

Optimizing Random Forests Propensity Score

HEINING CHAM, LANDON HURLEY, YUE TENG

MAY 24, 2016

2016 MODERN MODELING METHODS CONFERENCE

Outline

1. Random Forests Model Specifications
2. Simulation Study
3. Results
4. Discussion

Propensity Scores

$$e(\mathbf{X}_i) = \Pr(Z_i = t \mid \mathbf{X}_i)$$

1. Conditional probability that participant i is assigned to treatment group (t) given \mathbf{X}_i .
2. Coarsest function (one variable summary) of \mathbf{X}_i to equate the distributions of \mathbf{X}_i between treatment and control groups.

Propensity Scores : Properties

1. In practice, propensity score is unknown and needs to be estimated.
2. Incorrect propensity score estimation model produces biased average treatment effect (ATE) or average treatment effect on the treated (ATT) estimates (Drake, 1993).
3. We might not have a sufficient theoretical or empirical basis to specify the propensity score estimation model.

Using Random Forests to Estimate Propensity Scores

1. Random Forests is an **automatic** and **nonparametric** method to deal with regression problem with (1) many covariates, and (2) complex nonlinear and interaction effects of the covariates.
2. Austin (2012) and Lee, Stuart, and Lessler (2010) have investigated the performance of Random Forests for propensity score analysis.

Goal of Study

1. Austin (2012) and Lee et al. (2010) did not systematically investigate the effects of different Random Forests model specifications.
2. Here, we investigate the effects of different random forests model specifications on propensity score analysis.

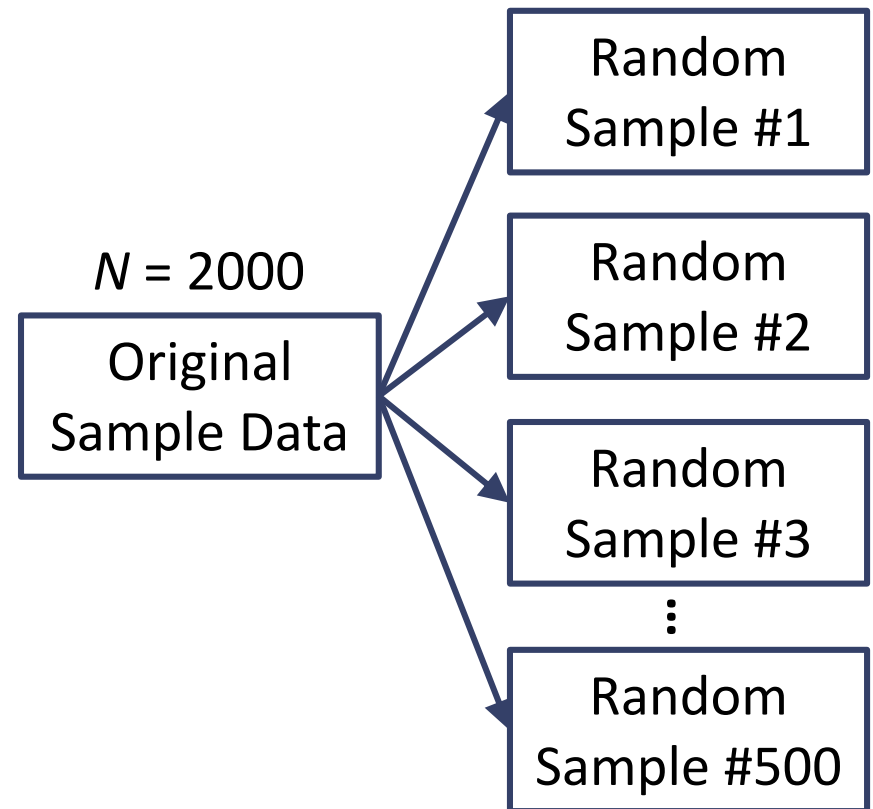
Step 1: Draw Multiple Random Samples

Strobl, Boulesteix, Zeileis, et al.
(2007) suggest:

1. Sampling without replacement
2. Random samples which are 0.632 times the sample size of the original data

This specification reduces the covariate selection bias towards covariates with many categories and continuous covariates in Random Forests.

$$N = 2000 \times 0.632 = 1264$$



Step 2:

Estimate Classification Tree Model In Each Sample

Depth of this tree = 3

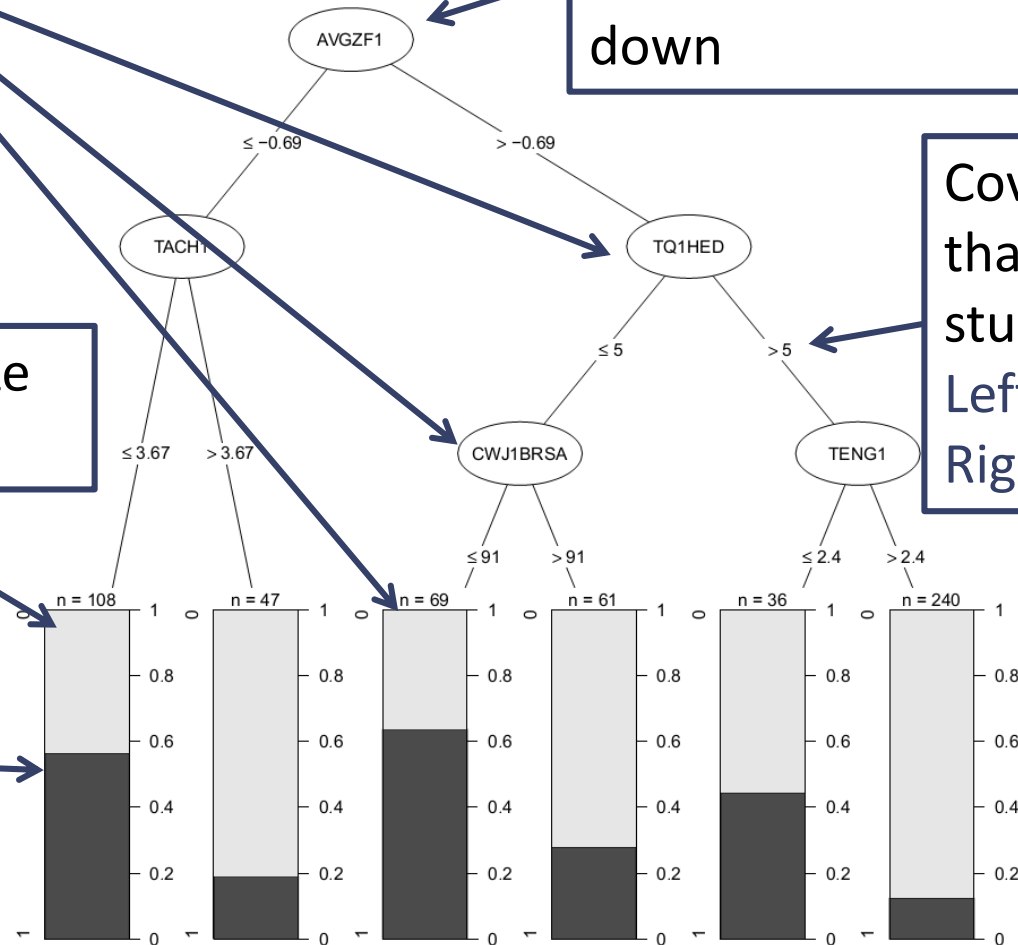
Node of covariate that classifies students into two levels bottom-down

