

(6) Heterogeneous Treatment Effects

Causal Data Science for Business Analytics

Christoph Ihl

Hamburg University of Technology

Monday, 27. May 2024



Introduction

Treatment Effect Heterogeneity: Motivation

- More comprehensive evaluation:
 - who wins or loses and by how much?
- This is useful along at least two dimensions:
 - **Informs action:**
 - More efficient allocation of public and private resources via targeting in the future:
 - Personalized policies, ads, medicine, ...
 - **Understanding:**
 - Heterogeneous effects can be suggestive for underlying mechanisms

Treatment Effect Heterogeneity: Definition

- Expected treatment effect in the target subpopulation with characteristics \mathbf{X}_i given by **Conditional Average Treatment Effect (CATE)**:
 - $\tau(\mathbf{x}) = \mathbb{E}[Y_i(1) - Y_i(0) | \mathbf{X}_i = \mathbf{x}]$
- $\mathbf{X}_i = \mathbf{H}_i \cup \mathbf{C}_i$
 - \mathbf{H}_i : motivated by the research question to understand specific effect heterogeneity in a pre-defined the target subpopulation.
 - \mathbf{C}_i : confounders that are required for identification.
- **Randomized experiments** - no confounders:
 - CATE defined with respect to considered heterogeneity variables: $\mathbf{X}_i = \mathbf{H}_i$
- **Measured Confounding** - distinguish two types of CATEs:
 - **Group ATE (GATE)** for groups G defined by H : $\tau(g) = \mathbb{E}[Y_i(1) - Y_i(0) | G_i = g]$
 - **Individualized ATE (IATE = CATE)**: $\tau(\mathbf{x}) = \mathbb{E}[Y_i(1) - Y_i(0) | \mathbf{X}_i = \mathbf{x}]$
 - most flexible/ personalized/ individualized effect prediction
 - Estimation step is affected by whether we are interested in GATEs or IATEs.

Treatment Effect Heterogeneity: Identification

- No need to establish new identification results:
 - All target parameters can be thought of as special cases of conditioning ITE on some function $f(\mathbf{X}_i = \mathbf{x})$
 - And by the Law of Iterated Expectations (LIE):

$$\begin{aligned}\mathbb{E}[Y_i(1) - Y_i(0) | f(\mathbf{X}_i) = f(\mathbf{x})] &= \mathbb{E}[\mathbb{E}[Y_i(1) - Y_i(0) | \mathbf{X}_i = \mathbf{x}, f(\mathbf{X}_i = \mathbf{x})] | f(\mathbf{X}_i = \mathbf{x})] \\ &= \mathbb{E}[\mathbb{E}[Y_i(1) - Y_i(0) | \mathbf{X}_i = \mathbf{x}] | f(\mathbf{X}_i = \mathbf{x})]\end{aligned}$$

- As $\mathbf{X}_i = \mathbf{H}_i \cup \mathbf{C}_i$ is assumed to contain all confounders, the inner expectation $\mathbb{E}[Y_i(1) - Y_i(0) | \mathbf{X}_i = \mathbf{x}]$ is identified in randomized experiments or under measured confounding
- \Rightarrow All aggregations with respect to a function $f(\mathbf{X}_i)$ are also identified.

Group Average Treatment Effects

Group ATEs: Examples

- **Group ATE (GATE):** $\tau(\mathbf{g}) = \mathbb{E}[Y_i(1) - Y_i(0) | \mathbf{G}_i = \mathbf{g}]$
- Examples for subgroups of interest:
 - Mutually exclusive *subgroups*, e.g.: $\mathbf{G} = \{\text{female}, \text{male}\}$,
 $\mathbf{G} = \{\text{age} < 50, \text{age} \geq 50\}$,
 $\mathbf{G} = \{\text{age} < 50 \ \& \ \text{female}, \text{age} \geq 50 \ \& \ \text{female}, \text{age} \geq 50 \ \& \ \text{male}, \dots\}$, ...
 - Single or low-dimensional *continuous variable*, e.g.: $\mathbf{G} = \text{age}$, $\mathbf{G} = \text{income}$, ...
 - Other *functions or small subsets* of \mathbf{X}_i
- Groups should be *pre-determined* and not be the result of data snooping

Group ATEs: Estimation

- Three strategies:
 1. Stratify the data and rerun the analysis for each subgroup.
 - Downside: Requires a lot of data and computation, can lead to high variance estimates for small subgroups.
 2. Specify an interaction term in an OLS regression model:
 - $Y_i = \beta_0 + \tau T_i + \beta_{G_i} G_i + \beta_{T_i G_i} T_i G_i + \beta_{X_i} X_i + \epsilon_i$
 - Downside: Requires a correct model specification, can be sensitive to misspecification.
 3. Double Machine Learning with AIPW model to estimate the GATEs directly.

Group ATEs: Double Machine Learning

- Previous lecture: ATE (AIPW) can be estimated as mean of a pseudo-outcome:

- $\tau_{ATE}^{AIPW} = \frac{1}{N} \sum_{i=1}^n \tilde{\tau}_{i,ATE}^{AIPW}$

- Pseudo-outcome is given by:

- $\tilde{\tau}_{i,ATE}^{AIPW} = \mu(1, \mathbf{X}_i) - \mu(0, \mathbf{X}_i) + \frac{T_i(Y_i - \mu(1, \mathbf{X}_i))}{\hat{e}_1(\mathbf{X}_i)} - \frac{(1-T_i)(Y_i - \mu(0, \mathbf{X}_i))}{\hat{e}_0(\mathbf{X}_i)}$

- Equivalent to a linear regression model with pseudo-outcome and constant:

- $\tilde{\tau}_{i,ATE}^{AIPW} = \alpha + \epsilon_i$ with $\hat{\alpha} = \tau_{ATE}^{AIPW}$

- Can be extended with heterogeneity variable(s) G_i :

- $\tilde{\tau}_{i,ATE}^{AIPW} = \alpha + \beta G_i + \epsilon_i$

- => Modelling the level of the effect, not the level of the outcome.

