

# Application of Machine Learning Algorithms to Detect Treatment Effect Heterogeneity for Three- Level Multisite Experiments

Wei Li, Walter Leite, & Jia Quan  
College of Education, University of Florida

April 13, 2024  
The AERA 2024 Conference

# Motivation (1)

- Multilevel randomized controlled trials (MRCTs) have been widely in education
  - Average treatment effect (ATE)
  - Heterogenous treatment effects (HTE)
    - Under what conditions for whom an intervention works
- Moderation analysis is traditionally used to evaluate HTE
  - An interaction between treatment and moderator (e.g., students' or schools' features)
  - Moderators are usually pre-specified – which covariates modify ATE?
- Recent development includes the use of machine learning (ML) methods to explore the HTEs
  - Identify subgroups with *significant effects*, post-analysis – what is the expected treatment effect for a group/individual with given characteristics
  - Select moderators from a potentially large number of covariates

## Motivation (2)

- MRCTs have nested data structures
  - E.g., students nested within teachers nested with schools
  - Cluster design – treatment at the school level
  - Block/multisite design – treatment at the student or teacher level
- Observations in the same clusters/sites are correlated rather than independent
  - Multilevel models, cluster robust SEs, bootstrap, etc. have been used to address the dependency
- Similarly, when applying ML methods to estimate CATE, applied researchers still need to consider the nested data structure
  - Most prior literature assumes the participants are independent
  - There is a lack of literature to guide educational researchers in appropriately applying ML methods for clustered data when evaluating HTEs

# Purpose

- Review the current available ML methods and tools that account for the nested data structure when explore HTEs
  - Focus on two ML methods – Cluster-Robust Causal Forest (CF) & Generic ML
- Demonstrate the application of these two methods using the dataset from a large multisite experimental study (Leite et al., 2023)
  - Provide recommendations to applied researchers on how to choose the appropriate methods and statistical package among alternative ML methods

# An Illustrative Example

- A large-scale multisite field experiment embedded within a virtual learning environment (VLE) for Algebra
  - Examine the effects of a video recommendation system
  - Students were randomly assigned to receive two types of video recommendations
    - Personalized video recommendations or generic recommendations
  - 2995 students nested within 54 teachers from 42 schools in three large school districts
  - Measures: 216 student- and teacher-level variables
    - Teacher survey – usage of VLE, instructional practice, perception of disruptions due to COVID, etc. (full survey available at <https://osf.io/h5tpn/>)
    - Student variables – gender, ethnicity, pretest score, absent days, etc.
    - Converted into 516 predictors, with 484 dummy-coded indicators
      - Some algorithm requires dummy-coding categorical variables

# Estimands of Interest

- Conditional average treatment effect (CATE)

$$\tau(x) = E[Y_{ij}(1) - Y_{ij}(0) | X = x] = u_1(x) - u_0(x) \quad (1)$$

- $X$  – a possibly high-dimensional vector of covariates
  - Require stable unit treatment values, unconfoundedness, and overlap
  - Best Linear Predictor (BLP) of CATE
  - Note that,  $ATE = \tau = E[\tau(x)]$ , site average or individual average; once we know CATE, we immediately know ATE
- Non-overlapping subgroup analysis
    - E.g., Sorted Group Average Treatment Effects (GATES)

$$\tau(G_k) = E[Y_{ij}(1) - Y_{ij}(0) | G_k] \quad (2)$$

- Moderator effect

$$\tau(m) = E[Y_{ij}(1) - Y_{ij}(0) | M = m] \quad (3)$$

- $M$  is a subset of predictors of  $X$

# Models: OLS with Teacher Fixed Effects and Interactions

- Traditional Moderator Analysis: Interaction approach

$$y_{ij} = \gamma_0 + \gamma_1 T_{ij} + \gamma_2 M_{ij} + \gamma_3 T_{ij} M_{ij} + u_j + T_{ij} u_j + r_{ij} \quad (4)$$