Covariate Value that classifies students:
Left: ≤ 0.5
Right: > 0.5

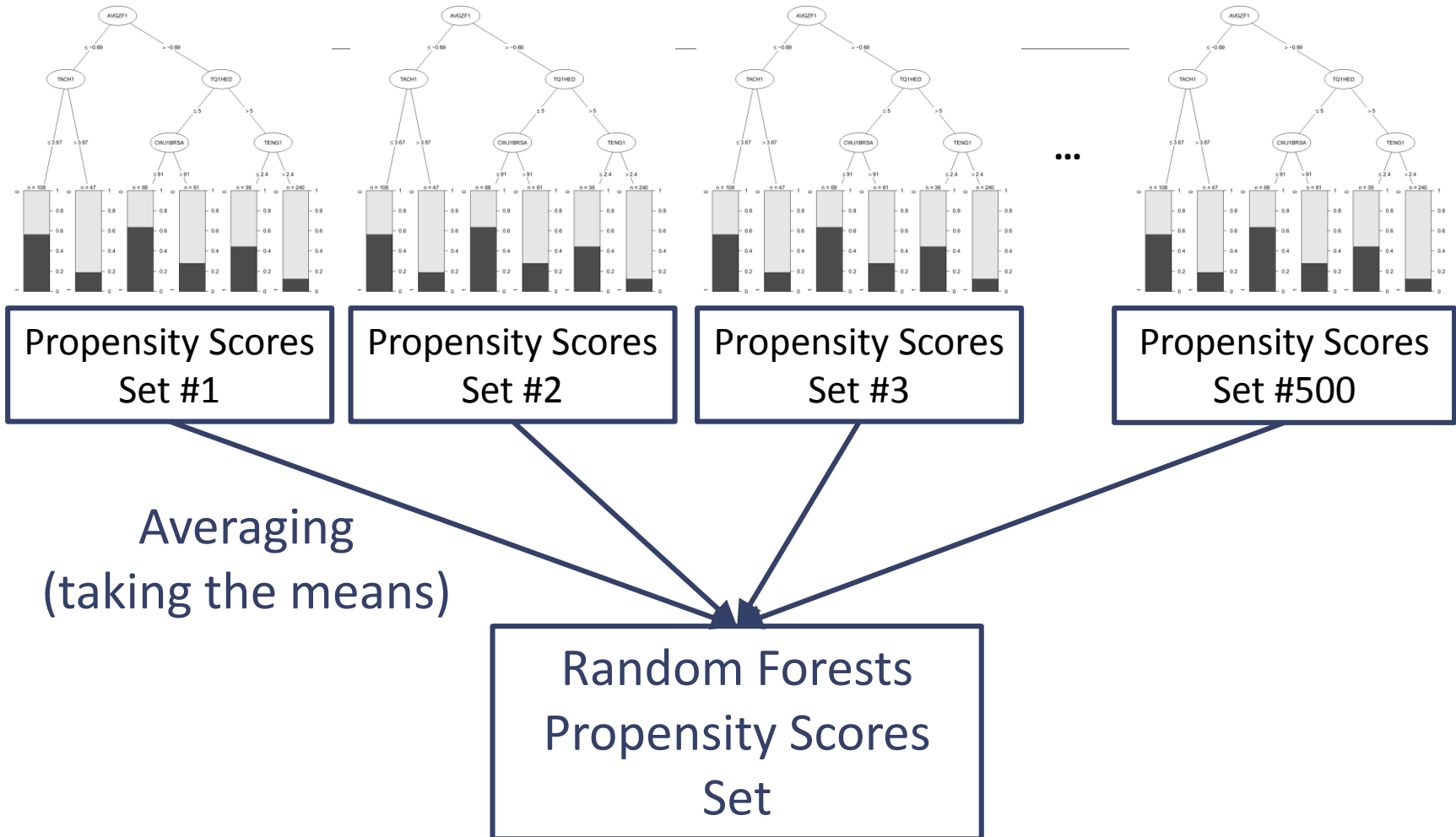
Terminal Node Size
($n = 108$)

Propensity Scores:

