \\ &\quad - E_{S^{te}, S^{est}} \left[(\mu(X_i, \Pi) - \hat{\mu}(X_i, S^{est}, \Pi))^2 \right] \\ &\quad - E_{S^{te}, S^{est}} \left[2(Y_i - \mu(X_i, \Pi))(\mu(X_i, \Pi) - \hat{\mu}(X_i, S^{est}, \Pi))^2 \right] \end{aligned} \tag{12}$$

Heterogeneous Treatment Effects

$$= -E_{(Y_i, X_i), S^{est}} \left[(Y_i - \mu(X_i, \Pi))^2 - Y_i^2 \right] - E_{X_i, S^{est}} \left[(\mu(X_i, \Pi) - \hat{\mu}(X_i, S^{est}, \Pi))^2 \right]$$

$$= -E_{(Y_i, X_i), S^{est}} \left[Y_i^2 - 2Y_i\mu(X_i, \Pi) + \mu^2(X_i, \Pi) - Y_i^2 \right] - E_{X_i, S^{est}} \left[(\mu(X_i, \Pi) - \hat{\mu}(X_i, S^{est}, \Pi))^2 \right]$$

$$= -E_{(Y_i, X_i), S^{est}} \left[-2Y_i\mu(X_i, \Pi) + \mu^2(X_i, \Pi) \right] - E_{X_i, S^{est}} \left[(\mu(X_i, \Pi) - \hat{\mu}(X_i, S^{est}, \Pi))^2 \right]$$

Note $E_{(Y_i, X_i), S^{est}}(Y_i) = E_{X_i, S^{est}}\mu(X_i, \Pi)$

$$= -E_{(Y_i, X_i), S^{est}} \left[\mu^2(X_i, \Pi) \right] - E_{X_i, S^{est}} \left[V(\hat{\mu}(X_i, S^{est}, \Pi)) \right]$$

The Honest Criterion

- ▶ How to estimate this quantities?
- ▶ First $E_{X_i, S^{est}} [V(\hat{\mu}(X_i, S^{est}, \Pi))]$

$$V(\hat{\mu}(X_i, S^{est}, \Pi)) = \frac{S_{Str}^2(l(x|\Pi))}{N^{est}(l(x|\Pi))}$$

$$\hat{E}_{X_i, S^{est}} [V(\hat{\mu}(X_i, S^{est}, \Pi)) | i \in S^{te}] = \sum_l p_l \frac{S_{Str}^2(l)}{N^{est}(l)}$$

$$= \sum_l \frac{1}{\#(l)} \frac{S_{Str}^2(l)}{N^{est}(l)}$$

$$= \frac{1}{N^{est}} \sum_{l \in \Pi} S_{Str}^2(l)$$

The Honest Criterion

- ▶ Next $E_{(Y_i, X_i), S^{est}} [\mu^2(X_i, \Pi)]$
- ▶ Note $V(\hat{\mu}|x, \Pi) = E(\hat{\mu}^2|x, \Pi) - [E(\hat{\mu}|x, \Pi)]^2$

$$\frac{S_{S^{tr}}^2(l(x|\Pi))}{N^{tr}(l(x|\Pi))} \approx \hat{\mu}^2(x|S^{tr}, \Pi) - \mu^2(x|\Pi)$$

$$\mu^2(x|\Pi) \approx \hat{\mu}^2(x|S^{tr}, \Pi) - \frac{S_{S^{tr}}^2(l(x|\Pi))}{N^{tr}(l(x|\Pi))}$$

$$\hat{E}_{X_i}(\mu^2(X_i|\Pi)) \approx \frac{1}{N^{tr}} \sum_{i \in S^{tr}} \hat{\mu}^2(x|S^{tr}, \Pi) - \sum_l \frac{1}{\#l} \frac{S_{S^{tr}}^2(l)}{N^{tr}/\#l}$$

The Honest Criterion

► Finally

$$\begin{aligned} -EMSE_{\mu}(\Pi) &= \frac{1}{N^{tr}} \sum_{i \in S^{tr}} \hat{\mu}^2(x|S^{tr}, \Pi) - \sum_l \frac{1}{N^{tr}} S_{S^{tr}}^2(l) - \frac{1}{N^{est}} \sum_{l \in \Pi} S_{S^{tr}}^2(l) \\ &= \frac{1}{N^{tr}} \sum_{i \in S^{tr}} \hat{\mu}^2(x|S^{tr}, \Pi) - \left(\frac{1}{N^{tr}} + \frac{1}{N^{est}} \right) \sum_{l \in \Pi} S_{S^{tr}}^2(l) \end{aligned}$$