- $y_{ij}$  - test score for student  $i$  in teach  $j$
- $T_{ij}$  - treatment indicator
- $M_{ij}$  - student-level moderator
- $u_j$  - teacher dummy variables
  - MLM - teacher-level random effect that follows a normal distribution
- $r_{ij}$  - level-1 error
- **$\gamma_3$  - moderator effects, not a causal effect** ([Dong et al., 2022](#))
- $\tau(M_{ij}) = \gamma_1 + \gamma_3 M_{ij}$ , **a causal effect**
- Assume clusters have an additive effect on the outcome
  - Same functions for all sites

## Models: S/T-learner

- Fit (*separate*) models to the treatment and control groups

$$\begin{aligned}E(Y|T = 1, X = x) &= f_1(x) \\E(Y|T = 0, X = x) &= f_0(x)\end{aligned}$$

- Then, CATE is

$$\begin{aligned}\tau(x) &= f_1(x) - f_0(x) \\Var(\hat{\tau}(x)) &= Var(\hat{f}_1(x)) + Var(\hat{f}_0(x))\end{aligned}$$

- To our best knowledge, no methods or tools consider the nested data structure for S/T-learner



# Models: Cluster-Robust RF

- Cluster-robust RF (Athey & Wager, 2019)

$$y_{ij} = \alpha_j(x) + \tau_j(x)T_{ij} + e_{ij}, \tau(x) = E[\tau_j(x)] \quad (5)$$

- $\alpha_j(x)$  - control group average for site (e.g., teacher)  $j$
- $\tau_j(x)$  - CATE in site (e.g., teacher)  $j$
- $\tau(x)$  - CATE across sites; site average
  - Give each cluster/site equal weight
  - Accurate for predicting effects on a new student from a new site
- Each cluster has its own main ( $\alpha_j(x)$ ) and treatment effect function ( $\tau_j(x)$ )

# Models: Generic ML

- Generic ML (Chernozhukov et al., 2023)

$$y_{ij} = b_0(x) + T_{ij}s_0(x) + e_{ij}, \tau(x) = s_0(x) \quad (6)$$

- $b_0(x) = E[y_{ij}|T_{ij} = 0, x]$  – baseline conditional average; mean for the control group across sites
- $s_0(x)$  - CATE across level-1 units (e.g., students); individual average
- Site dummy variables can be included when estimating  $b_0(x)$  and  $s_0(x)$
- An application of double/debiased ML
  - Utilize Neyman orthogonal moments and cross-fitting to address regularization bias and overfitting
- Focus on key features of CATE instead of CATE: e.g., BLP of CATE
  - **Sparsity** -  $s_0$  can be well-approximated by a function that only depends on a low-dimensional subset of  $X$

# Why/How to Consider Nested Data Structure

- In general, ML methods for CATE include three main steps:
  - Splitting the data into training and test sets – **cluster-based split?**
  - Use the training set and ML algorithms to build a prediction model – **will considering cluster membership improve prediction?**
  - Use the test set to estimate CATE/BLP and their standard errors (SEs) – **should we use cluster-robust SEs or something similar?**
- Based on our review of all the currently available methods and packages, only two algorithms – cluster-robust CF and the GenericML consider the nested data structure in at least one step
  - Not consider alternative methods, e.g., Bayesian additive regression trees (BART), Targeted MLE, Meta learners (e.g., S-, T- learners)

# How does Cluster-Robust CF Address the Nesting Effects?

- The cluster-robust CF algorithm considers the nested data structure in *all* three steps:
  - (1) for each  $b = 1, \dots, B$ , draw a **subsample of clusters** and then draw a random sample from each cluster as the training data;
  - (2) grow a tree via recursive partitioning on each such subsample of the data;
  - (3) make the out-of-bag predictions: to account for the potential within cluster dependency, **an observation  $i$  is considered to be out-of-bag if its cluster was not drawn in step (1)**
- Implemented through *grf* R package
  - SE of CATE – jackknife SE
  - Report cluster-robust SEs for BLP