Group ATEs: Advantages of DML

- Neyman-orthogonality of $\tilde{\tau}_{i\text{ATE}}^{\text{AIPW}}$ allows to apply standard statistical inference ([Semenova and Chernozhukov, 2021](#)).
- Computationally less expensive than subgroup analyses
 - Only one additional OLS, no new nuisance parameters).
- More flexible than specifying interaction terms in a linear model, as we flexibly adjust for confounding by ML methods.
- As $\tilde{\tau}_{i\text{ATE}}^{\text{AIPW}}$ is an unbiased signal, i.e. $\mathbb{E}[\tilde{\tau}_{i\text{ATE}}^{\text{AIPW}} | G_i = g] = \tau(g)$, to regress the pseudo-outcome $\tilde{\tau}_{i\text{ATE}}^{\text{AIPW}}$ on low-dimensional G_i we can either use
 - OLS or series regression ([Semenova and Chernozhukov, 2021](#)).
 - Kernel regression ([Fan et al., 2022](#); [Zimmert & Lechner, 2019](#)).

Group ATEs: Proof of DML

- Proof that $\mathbb{E}[\tilde{\tau}_{ATE}^{AIPW} \mid G_i = g] = \tau(g)$:

$$\begin{aligned}
 \mathbb{E}[\tilde{\tau}_{ATE}^{AIPW} \mid G_i = g] &= \mathbb{E} \left[\mu(1, \mathbf{X}_i) + \frac{T_i(Y_i - \mu(1, \mathbf{X}_i))}{e(\mathbf{X}_i)} - \mu(0, \mathbf{X}_i) - \frac{(1 - T_i)(Y_i - \mu(0, \mathbf{X}_i))}{1 - e(\mathbf{X}_i)} \mid G_i = g \right] \\
 &\stackrel{LIE}{=} \mathbb{E} \left[\underbrace{\mathbb{E} \left[\mu(1, \mathbf{X}_i) + \frac{T_i(Y_i - \mu(1, \mathbf{X}_i))}{e(\mathbf{X}_i)} \mid \mathbf{X}_i = \mathbf{x} \right]}_{\text{CAPO-AIPW} \Rightarrow \mathbb{E}[Y_i(1) \mid \mathbf{X}_i = \mathbf{x}]} - \underbrace{\mathbb{E} \left[\mu(0, \mathbf{X}_i) + \frac{(1 - T_i)(Y_i - \mu(0, \mathbf{X}_i))}{1 - e(\mathbf{X}_i)} \mid \mathbf{X}_i = \mathbf{x} \right]}_{\text{CAPO-AIPW} \Rightarrow \mathbb{E}[Y_i(0) \mid \mathbf{X}_i = \mathbf{x}]} \mid G_i = g \right] \\
 &= \mathbb{E} \left[\mathbb{E}[Y_i(1) \mid \mathbf{X}_i = \mathbf{x}] - \mathbb{E}[Y_i(0) \mid \mathbf{X}_i = \mathbf{x}] \mid G_i = g \right] \\
 &\stackrel{LIE}{=} \mathbb{E}[Y_i(1) - Y_i(0) \mid G_i = g] \\
 &= \tau(g)
 \end{aligned}$$

- Law of Iterated Expectations uses that G_i is a function of \mathbf{X}_i .

Group ATEs: Example based on DML

- Assess the effect of 401(k) program participation on net financial assets of 9,915 households in the US in 1991.
- First step (not shown): Estimate τ_{ATE}^{AIPW} using DoubleML.

```
1 # Get the individual ATEs (pseudo-outcomes)
2 data$ate_i <- dml_irm_forest[["psi_b"]] # get numerator of score function
3 mean_ate = mean(data$ate_i) # mean of pseudo outcomes = ATE
4
5 library(estimatr) # for linear robust post estimation
6 summary(lm_robust(ate_i ~ hown, data = data))
```

Estimates and significance testing of the effect of target variables

```
Estimate Std. Error t value Pr(>|t|)
e401      8206      1106   7.421 1.16e-13 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
Call:
lm_robust(formula = ate_i ~ hown, data = data)
```

Standard error type: HC2

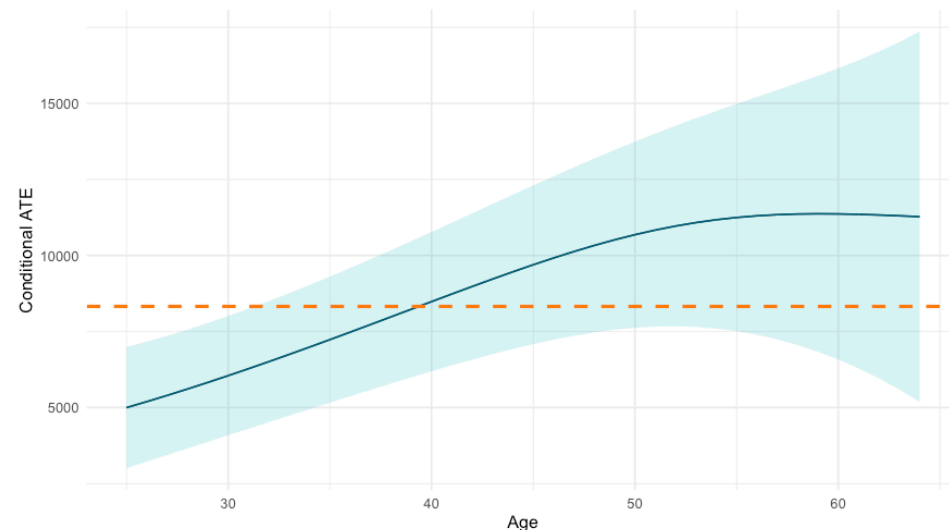
Coefficients:

	Estimate	Std. Error	t value	Pr(> t)	CI Lower	CI Upper	DF
(Intercept)	3477	711	4.890	1.025e-06	2083	4870	9913
hown	7445	1835	4.058	4.990e-05	3849	11041	9913