Proportion of retained students (dark bar) in the node



Step 3: Average All Classification Tree Propensity Scores Sets



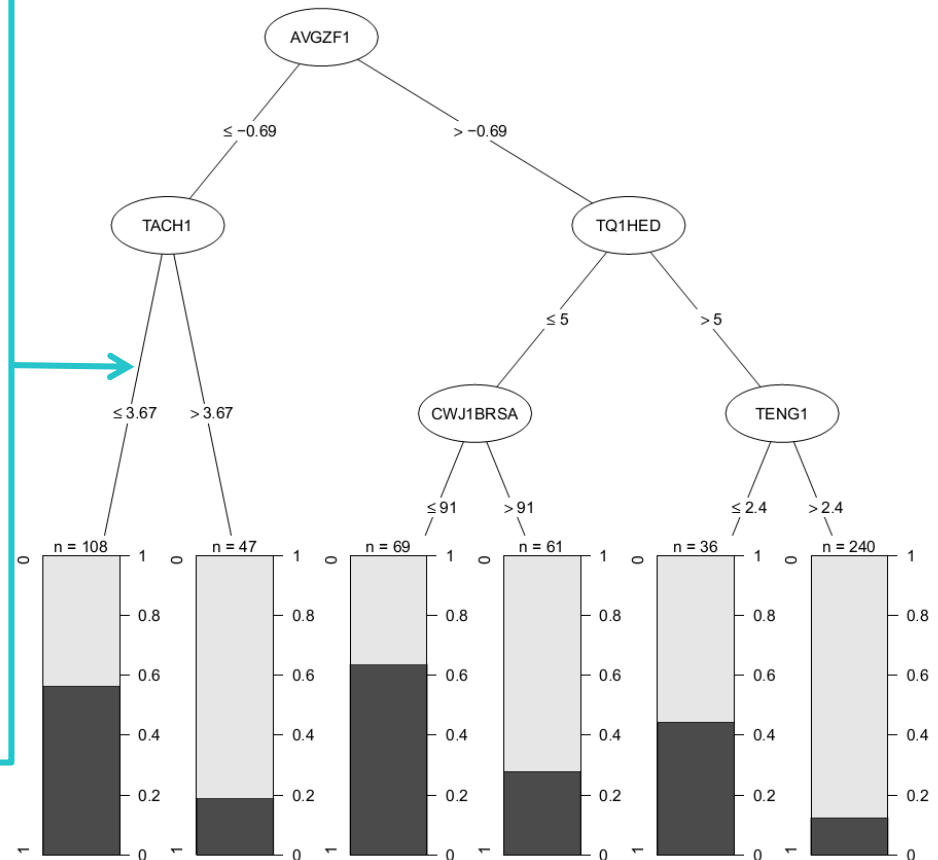
Model Specifications

#1: Decision rule to select the covariate and its value

(1) Gini Index

(2) Conditional Significance Test

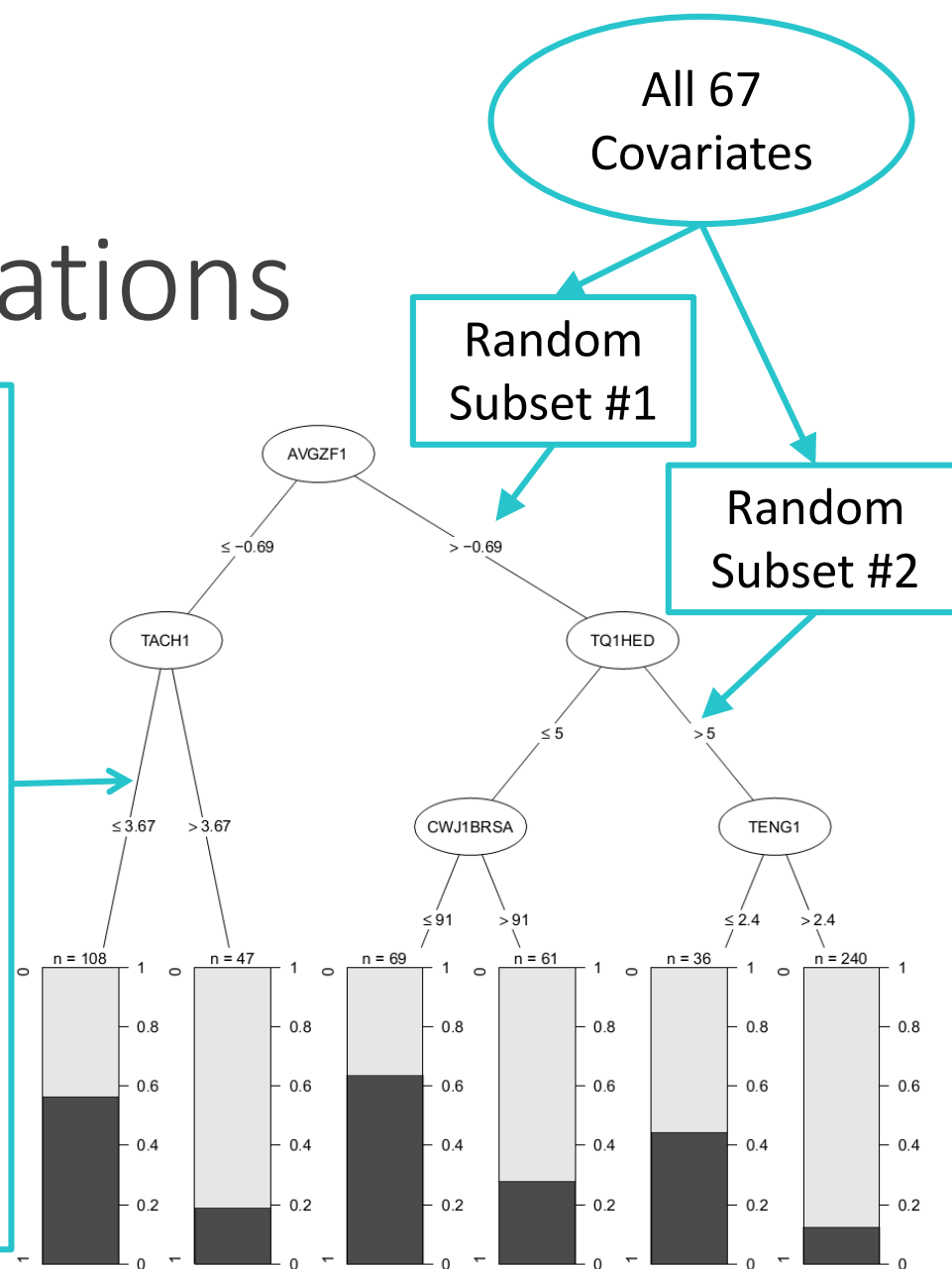
(2) is advantageous relative to (1) for reducing selection bias towards continuous and categorical covariates with many levels.



Model Specifications

#2: Random subset of covariates
Random subset is advantageous relative to all covariates for

- (a) Reducing sampling uncertainty of Random Forests propensity scores
- (b) Selecting covariates that are relatively less associated with grade retention but more associated with other covariates and the outcome



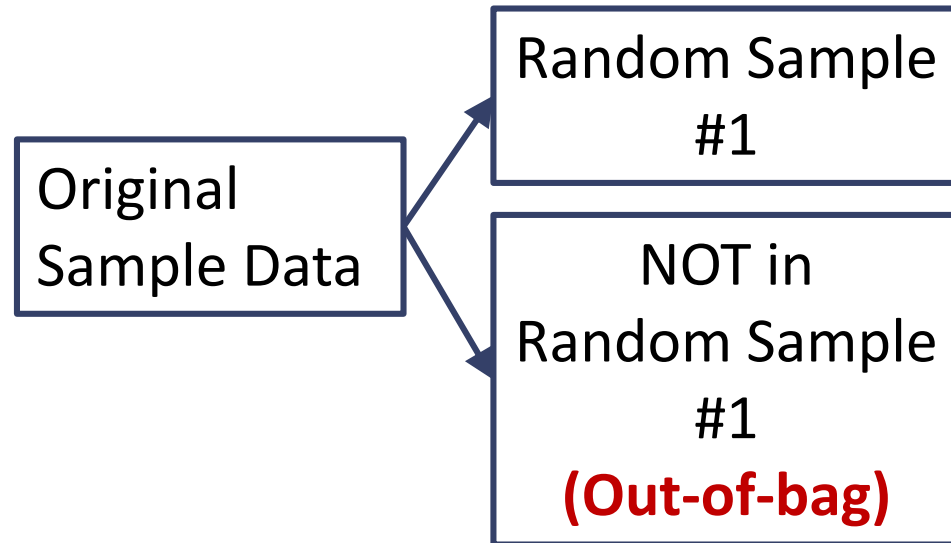
Model Specifications

#3: Data to calculate
Classification Tree
propensity scores

**(1) Full Sample Data
(Original Data)**

**(2) Data NOT in Random Sample
(Out-of-bag sample)**

**Out-of-bag sample maybe
advantageous relative to full
sample** for less biased average
treatment effect estimate.



Summary

1. **Conditional Significance Test** (vs. Gini Index)
 2. **Random Subset of Covariates** (vs. All Covariates / No Sampling)
 3. **Out-of-bag Sample** (vs. Original Sample / Full Sample)
- ✓ Specifications **in red** produce less biased **average treatment effect of the treated (ATT)** estimates.
 - ✓ The combination of these specifications **in red** will be optimal to produce the least biased **ATT** estimate.