# How does the GenericML algorithm Address the Nesting Effects? (1)

- The GenericML algorithm (Chernozhukov et al., 2023) estimates the best linear predictor (BLP) of CATE through the following steps:
  - (1) randomly split the data into training and test sets; **without consideration of clusters**
  - (2) estimates the CATE with any number of selected ML methods (e.g., random forest) using the training data; *can potentially consider clustering effects*
  - (3) use OLS regression to obtain the BLP of the CATE using the test data; **include site fixed effects (dummy variables or demean); easy to report cluster-robust SE from OLS estimation;**
- Note that
  - Random forest – Build B trees which place covariate splits that maximize the squared difference **in subgroups means**
  - Causal forest - Greedily places covariate splits that maximize the squared difference in **subgroup treatment effects**

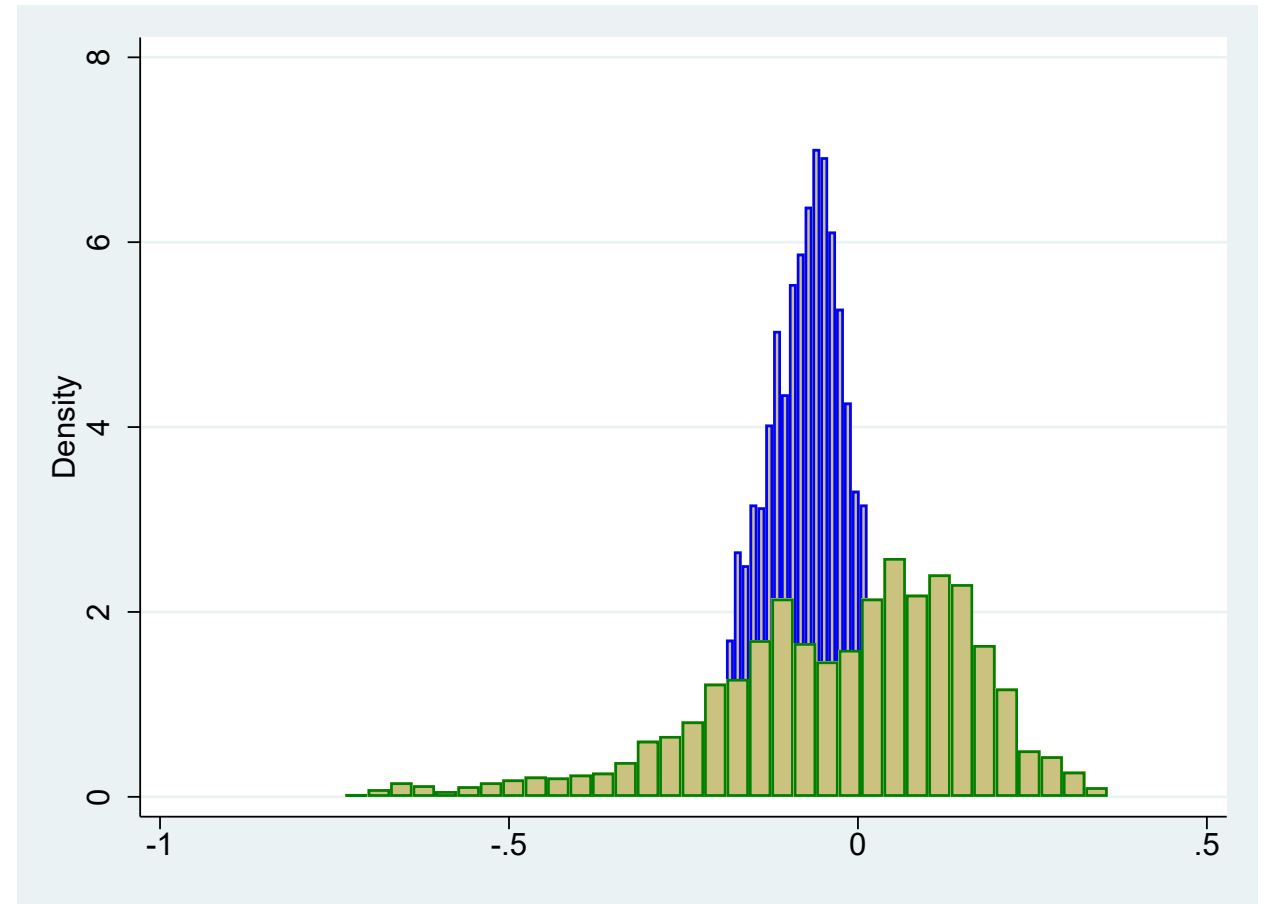
## How does the GenericML algorithm Address the Nesting Effects? (2)