Multiple R-squared: 0.00106 , Adjusted R-squared: 0.0009588

F-statistic: 16.47 on 1 and 9913 DF, p-value: 4.99e-05

```
1 library(np) # for kernel post estimation
2 age = data$age
3 ate_i = data$ate_i
4 np_model = npreg(ate_i ~ age) # kernel regression
5 plot(np_model) # plot the kernel regression
```



Metalearners

Predicting Individualized ATEs

- Group-level heterogeneity variables were hand-picked.
- Now predict individualized treatment effects based on all covariates \mathbf{X}_i :
 - **Individualized ATE (IATE = CATE):** $\tau(\mathbf{x}) = \mathbb{E}[Y_i(1) - Y_i(0) | \mathbf{X}_i = \mathbf{x}]$
 - Conditional expectation with unobserved outcome (counterfactuals)
- Given the assumptions of observed confounding, we can write the CATE as:
 - $\tau(\mathbf{x}) = \mathbb{E}[Y_i(1) - Y_i(0) | \mathbf{X}_i = \mathbf{x}] = \mathbb{E}[Y_i | T_i = 1, \mathbf{X}_i = \mathbf{x}] - \mathbb{E}[Y_i | T_i = 0, \mathbf{X}_i = \mathbf{x}]$
 - which can be approximated with ML.

S-Learner and T-Learner

- S-learner:

- 1. Use ML estimator of your choice to fit outcome model using \mathbf{X}_i AND T_i in the **full sample**: $\mu(T_i; \mathbf{X}_i)$.
- 2. Estimate CATE as $\tau(\mathbf{x}) = \mu(1; \mathbf{X}_i) - \mu(0; \mathbf{X}_i)$.

- T-learner:

- 1. Use ML estimator of your choice to fit model $\mu(1; \mathbf{X}_i)$ in **treated subsample**.
- 2. Use ML estimator of your choice to fit model $\mu(0; \mathbf{X}_i)$ in **control subsample**.
- 3. Estimate CATE as $\tau(\mathbf{x}) = \mu(1; \mathbf{X}_i) - \mu(0; \mathbf{X}_i)$.

S-Learner and T-Learner: Example

- Assess the effect of 401(k) program participation on net financial assets of 9,915 households in the US in 1991.
- Examples without proper cross-fitting.

```

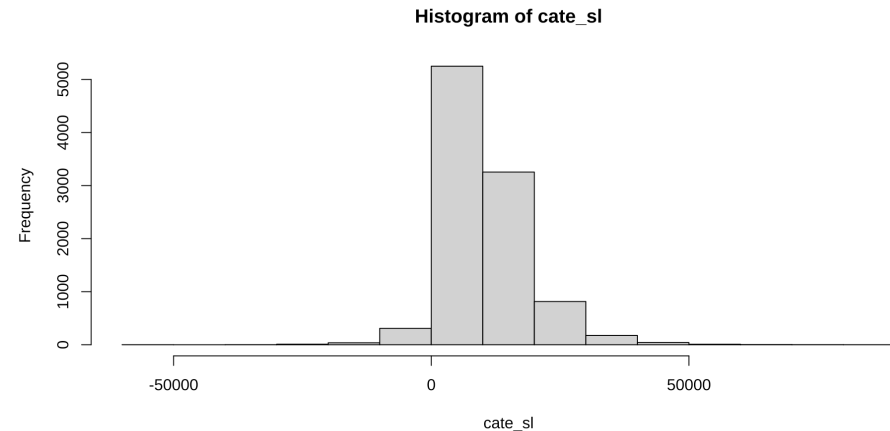
1 library(hdm) # for the data
2 library(grf) # generalized random forests, could also use mlr3
3
4
5 # Get data
6 data(pension)
7 # Outcome
8 Y = pension$net_tfa
9 # Treatment
10 T = pension$sp401
11 # Create main effects matrix
12 X = model.matrix(~ 0 + age + db + educ + fsize + hown + inc + male + mar
13
14 # Implement the S-Learner
15 TX = cbind(T,X)
16 rf = regression_forest(TX,Y)
17 T0X = cbind(rep(0,length(Y)),X)
18 T1X = cbind(rep(1,length(Y)),X)
19 cate_sl = predict(rf,T1X)$predictions - predict(rf,T0X)$predictions
20 hist(cate_sl)

```

```

1 # Implement the T-Learner
2 rfmu1 = regression_forest(X[T==1,],Y[T==1])
3 rfmu0 = regression_forest(X[T==0,],Y[T==0])
4 cate_t1 = predict(rfmu1, X)$predictions - predict(rfmu0, X)$predictions
5 hist(cate_t1)

```



S-Learner and T-Learner: Disadvantage

- The prediction problems **do not know of joint goal** to approximate a difference:

- $\mu(1; \mathbf{X}_i)$ minimizes $\text{MSE}(\mu(1; \mathbf{x})) = \mathbb{E}[(\mu(1; \mathbf{x}) - \mu(1; \mathbf{X}_i))^2]$.
- $\mu(0; \mathbf{X}_i)$ minimizes $\text{MSE}(\mu(0; \mathbf{x})) = \mathbb{E}[(\mu(0; \mathbf{x}) - \mu(0; \mathbf{X}_i))^2]$.
- BUT they **should aim to minimize**:

$$\begin{aligned}
 \text{MSE}(\tau(\mathbf{x})) &= \mathbb{E}[(\tau(\mathbf{x}) - \tau(\mathbf{X}_i))^2] \\
 &= \mathbb{E}[(\mu(1, \mathbf{x}) - \mu(0, \mathbf{x})) - (\mu(1, \mathbf{X}_i) - \mu(0, \mathbf{X}_i))]^2] \\
 &= \mathbb{E}[(\mu(1, \mathbf{x}) - \mu(1, \mathbf{X}_i))^2] + \mathbb{E}[(\mu(0, \mathbf{x}) - \mu(0, \mathbf{X}_i))^2] \\
 &\quad - 2\mathbb{E}[(\mu(1, \mathbf{x}) - \mu(1, \mathbf{X}_i))(\mu(0, \mathbf{x}) - \mu(0, \mathbf{X}_i))] \\
 &= \text{MSE}(\mu(1, \mathbf{x})) + \text{MSE}(\mu(0, \mathbf{x})) - 2\text{MCE}(\mu(1, \mathbf{x}), \mu(0, \mathbf{x}))
 \end{aligned}$$