Simulation Study Design (1)

Constructed loosely based on Im, Hughes, Kwok, Puckett, and Cerda (2013)

1. Covariates

- a) Covariate Types: **16 binary, 40 standard normal, 8 ordered-categorical (0 to 6)**
- b) Covariate Correlations: **Low and High**

2. Propensity Score Model (Logistic Regression)

- a) **Linear** (in the logit metric) and **Nonlinear** (added interaction and quadratic effects)
- b) Magnitude of Regression Coefficients: **Low and High**

3. Treatment-Outcome Model (Linear Regression)

- a) Magnitude of *ATT*: **Zero and Non-zero (moderate effect size)**
- b) Magnitude of Regression Coefficients: **Low and High**

4. Sample Sizes : **600 and 2000**

Simulation Study Design (2)

4. Benchmark Methods
 - **Uncorrected**, **ANCOVA** (True), **Logistic** Regression Propensity Score (True)
5. Decision Rule to Select Covariate and its Value
 - **Gini** Index, Conditional Significance Test (**Sig**)
6. Random Sampling of Covariates for Selection
 - No Sampling (**NS**; All Covariates), 8 Covariates (**S8**), 4 Covariates (**S4**)
7. Methods for Estimating Propensity Scores
 - Full Sample (**F**; Original Sample), Out-of-bag Sample (**O**)
8. Methods of Equating Groups on Propensity Scores:
 - Nearest Neighbor Matching (**Matching**), Weighting by Odds (**Weighting**)
 - Both methods estimate **ATT**.

Estimation results were consistent across

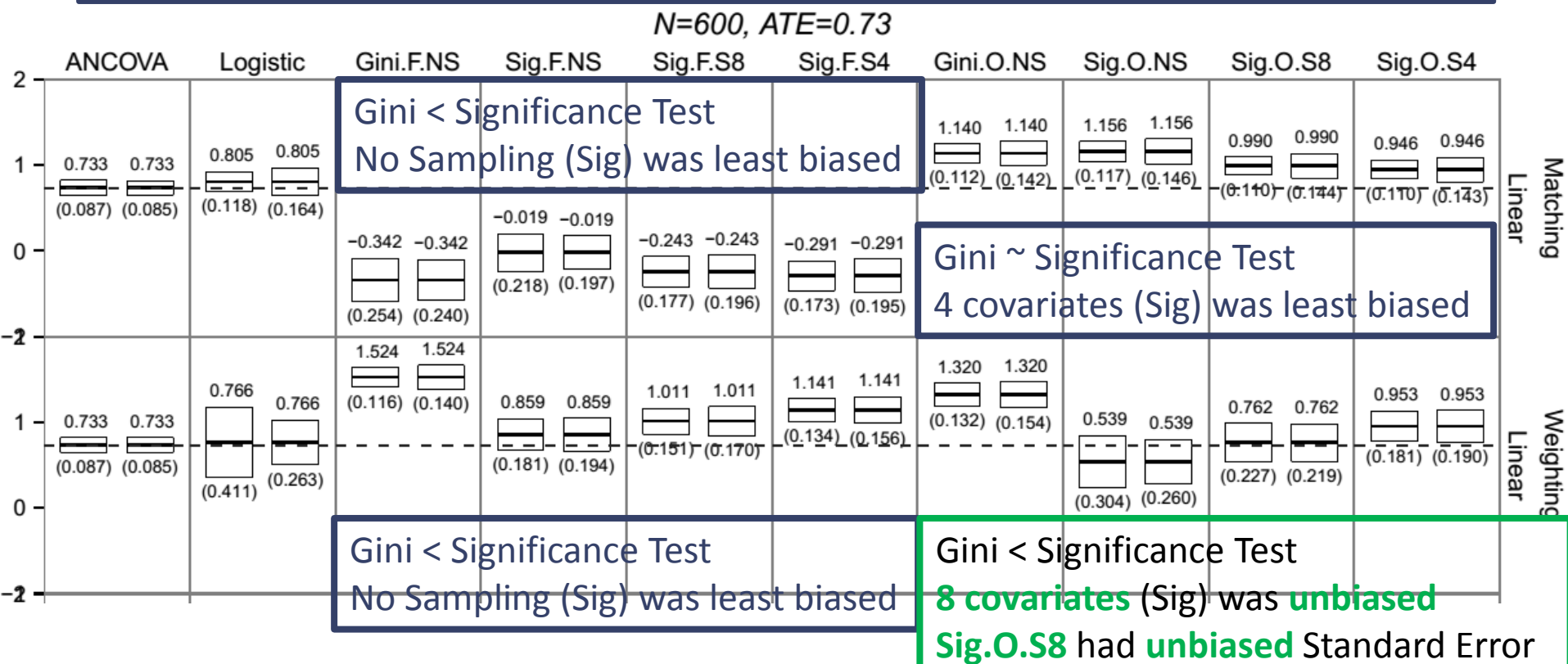
- Linear and Nonlinear Propensity Score Model
- Magnitude of Average Treatment Effect (Zero and Non-zero)

ANCOVA produced **unbiased** ATT, **unbiased** and **smallest** Standard Error

Logistic Regression Propensity Scores produced **unbiased** ATT estimate

Matching: Standard Error was **overestimated**

Weighting: Standard Error was **underestimated**



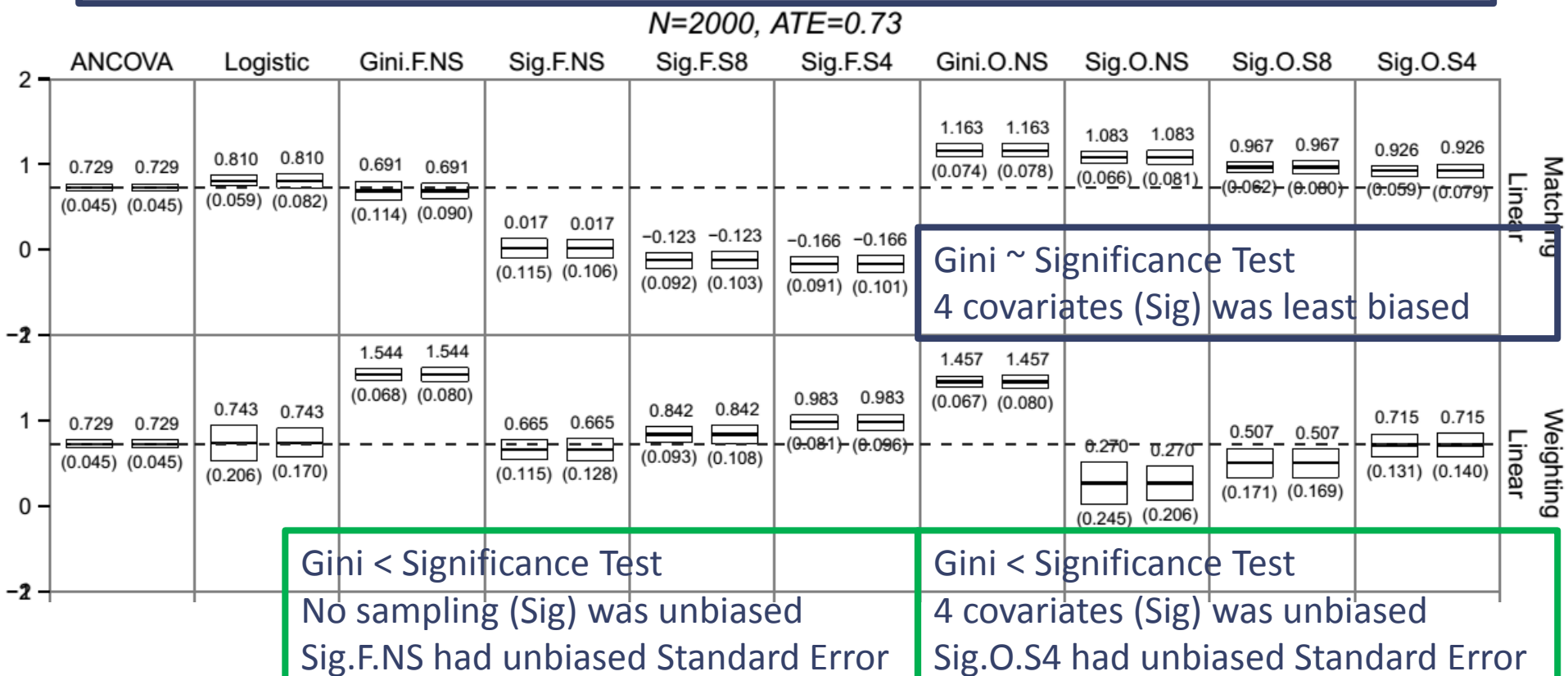
$N = 600$: Sig.O.S8 (Weighting)