Honest Inference for Treatment Effects

- ▶ Given a tree Π , define for all x and both treatment levels w the population average outcome

$$\mu(w, x | \Pi) = E[Y_i(w) | X_i \in l(x | \Pi)]$$

- ▶ The Average Treatment Effect

$$\tau(x | \Pi) = E[Y_i(1) - Y_i(0) | X_i \in l(x | \Pi)]$$

$$= \mu(1, x | \Pi) - \mu(0, x | \Pi)$$

Honest Inference for Treatment Effects

- The estimated counterparts are

$$\hat{\mu}(w, x|S, \Pi) = \frac{1}{\#(\{i \in S_w : X_i \in l(x|\Pi)\})} \sum_{i \in S_w : X_i \in l(x|\Pi)} Y_i^{obs} \quad (13)$$

$$\hat{\tau}(X, S, \Pi) = \hat{\mu}(1, x|S, \Pi) - \hat{\mu}(0, x|S, \Pi) \quad (14)$$

- Define the MSE for treatment effects as

$$MSE_{\tau}(S^{te}, S^{est}, \Pi) = \frac{1}{\#(S^{te})} \sum_{i \in S^{te}} \left\{ (\tau_i - \hat{\tau}(X_i|S^{est}, \Pi))^2 - \tau_i^2 \right\}$$

Honest Inference for Treatment Effects

Adapt $EMSE_{\mu}$ to estimate $EMSE_{\tau}$

$$-EMSE_{\mu}(\hat{S}^{tr}, S^{est}, \Pi) = \frac{1}{N^{tr}} \sum_{i \in S^{tr}} \hat{\mu}^2(X_i | S^{tr}, \Pi) - \left(\frac{1}{N^{tr}} + \frac{1}{N^{est}} \right) \sum_{l \in \Pi} S_{S^{tr}}^2(l)$$

for HTE

$$-EMSE_{\tau}(\hat{S}^{tr}, S^{est}, \Pi) = \frac{1}{N^{tr}} \sum_{i \in S^{tr}} \hat{\tau}^2(X_i | S^{tr}, \Pi) - \left(\frac{1}{N^{tr}} + \frac{1}{N^{est}} \right) \sum_{l \in \Pi} \left(\frac{S_{S^{tr}^{treat}}^2(l)}{p} + \frac{S_{S^{tr}^{control}}^2(l)}{(1-p)} \right)$$

Observational Studies with Unconfoundedness

► Athey and Imbens (2016):

“The proposed methods can be adapted to observational studies under the assumption of unconfoundedness. In that case we need to modify the estimates within leaves to remove the bias from simple comparisons of treated and control units. There is a large literature on methods for doing so,, for example, we can do so by propensity score weighting. Efficiency will improve if we renormalize the weights within each leaf and within the treatment and control group when estimating treatment effects”

Causal Forests

- ▶ Trees can be noisy. We can use forests
 - ▶ Draw a sample bootstrap of size s
 - ▶ Split the sample into Tr and Est
 - ▶ Use Tr to grow the tree
 - ▶ Use Est to estimate the leaf-specific effects
- ▶ Advantages
 - ▶ Consistent for $\tau(x)$
 - ▶ Asymptotically Normal
 - ▶ “Auto” search for HTE
- ▶ Disadvantage
 - ▶ Sample splitting (noisier estimates)

Review & Next Steps

- ▶ Problem: we never observe t_i unlike prediction that we observe Y_i
- ▶ Causal Trees search for leaves with
 - ▶ HTE across leaves
 - ▶ precisely-estimated leaf effects
- ▶ Key is the honest Criterion
- ▶ Work well with RCTs
- ▶ With selection on observables, recommendation is propensity forests?
- ▶ Next class: Causal forests demo
- ▶ Questions? Questions about software?

Further Readings

- ▶ Athey, S., & Imbens, G. (2016). Recursive partitioning for heterogeneous causal effects. *Proceedings of the National Academy of Sciences*, 113(27), 7353-7360.
- ▶ Lundberg, I (2017). Causal forests. A tutorial in high dimensional causal inference. Mimeo
- ▶ Taddy, M. (2019). *Business data science: Combining machine learning and economics to optimize, automate, and accelerate business decisions*. McGraw Hill Professional.