- **Lechner (2018)** calls the additional term **Mean Correlated Error (MCE)**: correlated errors matter less
- Example - both make same error: $\mu(1; \mathbf{X}_i) = \mu(1; \mathbf{X}_i) + 2$ and $\mu(0; \mathbf{X}_i) = \mu(0; \mathbf{X}_i) + 2$
 - But their CATE would still be on point: $\text{MSE}(\tau(\mathbf{x})) = 4 + 4 - 2(2 \cdot 2) = 0$
- Example - errors go in different direction: $\mu(1; \mathbf{X}_i) = \mu(1; \mathbf{X}_i) + 2$ and $\mu(0; \mathbf{X}_i) = \mu(0; \mathbf{X}_i) - 2$
 - But their CATE would be off: $\text{MSE}(\tau(\mathbf{x})) = 4 + 4 - 2(2 \cdot (-2)) = 16$

Two Approaches to Improvements

1. **Modify** supervised ML methods to target causal effect estimation

- Method specific, e.g.:
 - Causal tree ([Athey and Imbens, 2016](#))
 - Causal forest ([Athey, Tibshirani & Wager, 2019](#))
- Not covered here (does not scale very well to high-dimensional data)

2. **Combine** supervised ML methods to target causal effect estimation

- Generic approach – [Metalearners](#), e.g.:
 - X-learner ([Künzel et al., 2019](#))
 - not covered here; handles sample imbalance, but not doubly robust
 - R-learner
 - DR-learner

What are Metalearners?

- Metalearners **combine multiple supervised ML steps** in a pipeline that outputs predicted CATEs.
- The common ones require the following steps:
 1. **Estimate nuisance parameters** using suitable ML method.
 2. Plug them into a clever **minimization problem** targeting CATE.
 3. **Solve the minimization problem** using suitable ML method.
 4. **Predict** CATE using the model learned in 3.
- Most popular ML methods are suitable and can be applied in steps 1, 3 and 4.
- Like for standard prediction methods, **statistical inference is usually not available**.

R-learner: Idea

- **Partially linear model**, but now allowing for treatment effects that vary with \mathbf{X} :

$$\blacksquare Y_i = \underbrace{\tau(\mathbf{X}_i)}_{\mu(\mathbf{X}_i)} T_i + g(\mathbf{X}_i) + \epsilon_{Y_i}, \quad \mathbb{E}(\underbrace{\epsilon_{Y_i}}_{e(\mathbf{X}_i)} | T_i, \mathbf{X}_i) = 0$$

$$\blacksquare \Rightarrow \underbrace{Y_i - \mathbb{E}[Y_i | \mathbf{X}_i]}_{\text{outcome residual}} = \tau(\mathbf{X}_i) \underbrace{(T_i - \mathbb{E}[T_i | \mathbf{X}_i])}_{\text{treatment residual}} + \epsilon_{Y_i}$$

- This motivates the R-learner of **Nie and Wager, 2020**:

- $\tau_{\text{RL}}(\mathbf{x}) = \arg \min_{\tau} \sum_{i=1}^n (Y_i - \hat{\mu}(\mathbf{X}_i) - \tau(\mathbf{X}_i)(T_i - \hat{e}(\mathbf{X}_i)))^2$
- with cross-fitted high-quality nuisance parameters from first step.
- But how to estimate it?

R-learner with Linear ML-Methods

- CATE as linear function $\tau(\mathbf{X}_i) = \beta' \mathbf{X}_i$:

$$\begin{aligned}\hat{\beta}_{\text{RL}} &= \arg \min_{\beta} \sum_{i=1}^N (Y_i - \mu(\mathbf{X}_i) - \underbrace{\beta' (T_i - e(\mathbf{X}_i))}_{=\tilde{\mathbf{X}}_i} \mathbf{X}_i)^2 \\ &= \arg \min_{\beta} \sum_{i=1}^N (Y_i - \mu(\mathbf{X}_i) - \beta' \tilde{\mathbf{X}}_i)^2\end{aligned}$$

- $\tilde{\mathbf{X}}_i = (T_i - e(\mathbf{X}_i))\mathbf{X}_i$ are modified / pseudo-covariates.
- $\tau_{\text{RL}}(\mathbf{x}) = \hat{\beta}_{\text{RL}} \mathbf{x} \neq \hat{\beta}_{\text{RL}} \tilde{\mathbf{x}}$ is the estimated CATE for a specific \mathbf{x} .
- All linear shrinkage estimators (Lasso and friends) can be applied, nuisance parameters can still be estimated with non-linear ML.

R-learner with Generic ML-Methods

- If we are not willing to impose linearity of the CATE, we can rewrite the R-learner:

$$\begin{aligned}
 \tau_{\text{RL}}(\mathbf{x}) &= \arg \min_{\tau} \sum_{i=1}^n (Y_i - \mu(\mathbf{X}_i) - \tau(\mathbf{X}_i)(T_i - e(\mathbf{X}_i)))^2 \\
 &= \arg \min_{\tau} \sum_{i=1}^n \frac{(T_i - e(\mathbf{X}_i))^2}{(T_i - e(\mathbf{X}_i))^2} (Y_i - \mu(\mathbf{X}_i) - \tau(\mathbf{X}_i)(T_i - e(\mathbf{X}_i)))^2 \\
 &= \arg \min_{\tau} \sum_{i=1}^n (T_i - e(\mathbf{X}_i))^2 \left(\frac{Y_i - \mu(\mathbf{X}_i) - \tau(\mathbf{X}_i)(T_i - e(\mathbf{X}_i))}{T_i - e(\mathbf{X}_i)} \right)^2 \\
 &= \arg \min_{\tau} \sum_{i=1}^n (T_i - e(\mathbf{X}_i))^2 \left(\frac{Y_i - \mu(\mathbf{X}_i)}{T_i - e(\mathbf{X}_i)} - \tau(\mathbf{X}_i) \right)^2
 \end{aligned}$$

- Supervised ML methods that can deal with weighted minimization (e.g. neural nets, random forest, boosting, ...) with
 - weights: $(T_i - e(\mathbf{X}_i))^2$.
 - pseudo-outcome: $\frac{Y_i - \mu(\mathbf{X}_i)}{T_i - e(\mathbf{X}_i)}$.
 - the unmodified covariates: \mathbf{X}_i .