$N = 2000$: Gini.F.NS (Matching), Sig.F.NS (Weighting),
Sig.O.S4 (Weighting)

Gini > Significance Test

No Sampling (Sig) was least biased

Gini.F.NS had **unbiased** ATT, **underestimated** Standard Error



1. $N = 600$: **Sig.O.S8 (Weighting)** had **slightly inflated α** .
2. $N = 2000$: **Gini.F.NS (Matching)** had **inflated α** .

Sig.F.NS and **Sig.O.S4 (Weighting)** had **correct α** .

$N = 600$, Linear	Coverage Rate	α	Power	
ANCOVA	0.942	0.041	1.000	Correct
Logistic (Matching)	0.983	0.025	1.000	Conservative
Logistic (Weighting)	0.803	0.187	0.783	Inflated
Sig.O.S8 (Weighting)	0.922	0.087	0.869	Slightly Inflated
$N = 2000$, Linear	Coverage Rate	α	Power	
ANCOVA	0.949	0.053	1.000	
Logistic (Matching)	0.905	0.073	1.000	
Gini.F.NS (Matching)	0.863	0.142	1.000	
Logistic (Weighting)	0.900	0.115	0.902	
Sig.F.NS (Weighting)	0.964	0.035	0.996	
Sig.O.S4 (Weighting)	0.950	0.040	0.962	

Discussion (1)

Nearest Neighbor Matching – Hypotheses **NOT** supported.

Some Explanations:

1. Model specifications were **NOT** optimal for matching.
2. Matching setting is sensitive to propensity score model misspecifications (Zhao, 2008).
 - Especially with matching without replacement of control group participants.

Discussion (2)

Weighting by Odds – More hypotheses supported.

1. The **optimal** specification:
 - **(a) Conditional Significance Test**
 - **(b) Random Covariates Subset**
 - **(c) Out-of-bag Sample**
2. Number of covariates in subset was **SENSITIVE** to sample size.
3. Lee et al.'s (2010) results showed specification was insensitive to number of covariates in subset
 - Their model had fewer covariates, and more covariate pairs were uncorrelated.
4. **↑ Total Number of Covariates, ↑ Sensitivity of Number of Covariates in Subset**

Discussion (3)

Austin (2012) found Random Forests had biased estimates.

Some Explanations:

1. Austin investigated a different the weighting for to estimate average treatment effect (ATE), not ATT
2. When estimating ATE, random forests may require different specifications to produce optimal specifications.
3. Weighting method is sensitive to propensity score model specification.

Follow-up Question:

Can Absolute Standardized Mean Difference (ASMD) search the optimal model specification?

$$ASMD = |\bar{X}_t - \bar{X}_c| / s_t^b$$

↑ standardized mean difference *between* different propensity scores, **did not** necessarily relate to ↓ bias of ATT estimate.

1. **Between** different random forests model specifications, does standardized mean difference of covariates relate to the magnitude of bias of ATT estimates?
2. **Within** a random forests model specification across repeated samples, does standardized mean difference of covariates relate to the magnitude of bias of ATT estimates?

Summary of Results (1)

Between different propensity score estimations:

1. Logistic Regression

Weighting *ASMDs* > **Matching** *ASMDs*

- **Both** had satisfactory *ASMDs*
- Suggested that *ASMDs* should **NOT** be compared between equating methods

Summary of Results (1)

Between different propensity score estimations:

2. Random Forests

- **Optimal specifications** had **small** and **satisfactory ASMDs**
- Specifications produced **VERY biased ATEs (> 60%)** had **large ASMDs**
- **But**, specifications produced **biased ATEs (10 – 60%)** had **small** and **satisfactory ASMDs**

Summary of Results (2)

Within a Random Forests specification across replications

1. **NO** substantial correlations between *ASMDs* and the magnitude of *ATT* bias (< 0.3)
2. Potential reason: Reduction in range, *ATT* bias and *ASMDs* are reduced, severely attenuating the correlations.

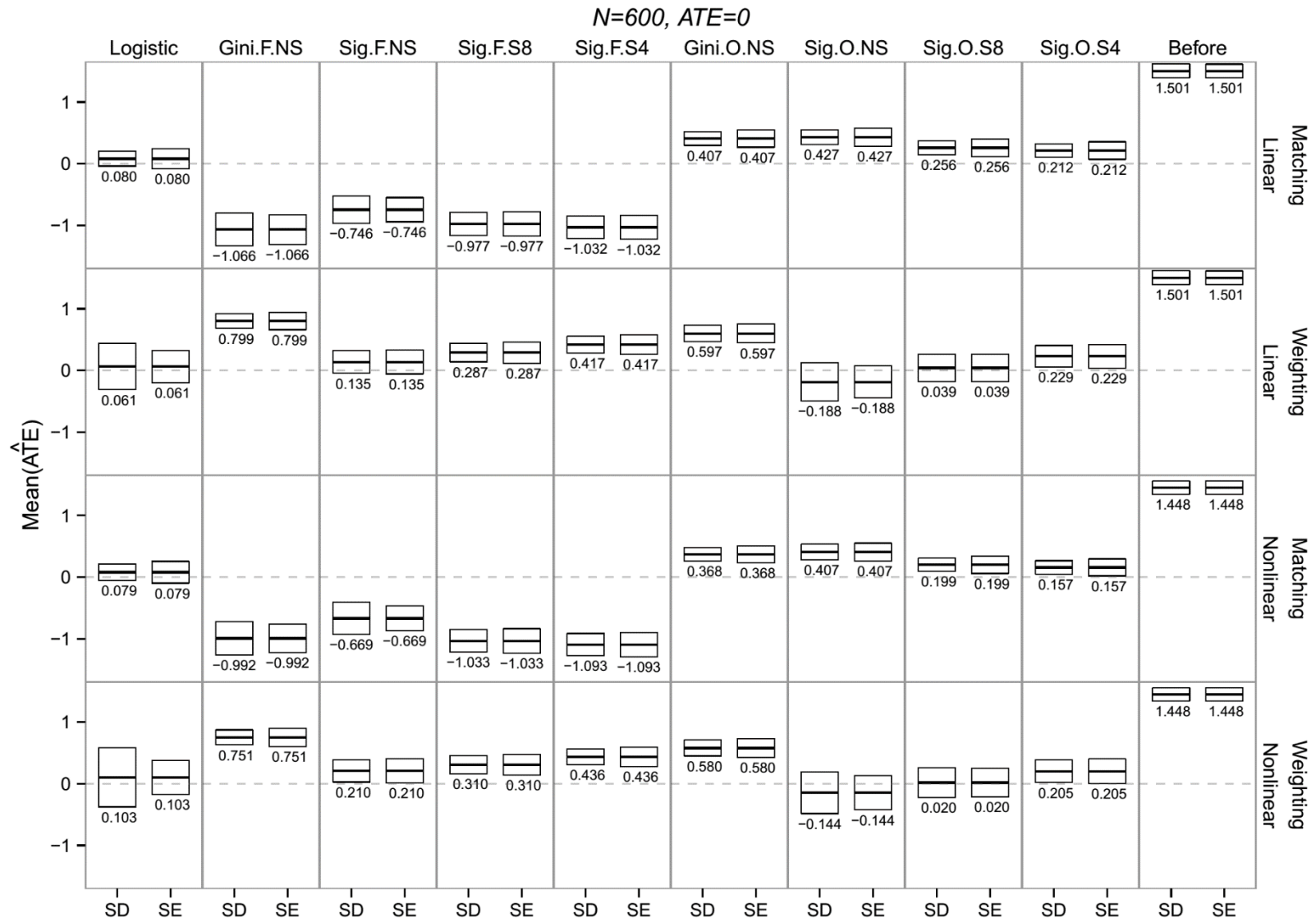
Implications: **NEW Procedures** to determine the optimal Random Forests specifications are needed.