- Implemented through the *GenericML* R package
  - Estimate the Sorted group average treatment effects (GATEs): creating five groups of participants using quintiles of the CATE distribution
  - Perform classification analysis (CLAN) to explore the relationships between covariates and the CATE
  - Report cluster-robust SEs for BLP, GATEs, and CLAN
    - OLS estimation easy to deal with clustering

# Results: Cluster-Robust RF

Method	ATE	SE
CF w/o clustering	-0.029	0.029
Cluster-Robust CF	-0.058	0.044

Note. We used lasso, elastic net, support vector machine, XGBoost, and random forests (RF). RF is the best learner.



Blue: CATE estimates from cluster-robust CF  
Yellow: CATE estimates from CF w/o clustering

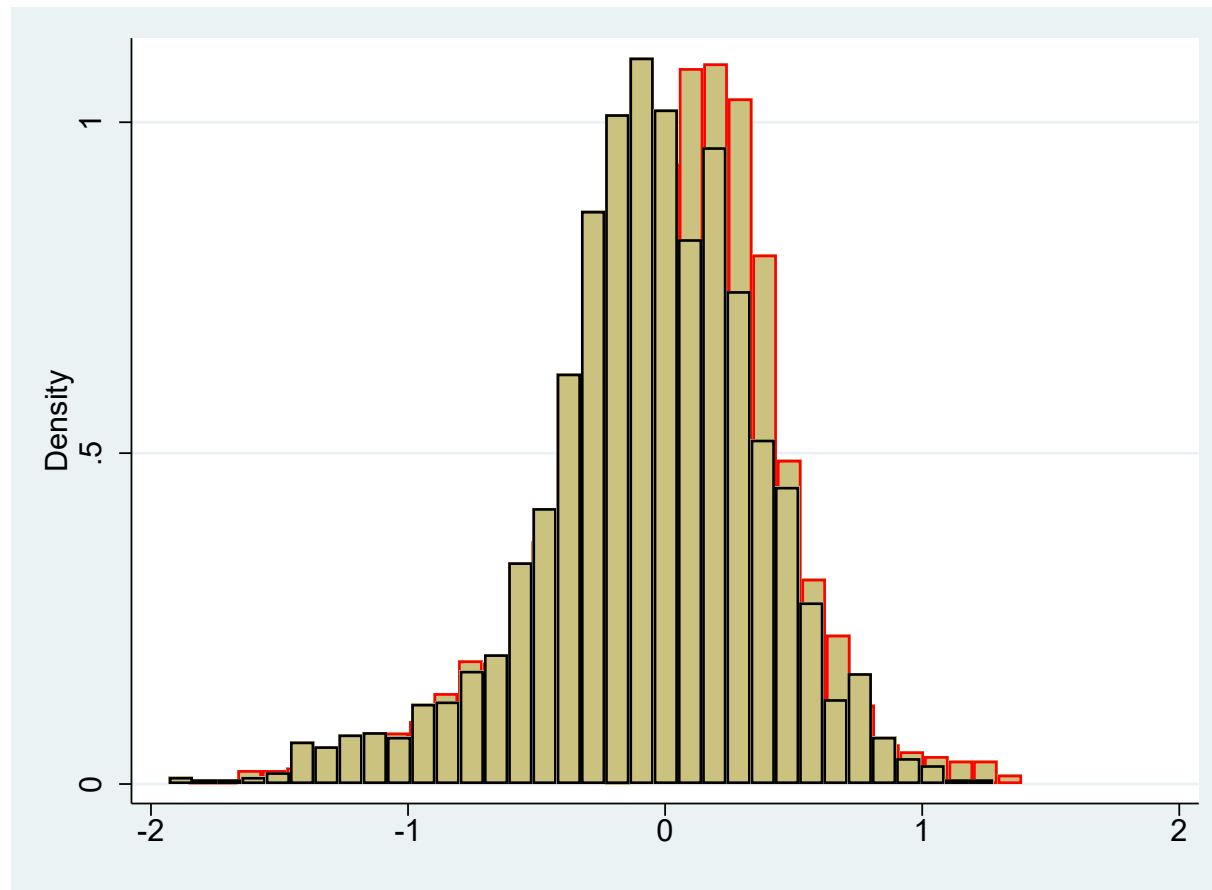
# Results: GenericML (RF)