DR-learner

- Recall the pseudo-outcome of the AIPW-ATE from previous lecture and condition on \mathbf{X}_i (same “trick” as for GATE estimation):

$$\tau_{\text{DR}}(\mathbf{x}) = \mathbb{E} \left[\underbrace{\mu(1, \mathbf{X}_i) - \mu(0, \mathbf{X}_i)}_{\tilde{\tau}_{\text{ATE}}^{\text{AIPW}}} + \frac{T_i(Y_i - \mu(1, \mathbf{X}_i))}{e_1(\mathbf{X}_i)} - \frac{(1-T_i)(Y_i - \mu(0, \mathbf{X}_i))}{e_0(\mathbf{X}_i)} \mid \mathbf{X}_i = \mathbf{x} \right]$$

$$\tau_{\text{DR}}(\mathbf{x}) = \mathbb{E} \left[\tilde{\tau}_{\text{ATE}}^{\text{AIPW}} \mid \mathbf{X}_i = \mathbf{x} \right]$$

- DR-learner** by Kennedy (2020) uses $\tilde{\tau}_{\text{ATE}}^{\text{AIPW}}$ in a generic ML problem:

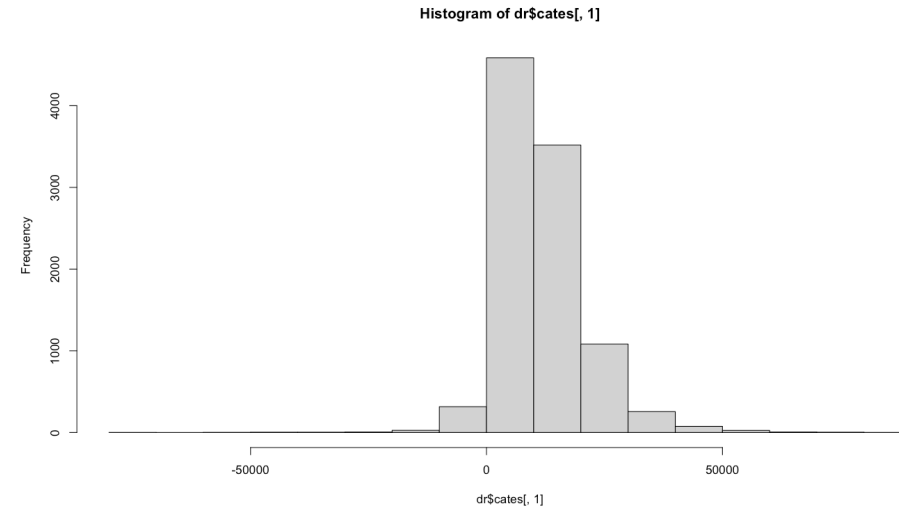
$$\tau_{\text{RL}}^{\hat{}}(\mathbf{x}) = \arg \min_{\tau} \sum_{i=1}^N \left(\tilde{\tau}_{\text{ATE}}^{\text{AIPW}} - \tau(\mathbf{X}_i) \right)^2$$

- Cross-fitting**: in 4 subsamples (1) train a model for $e(\hat{\mathbf{X}}_i)$, (2) train a model for $\mu(\hat{\mathbf{X}}_i)$, (3) construct $\tilde{\tau}_{\text{ATE}}^{\text{AIPW}}$ and regress on \mathbf{X}_i , (4) predict $\tau_{\text{RL}}^{\hat{}}(\mathbf{x})$. Then rotate.

DR-learner: Example

- Assess the effect of 401(k) program participation on net financial assets of 9,915 households in the US in 1991.

```
1 library(hdm) # for the data
2 library(causalDML) # generalized random forests, could also use mlr3
3
4
5 # Get data
6 data(pension)
7 # Outcome
8 Y = pension$net_tfa
9 # Treatment
10 T = pension$p401
11 # Create main effects matrix
12 X = model.matrix(~ 0 + age + db + educ + fsize + hown + inc + male + mar
13
14 # Implement the DR-Learner
15 dr = dr_learner(Y,T,X,
16               ml_w = list(create_method("forest_grf")),
17               ml_y = list(create_method("forest_grf")),
18               ml_tau = list(create_method("forest_grf"))
19 )
20
21 # DR-learner distribution of B-A
22 hist(dr$scates[,1])
```



HTE Evaluation

How to evaluate estimated CATEs?

1. **Descriptive**: histogram, kernel density plots, box plots, etc. ...
 2. **Inference**: test whether effect heterogeneity is systematic or just noise.
 3. Explore what drives the heterogeneous effects.
- Challenges with inference:
 - Unique to causal ML: Due to missing counterfactual, we cannot benchmark predicted against effect => no classic out-of-sample testing.
 - Shared with supervised ML: statistical inference for predicted CATE is not available or at least challenging.
 - Approach to inference:
 - Rather than (consistent) estimation of & inference on the individual CATEs directly, derive summary statistics of their (noisy) distribution.
 - Test joint hypothesis that there is effect heterogeneity & the applied estimation method is able to detect it at least partially.
 - We discuss the three methods proposed by [Chernozhukov et al. \(2017-2023\)](#):
 - Best Linear Predictor (**BLP**).
 - High-vs.-low Sorted Group Average Treatment Effect (**GATES**).
 - Classification Analysis (**CLAN**) to explore what drives the heterogeneous effects.