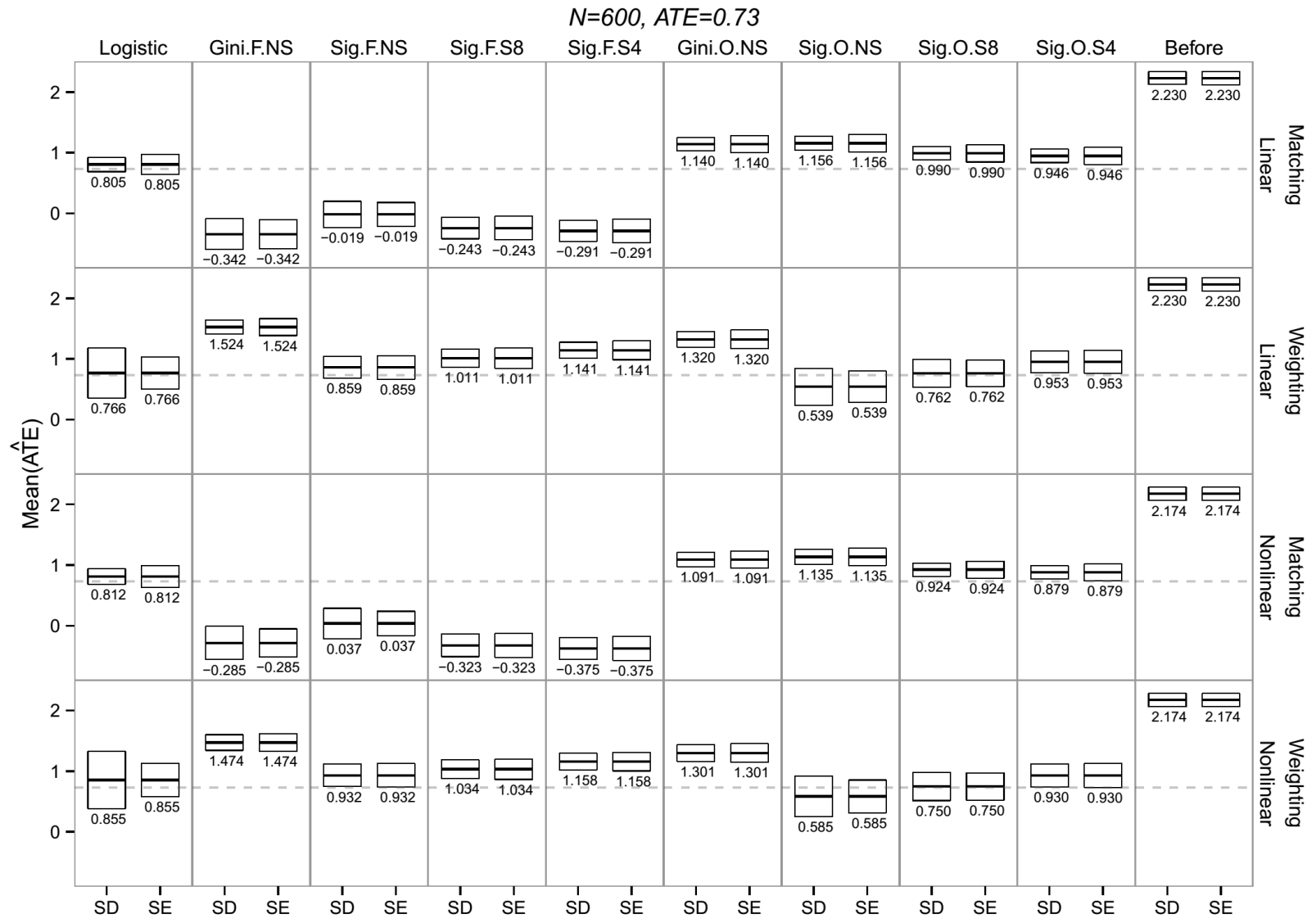
Concluding Remarks: A Dilemma

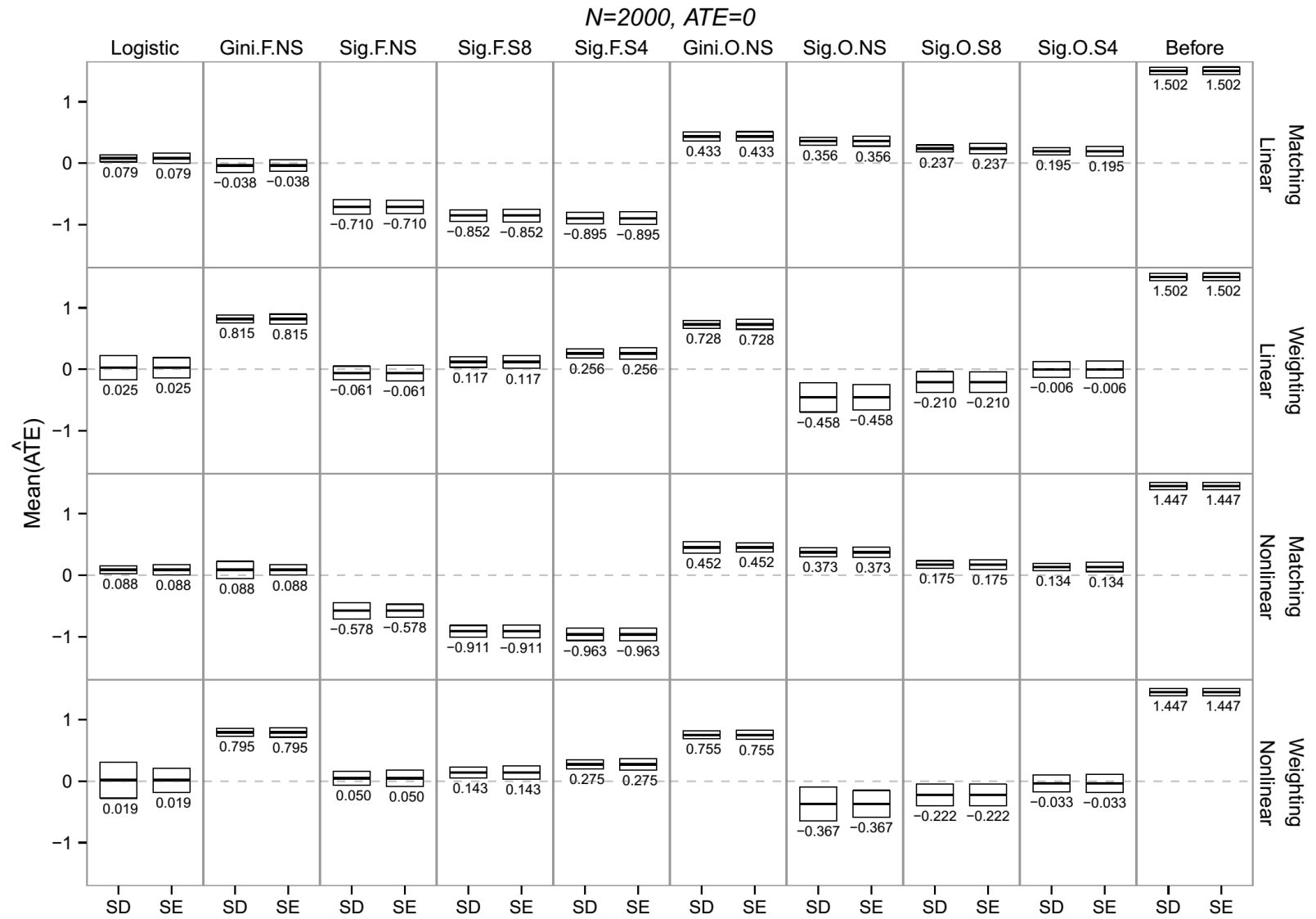
1. In our abstract, we hypothesized that the equating methods which are less dependent on misspecified propensity score estimation models produced unbiased ATT.
2. Our preliminary analyses show that these equating methods work.
3. BUT, if it is so, what is the point of using random forests to estimate propensity scores?