GenericML		Estimate	P-Value
With Teacher Fixed Effects	ATE	-0.030	0.434
	Treatment Heterogeneity	0.917	0.000
Without Teacher Fixed Effects	ATE	-0.025	0.535
	Treatment Heterogeneity	0.976	0.000

GenericML	GCATE	Estimate	P-Value
Without Teacher Fixed Effects	Group 1	-0.657	0.000
	Group 2	-0.199	0.027
	Group 3	0.003	0.955
	Group 4	0.175	0.051
	Group 5	0.525	0.000
	Group 5 - Group 1	1.179	0.000
With Teacher Fixed Effects	Group 1	-0.627	0.000
	Group 2	-0.185	0.030
	Group 3	-0.005	0.952
	Group 4	0.167	0.056
	Group 5	0.482	0.000
	Group 5 - Group 1	1.099	0.000



# Results: BLP of CATE from GenericML (RF)



Black: BLP of CATE estimates without teacher fixed effects

Red: BLP of CATE estimates with teacher fixed effects

# Discussion and Conclusion

- Clustering effects should be considered when estimating CATE
  - No perfect solution now
  - It seems CF is preferred
- Differences of CATE estimates between cluster-robust CF and Generic ML
  - CATE or BLP of CATE
  - Site average or individual average
- Future directions
  - Alternative methods: BART, R-learners, TMLE
  - Simulations study
  - Cluster design

Questions or Comments?

Thank you!  
[wei.li@coe.ufl.edu](mailto:wei.li@coe.ufl.edu)

# Appendix: ML and CATE

# CR-CF

$$\hat{\tau}_j = \frac{1}{n_j} \sum_{\{i:A_i=j\}} \hat{\Gamma}_i, \quad \hat{\tau} = \frac{1}{J} \sum_{j=1}^J \hat{\tau}_j, \quad \hat{\sigma}^2 = \frac{1}{J(J-1)} \sum_{j=1}^J (\hat{\tau}_j - \hat{\tau})^2, \quad (8)$$
$$\hat{\Gamma}_i = \hat{\tau}^{(-i)}(X_i) + \frac{Z_i - \hat{e}^{(-i)}(X_i)}{\hat{e}^{(-i)}(X_i)(1 - \hat{e}^{(-i)}(X_i))} \left( Y_i - \hat{m}^{(-i)}(X_i) - \left( Z_i - \hat{e}^{(-i)}(X_i) \right) \hat{\tau}^{(-i)}(X_i) \right).$$

# BLP of CATE (Chernozhukov, 2018)

Conditional average treatment effect (CATE):

$$g_0(1, X) - g_0(0, X)$$

- captures (learnable) heterogeneity in treatment effects under unconfoundedness
- generally high-dimensional nonparametric object - inference impractical (impossible?)

Another potential summary is best linear predictor (BLP) of CATE given pre-specified (low-dimensional) vector  $W$

- other summaries possible; Semenova and Chernozhukov (2021)

Inference for BLP is possible using orthogonal score for ATE

6.1. **Implementation Algorithm.** We describe an algorithm based on the first identification strategy and provide some specific implementation details for the empirical example.