Best Linear Predictor (BLP) – Definition

- BLP is defined as the solution of the **hypothetical regression of the true CATE on the demeaned predicted CATE**:

Definition “Best Linear Predictor (BLP)”

The best linear predictor of $\tau(\mathbf{X}_i)$ by $\tau(\mathbf{X}_i)$ is the solution to:

$$(\beta_1, \beta_2) = \arg \min_{\tilde{\beta}_1, \tilde{\beta}_2} \mathbb{E} \left[\left(\tau(\mathbf{X}_i) - \tilde{\beta}_1 - \tilde{\beta}_2 (\tau(\mathbf{X}_i) - \mathbb{E}[\tau(\mathbf{X}_i)]) \right)^2 \right]$$

- which, if exists, is defined as

$$\mathbb{E}[\tau(\mathbf{X}_i) | \tau(\mathbf{X}_i)] := \underbrace{\beta_1}_{\text{demeaned prediction}} + \beta_2 (\tau(\mathbf{X}_i) - \mathbb{E}[\tau(\mathbf{X}_i)])$$

- where

$$\begin{aligned} \beta_1 &= \mathbb{E}[\tau(\mathbf{X}_i)] = \text{ATE (because of the demeaning)} \\ \beta_2 &= \frac{\text{Cov}[\tau(\mathbf{X}_i), \tau(\mathbf{X}_i)]}{\text{Var}[\tau(\mathbf{X}_i)]} \end{aligned}$$

BLP – Interpretation

- $\beta_2 = \frac{\text{Cov}[\tau(\mathbf{X}_i), \tau(\mathbf{X}_i)]}{\text{Var}[\tau(\mathbf{X}_i)]} = 1$ if $\tau(\mathbf{X}_i) = \tau(\mathbf{X}_i)$ (what we would like to see)
- $\beta_2 = 0$ if $\text{Cov}[\tau(\mathbf{X}_i), \tau(\mathbf{X}_i)] = 0$, which can have **two reasons**:
 - 1. $\tau(\mathbf{X}_i)$ is **constant** (no heterogeneity to detect).
 - 2. $\tau(\mathbf{X}_i)$ is **not constant** but the estimator is not capable of finding it (bad estimator and/or not enough observations).
- Therefore, testing $H_0 : \beta_2 = 0$ is a **joint test** of
 - i. existence of heterogeneity and
 - ii. the estimators capability to find it.

BLP – Identification Strategy A

Strategy A: **Weighted residual BLP**

- $(\beta_1, \beta_2) = \arg \min_{\tilde{\beta}_1, \tilde{\beta}_2} \mathbb{E} \left[\omega(\mathbf{X}_i) \left(Y_i - \tilde{\beta}_1 (T_i - e(\mathbf{X}_i)) - \tilde{\beta}_2 (T_i - e(\mathbf{X}_i)) (\tau(\mathbf{X}_i) - \mathbb{E}[\tau(\mathbf{X}_i)]) - \alpha \mathbf{X}_i^C \right) \right]$
- where:
 - $\omega(\mathbf{X}_i) = \frac{1}{e(\mathbf{X}_i)(1-e(\mathbf{X}_i))}$
 - \mathbf{X}_i^C is not required for identification, but contains optional functions of \mathbf{X}_i to reduce estimation noise, e.g. $[1, \mu(0, \mathbf{X}_i), e(\mathbf{X}_i), e(\mathbf{X}_i)\tau(\mathbf{X}_i)]$
- See Appendix A in [Chernozhukov et al. \(2017-2023\)](#) for a detailed derivation.

BLP – Identification Strategy B

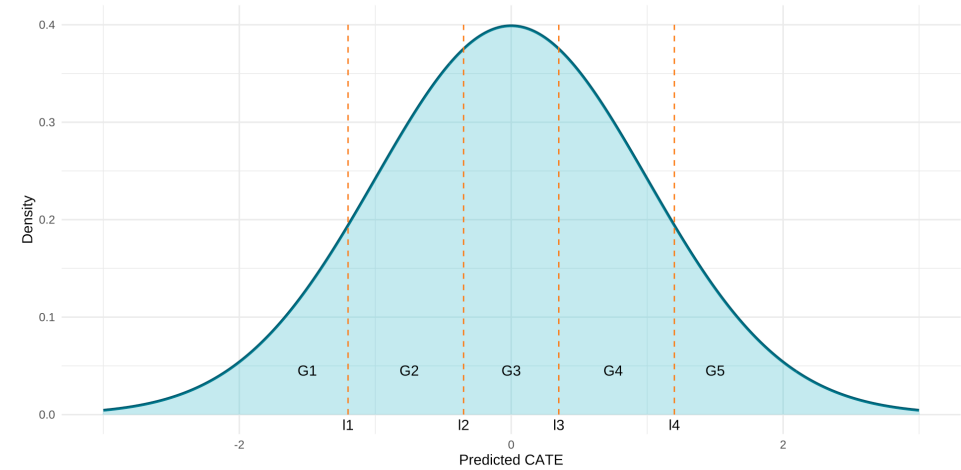
Strategy B: **Horvitz-Thompson BLP**

- $(\beta_1, \beta_2) = \arg \min_{\tilde{\beta}_1, \tilde{\beta}_2} \mathbb{E} \left[\left(H_i Y_i - \tilde{\beta}_1 - \tilde{\beta}_2 (\tau(\mathbf{X}_i) - \mathbb{E}[\tau(\mathbf{X}_i)]) - \alpha H_i \mathbf{X}_i^C \right)^2 \right]$
- where:
 - $H_i = \frac{T_i - e(\mathbf{X}_i)}{e(\mathbf{X}_i)(1 - e(\mathbf{X}_i))}$ are the Horvitz-Thompson (IPW) weights.
 - $H_i Y_i$ serves as a **pseudo-outcome**.
 - \mathbf{X}_i^C is not required for identification, but contains optional functions of \mathbf{X}_i to reduce estimation noise, e.g. $[1, \mu(0, \mathbf{X}_i), e(\mathbf{X}_i), e(\mathbf{X}_i)\tau(\mathbf{X}_i)]$
- See Appendix A in **Chernozhukov et al. (2017-2023)** for a detailed derivation.

Sorted Group Average Treatment Effect (GATES)

- Idea:

- slice the distribution of $\tau(\mathbf{X}_i)$ into K parts and compare the average treatment effect of individuals within each slice.
- if $\tau(\mathbf{X}_i)$ is a good approximation of $\tau(\mathbf{X}_i)$, then we expect to observe the following monotonicity:
 $\gamma_1 \leq \gamma_2 \leq \dots \leq \gamma_K$.