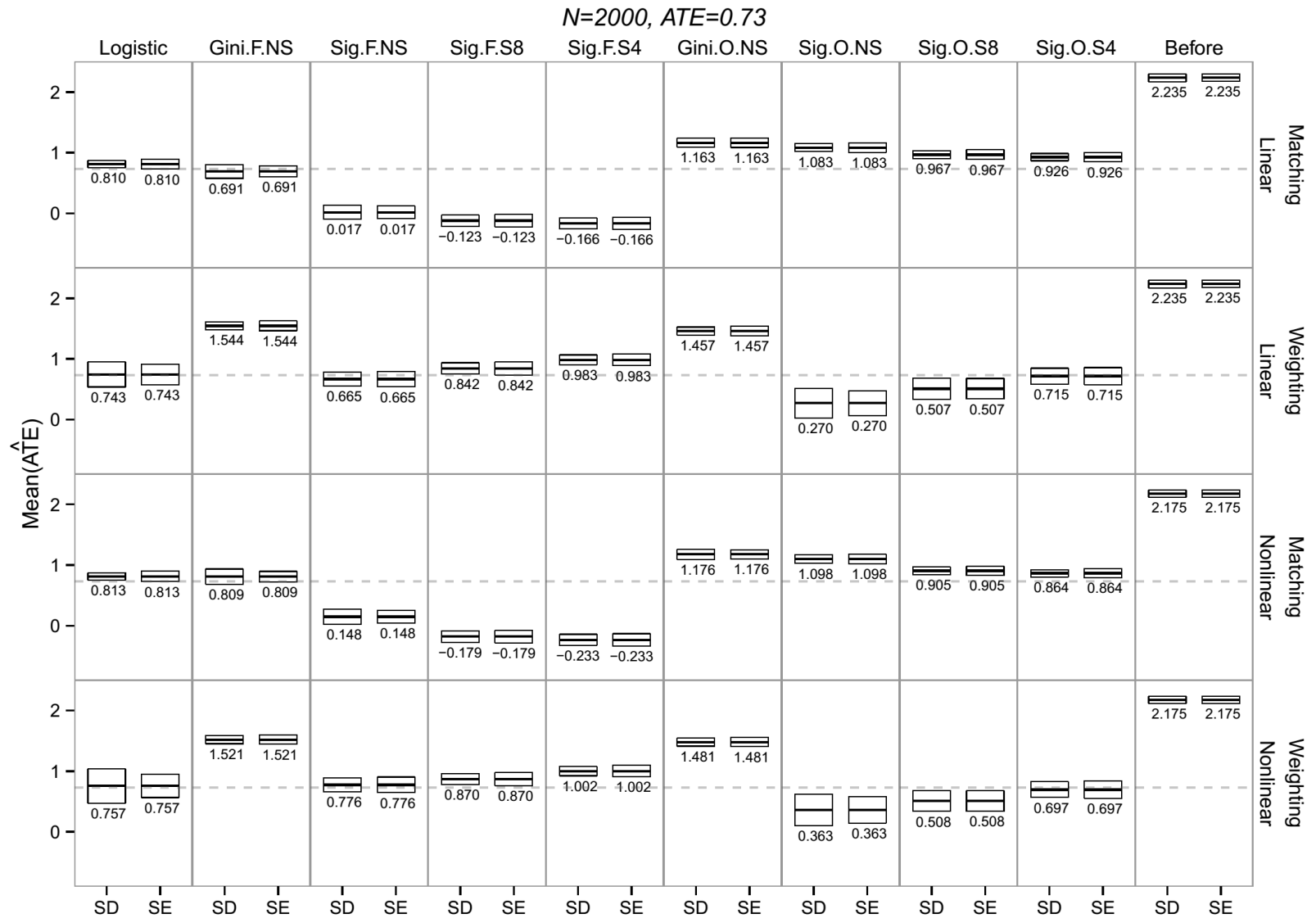
Appendix

FIGURES OF RESULTS



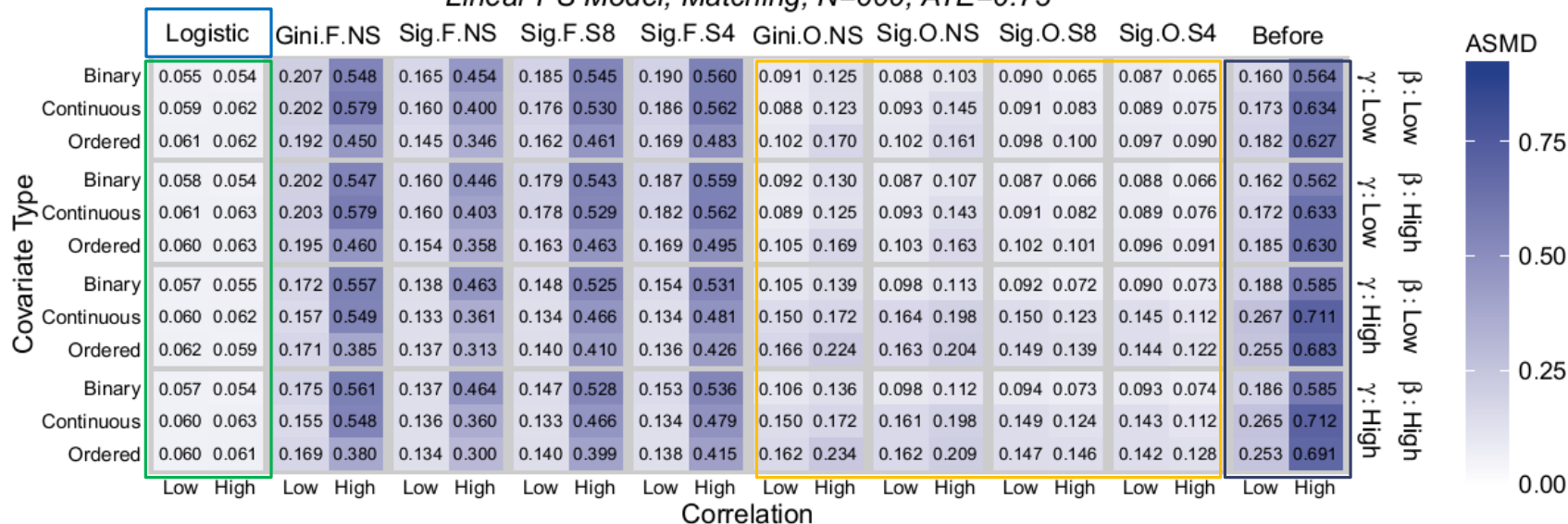




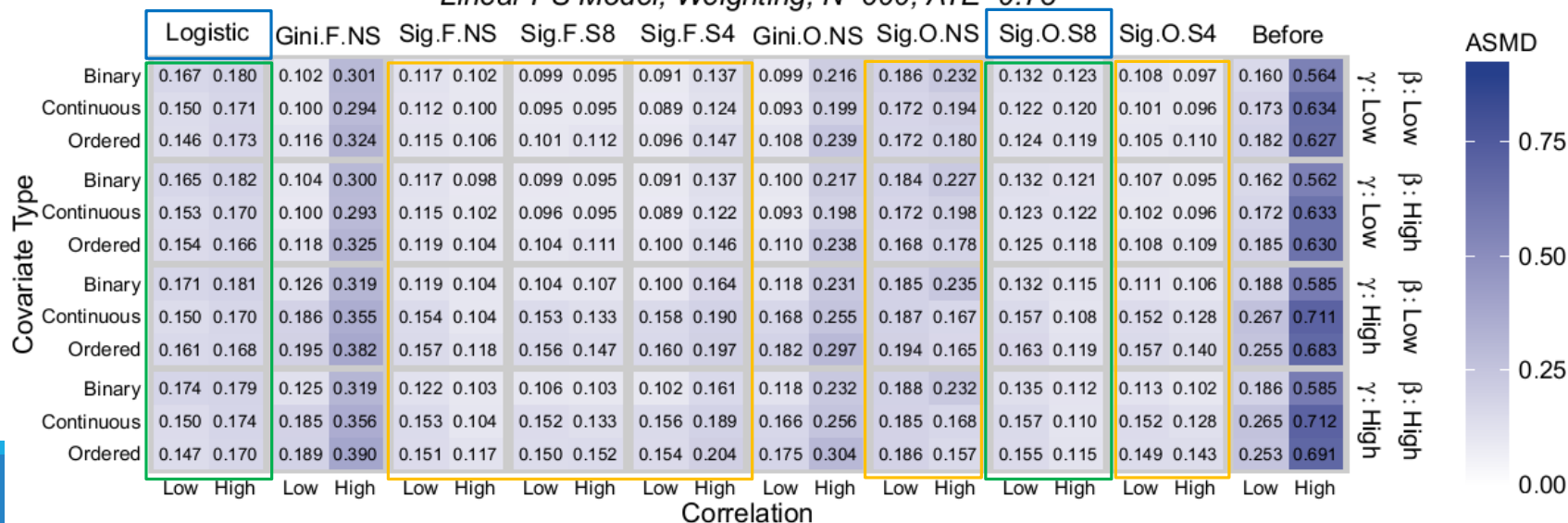