**Algorithm 1 (Inference Algorithm).** The inputs are given by the data on units  $i \in [N] = \{1, \dots, N\}$ .

Step 0. Fix the number of splits  $S$  and the significance level  $\alpha$ , e.g.  $S = 100$  and  $\alpha = 0.05$ .

Step 1. Compute the propensity scores  $p(Z_i)$  for  $i \in [N]$ .

Step 2. Consider  $S$  splits in half of the indices  $i \in \{1, \dots, N\}$  into the main sample,  $M$ , and the auxiliary sample,  $A$ . Over each split  $s = 1, \dots, S$ , apply the following steps:

- Tune and train each ML method separately to learn  $B(\cdot)$  and  $S(\cdot)$  using  $A$ . For each  $i \in M$ , compute the predicted baseline effect  $B(Z_i)$  and predicted treatment effect  $S(Z_i)$ . If there is zero variation in  $B(Z_i)$  and  $S(Z_i)$  add Gaussian noise with small variance to the proxies, e.g., a 1/20-th fraction of the sample variance of  $Y$ .
- Estimate the BLP parameters by weighted OLS in  $M$ , i.e.,

$$Y_i = \hat{\alpha}' X_{1i} + \hat{\beta}_1(D_i - p(Z_i)) + \hat{\beta}_2(D_i - p(Z_i))(S_i - \mathbb{E}_{N,M} S_i) + \hat{\epsilon}_i, \quad i \in M$$

such that  $\mathbb{E}_{N,M}[w(Z_i)\hat{\epsilon}_i X_i] = 0$  for  $X_i = [X'_{1i}, D_i - p(Z_i), (D_i - p(Z_i))(S_i - \mathbb{E}_{N,M} S_i)]'$ , where  $w(Z_i) = \{p(Z_i)(1 - p(Z_i))\}^{-1}$  and  $X_{1i}$  includes a constant,  $B(Z_i)$  and  $S(Z_i)$ .

- Estimate the GATES parameters by weighted OLS in  $M$ , i.e.,

$$Y_i = \hat{\alpha}' X_{1i} + \sum_{k=1}^K \hat{\gamma}_k \cdot (D_i - p(Z_i)) \cdot 1(S_i \in I_k) + \hat{v}_i, \quad i \in M,$$

such that  $\mathbb{E}_{N,M}[w(Z_i)\hat{v}_i W_i] = 0$  for  $W_i = [X'_{1i}, \{(D_i - p(Z_i))1(S_i \in I_k)\}_{k=1}^K]'$ , where  $w(Z_i) = \{p(Z_i)(1 - p(Z_i))\}^{-1}$ ,  $X_{1i}$  includes a constant,  $B(Z_i)$  and  $S(Z_i)$ ,  $I_k = [\ell_{k-1}, \ell_k)$ , and  $\ell_k$  is the  $(k/K)$ -quantile of  $\{S_i\}_{i \in M}$ .

- Estimate the CLAN parameters in  $M$  by

$$\hat{\delta}_1 = \mathbb{E}_{N,M}[g(Y_i, Z_i) | S_i \in I_1] \quad \text{and} \quad \hat{\delta}_K = \mathbb{E}_{N,M}[g(Y_i, Z_i) | S_i \in I_K],$$

where  $I_k = [\ell_{k-1}, \ell_k)$  and  $\ell_k$  is the  $(k/K)$ -quantile of  $\{S_i\}_{i \in M}$ .