Definition "Sorted Group Average Treatment Effect (GATES)"

$$\gamma_k := \mathbb{E}[\tau(\mathbf{X}_i) | G_k], k = 1, \dots, K$$

- where $G_k = \{\tau(\mathbf{X}_i) \in I_k\}$ with $I_k = [l_{k-1}, l_k)$ and $-\infty = l_0 < l_1 < \dots < l_K = \infty$.

GATES – Identification

Strategy A: Weighted residual GATES

- $(\gamma_1, \dots, \gamma_K) = \arg \min_{\tilde{\gamma}_1, \dots, \tilde{\gamma}_K} \mathbb{E} \left[\omega(\mathbf{X}_i) \left(Y_i - \sum_k \tilde{\gamma}_k (T_i - e(\mathbf{X}_i)) \mathbb{1}[G_k] - \alpha \mathbf{X}_i^C \right)^2 \right]$
 - where $\omega(\mathbf{X}_i) = \frac{1}{e(\mathbf{X}_i)(1-e(\mathbf{X}_i))}$.

Strategy B: Horvitz-Thompson GATES

- $(\gamma_1, \dots, \gamma_K) = \arg \min_{\tilde{\gamma}_1, \dots, \tilde{\gamma}_K} \mathbb{E} \left[\left(H_i Y_i - \sum_k \tilde{\gamma}_k \mathbb{1}[G_k] - \alpha H_i \mathbf{X}_i^C \right)^2 \right]$
 - where $H_i Y_i$ serves as a **pseudo-outcome** and $H_i = \frac{T_i - e(\mathbf{X}_i)}{e(\mathbf{X}_i)(1-e(\mathbf{X}_i))}$ being the Horvitz-Thompson (IPW) weights.
- \mathbf{X}_i^C is not required for identification, but contains optional functions of \mathbf{X}_i to reduce estimation noise, e.g. $[1, \mu(0, \mathbf{X}_i), e(\mathbf{X}_i), e(\mathbf{X}_i)\tau(\mathbf{X}_i)]$
- See Appendix A in [Chernozhukov et al. \(2017-2023\)](#) for a detailed derivation.

Classification Analysis (CLAN)

Classification Analysis (CLAN) can be implemented by simple mean comparisons of covariates in extreme GATES groups:

Definition “Classification Analysis (CLAN)”

Classification Analysis (CLAN) compares the covariate values of the **least affected group** G1 with the **most affected group** GK defined for the GATES:

- $\delta_K - \delta_1$

where

- $\delta_k = \mathbb{E}[X_i | G_k] = \frac{1}{n_k} \sum_{i=1}^n X_i \mathbb{1}[G_k].$

BLP, GATES & CLAN – Implementation

- R package [GenericML](#) by [Welz, Alfons, Demirer, and Chernozhukov \(2022\)](#).
 - **Algorithm:**
 - **IN:** Data = $(Y_i, \mathbf{X}_i, T_i)_{i=1}^N$, significance level α , a suite of ML methods, number of splits S .
 - **OUT:** p – values and $(1 - 2\alpha)$ confidence intervals of point estimates of each target parameter in GATES, BLP, and CLAN.
1. Compute propensity scores $e(\mathbf{X}_i)$.
 2. Do S splits of $\{1, \dots, N\}$ into disjoint sets A and M of same size.
 3. **for** each ML method and each split $s = 1, \dots, S$, **do**
 - a. Tune and train each ML method to learn $\mu(0, \mathbf{X}_i)$ and $\tau(\mathbf{X}_i)$ on A .
 - b. On M , use $\mu(0, \mathbf{X}_i)$ and $\tau(\mathbf{X}_i)$ to estimate the BLP, GATES, CLAN target parameters.
 - c. Compute some performance measures for the ML methods.
 4. Choose the best ML method based on the medians of the performance measures.
 5. Calculate the medians of the confidence bounds, p -values, and point estimates of each target parameter.
 6. Adjust the confidence bounds and p -values.

More References

- CATE Prediction Methods:
 - BART ([Hahn, Murray & Carvalho, 2020](#)).
 - Causal Boosting/MARS, ... ([Powers, Qian, Jung, Schuler, Shah, Hastie & Tibshirani, 2019](#)).
 - Dragonnet ([Shi, Blei & Veitch, 2019](#)).
 - Modified Causal Forest ([Lechner & Mareckova, 2022](#)).
 - Orthogonal Random Forest ([Oprescu, Syrgkanis & Wu, 2019](#)).
 - TARNet ([Shalit, Johansson & Sontag 2019](#)).
 - X-learner ([Künzel, Sekhon, Bickel & Yu, 2019](#)).
- HET Evaluation:
 - Rank-Weighted Average Treatment Effect (RATE) ([Yadlowsky et al., 2021](#)).
 - Calibration Error for Heterogeneous Treatment Effects ([Xu & Yadlowsky, 2022](#)).
 - More on GATES in experiments ([Imai & Li, 2022–2024](#)).

Optimal Policy Learning

Optimal Policy Learning – Goal

- From evaluation (What works for whom?) towards data-driven (personalized) treatment recommendations:
 - How to optimally treat whom?
- **Notation:**
 - Binary treatment indicator: $T_i \in \{0, 1\}$
 - Potential outcome (PO) under treatment t : $Y_i(t)$
 - Exogenous covariate(s): \mathbf{X}_i
 - Conditional Average PO: $\mu_t(\mathbf{x}) := \mathbb{E}[Y(t) \mid \mathbf{X}_i = \mathbf{x}]$
 - Conditional Average Treatment Effect (CATE): $\tau(\mathbf{x}) := \mu_1(\mathbf{x}) - \mu_0(\mathbf{x})$
- **Additional notation:**
 - Policy rule for \mathbf{x} (conditional treatment choice): $\pi(\mathbf{X}_i) \in \{0, 1\}$.
 - PO under policy $\pi(\mathbf{X}_i)$: $Y_i(\pi(\mathbf{X}_i))$.
 - Value function (average PO under policy $\pi(\mathbf{X}_i)$): $Q(\pi) := \mathbb{E}[Y_i(\pi(\mathbf{X}_i))]$.
- **Goal:** Find the optimal policy π^* that maximizes the value function $Q(\pi)$.