# Group Average Treatment Effects (GATEs)

Group average treatment effects (GATEs):

- Let  $G$  be an indicator for belonging to some group of interest (e.g. an education category)
- $\text{GATE} = E[g_0(1, X) - g_0(0, X) | G = 1]$
- Can use to summarize heterogeneity along pre-specified directions of interest
- Average treatment effect on the treated (ATET) is a special case

For  $\psi_1(Y, Z, X)$  defined above, orthogonal moment for GATE is

$$E \left[ \frac{G}{p_G} \psi_1(Y, Z, X) \right] = 0$$

- Nuisance functions:  $E[Z|X] = m_0(X)$ ;  $E[Y|Z, X] = g_0(Z, X)$ ;  $p_G = E[G]$

1. Partition sample indices into random folds of approximately equal size:  $\{1, \dots, n\} = \cup_{k=1}^K l_k$ . For each  $k = 1, \dots, K$ , compute estimators  $\hat{p}_{[k]}$ ,  $\hat{g}_{[k]}$ , and  $\hat{m}_{[k]}$  of  $E[G]$  and the conditional expectation functions  $g_0(Z, X) = E[Y|Z, X]$  and  $m_0(X) = E[Z|X]$  leaving out the  $k^{\text{th}}$  block of data and enforcing  $\epsilon \leq \hat{m}_{[k]} \leq 1 - \epsilon$ .
2. For each  $i \in l_k$ , let

$$\hat{\psi}(Y_i, Z_i, X_i, G_i; \alpha) = \frac{G_i}{\hat{p}_{[k]}} \left( \hat{g}_{[k]}(1, X_i) - \hat{g}_{[k]}(0, X_i) + \frac{Z_i(Y_i - \hat{g}_{[k]}(1, X_i))}{\hat{m}_{[k]}(X_i)} - \frac{(1 - Z_i)(Y_i - \hat{g}_{[k]}(0, X_i))}{1 - \hat{m}_{[k]}(X_i)} \right) - \frac{G_i}{\hat{p}_{[k]}} \alpha.$$

Compute the estimator  $\hat{\alpha}$  as the solution to  $E_n[\hat{\psi}(W_i; \alpha)] = 0$  which yields

$$\hat{\alpha} = \frac{E_n \left[ \frac{G_i}{\hat{p}_{[k]}} \left( \hat{g}_{[k]}(1, X_i) - \hat{g}_{[k]}(0, X_i) + \frac{Z_i(Y_i - \hat{g}_{[k]}(1, X_i))}{\hat{m}_{[k]}(X_i)} - \frac{(1 - Z_i)(Y_i - \hat{g}_{[k]}(0, X_i))}{1 - \hat{m}_{[k]}(X_i)} \right) \right]}{E_n \left[ \frac{G_i}{\hat{p}_{[k]}} \right]}$$

3. Let

$$\hat{\varphi}(Y_i, Z_i, X_i, G_i) = \frac{\hat{\psi}(Y_i, Z_i, X_i, G_i; \hat{\alpha})}{E_n \left[ \frac{G_i}{\hat{p}_{[k]}} \right]}$$

Construct standard errors via

$$\sqrt{\hat{V}/n}, \quad \hat{V} = E_n[\hat{\varphi}(Y_i, Z_i, X_i, G_i)^2]$$

and use standard normal critical values for inference.

- Variable importance

- CF-noncluster: [1] "pretest" "Q34\_1\_feb" "Q13\_8\_feb\_2"  
"Q25\_apr" "EOC\_scale\_score"

- Cluster-robust CF: [1] "pretest" "EOC\_scale\_score" "absent\_days"  
"EOC\_achieve\_level" "mean\_num\_received"



# Permutation Importance

Variable Name	Rank	Survey Item	Parent Code
Q10_7_feb	1	How frequently did you check each of these Algebra Nation reports during the past month? - Video recommendation views	Fidelity of Implementation
Q93_13_may	2	Thinking about your ability to provide high-quality instruction during Spring 2021, how challenging do you find: - Balancing personal and work life	Organizational, Personal?
years_teaching	3	How many years have you been teaching (not including this current school year)? (a variable from Feb teacher survey)	Experience
Q20_2_feb	4	During the past month, did you use Algebra Nation Check Your Understanding quizzes using any of the following methods? - Assigned to groups/centers.	Fidelity of Implementation
Q69_apr	5	When a low-achieving child progresses in mathematics, it is usually due to extra attention given by me.	Teacher Efficacy