Optimal Policy Alternatives

- 1. Assign individuals to treatment with higher PO under treatment than without?
 - $\pi^* = \mathbb{1}[Y_i(1) > Y_i(0)] = \mathbb{1}[Y_i(1) - Y_i(0) > 0] = \mathbb{1}[\tau_i > 0]$
 - **Fundamental problem of causal inference:** counterfactuals unknown.
- 2. Assign individuals to treatment with higher CATE than without?
 - $\pi^* = \mathbb{1}[Y_i(1) > Y_i(0) | \mathbf{X}_i = \mathbf{x}] = \mathbb{1}[\tau(\mathbf{X}_i = \mathbf{x}) > 0]$
 - **Problem:** minimizing $\text{MSE}_{\text{CATE}} = \mathbb{E}[(\tau(\mathbf{x}) - \tau(\mathbf{x})^2)]$ does not necessarily improve downstream policy rule learning (Qian & Murphy, 2011).
 - Similar to the case where MSE minimization in treated and control groups separately is not the best strategy to minimize CATE MSE.
- 3. Instead: $\pi^* = \arg \min_{\pi} \mathbb{E}[Y_i(\pi(\mathbf{X}_i))] = \arg \min_{\pi} Q(\pi(\mathbf{X}_i))$

Optimal Policy Objective Function

- Objective function can have many different forms but one has proven very useful in the context of ML policy learning:
 - Comparing the value function against a **benchmark policy** that assigns treatments via **fair coin flip**:
 - 50-50 chance of being treated: $\pi^{\text{coin}} \sim \text{Bernoulli}(0, 5)$.

$$\begin{aligned}
 \pi^* &= \arg \max_{\pi} Q(\pi) = \arg \max_{\pi} \mathbb{E}[Y(\pi)] = \arg \max_{\pi} \mathbb{E}[Y(\pi) - 0.5\mathbb{E}[Y(1)] + 0.5\mathbb{E}[Y(0)]] \\
 &= \arg \max_{\pi} \mathbb{E}[\pi Y(1) + (1 - \pi)Y(0)] - 0.5\mathbb{E}[Y(1)] - 0.5\mathbb{E}[Y(0)] \\
 &= \arg \max_{\pi} \mathbb{E}[(\pi - 0.5)Y(1)] + \mathbb{E}[(0.5 - \pi)Y(0)] = \arg \max_{\pi} \mathbb{E}[(\pi - 0.5)(Y(1) - Y(0))] \\
 &= \arg \max_{\pi} 2\mathbb{E}[(\pi - 0.5)(Y(1) - Y(0))] \\
 &= \arg \max_{\pi} \mathbb{E}[(2\pi - 1)(Y(1) - Y(0))] \\
 &\stackrel{\text{LIE}}{=} \arg \max_{\pi} \mathbb{E}[(2\pi - 1)\tau(\mathbf{X}_i)] \\
 &= \arg \max_{\pi} \underbrace{\mathbb{E}[|\tau(\mathbf{X}_i)| \text{sign}(\tau(\mathbf{X}_i))]}_{A(\pi)} \underbrace{(2\pi(\mathbf{X}_i) - 1)}_{\text{}}
 \end{aligned}$$

- where $(2\pi(\mathbf{X}_i) - 1) \in \{-1, 1\}$ is one if policy assigns treatment and minus one if not.

Optimal Policy Objective Function – Intuition

- $A(\pi) := \mathbb{E}[|\tau(\mathbf{X}_i)|\text{sign}(\tau(\mathbf{X}_i))(2\pi(\mathbf{X}_i) - 1)]$ measures the **advantage** of a policy compared to random allocation:
 - If $\text{sign}(\tau(\mathbf{X}_i))(2\pi(\mathbf{X}_i) - 1) = 1$, i.e. if the policy picks the better treatment for \mathbf{X}_i , we **earn the absolute value of the CATE**.
 - If $\text{sign}(\tau(\mathbf{X}_i))(2\pi(\mathbf{X}_i) - 1) = -1$, i.e. if the policy picks the worse treatment for \mathbf{X}_i , we **lose the absolute value of the CATE**.
- We need to **get it right for those with biggest CATEs**, those with CATEs close to zero are negligible.
- This shows the **difference to CATE MSE minimization**, where we need to find good approximations everywhere.

Optimal Policy Identification & Estimation

- Potential outcomes or CATE functions unknown, need to be identified before optimization.
- **Athey and Wager (2021)** recommend the **pseudo-outcome** (again) because of all the nice properties:
 - $\tilde{\tau}_{i,ATE}^{AIPW} = \mu(1, \mathbf{X}_i) - \mu(0, \mathbf{X}_i) + \frac{T_i(Y_i - \mu(1, \mathbf{X}_i))}{\hat{e}_1(\mathbf{X}_i)} - \frac{(1-T_i)(Y_i - \mu(0, \mathbf{X}_i))}{\hat{e}_0(\mathbf{X}_i)}$
- **Binary weighted classification problem**: classify the sign of the CATE while favoring correct classifications with larger absolute CATEs.

$$\blacksquare \pi^* = \arg \max_{\pi \in \Pi} \left\{ \frac{1}{N} \sum_{i=1}^N \overbrace{|\hat{Y}_{i,ATE}|}^{\text{weight}} \underbrace{\text{sign}(\hat{Y}_{i,ATE})}_{\text{to be classified}} \overbrace{\frac{1}{(2\pi(\hat{X}_i) - 1)}}^{\text{function to be learned}} \right\}$$

- Possible methods: e.g. decision trees/forests, logistic lasso, SVM, etc.

Thank you for your attention!



 startupengineer.io/authors/ihl

 [christoph-ihl](https://www.linkedin.com/in/christoph-ihl)

 [christophihl](https://github.com/christophihl)

 [Ihluminate](https://twitter.com/Ihluminate)