Mean *ASMD*

Linear PS Model, Matching, N=600, ATE=0.73

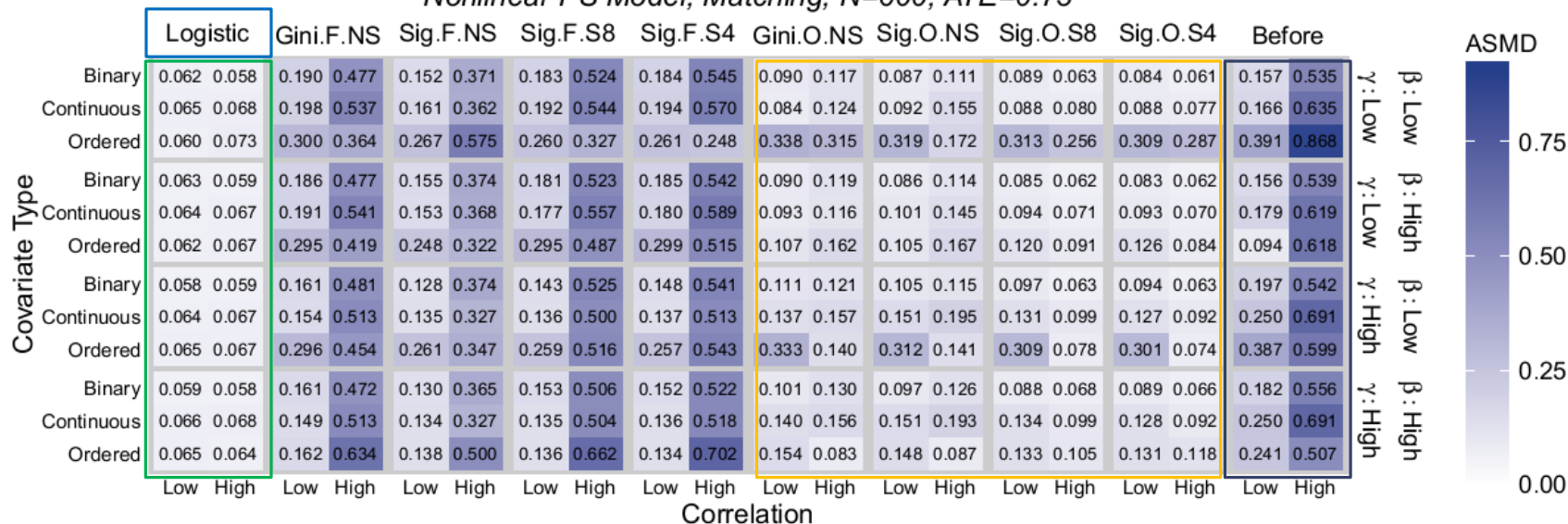


Linear PS Model, Weighting, N=600, ATE=0.73

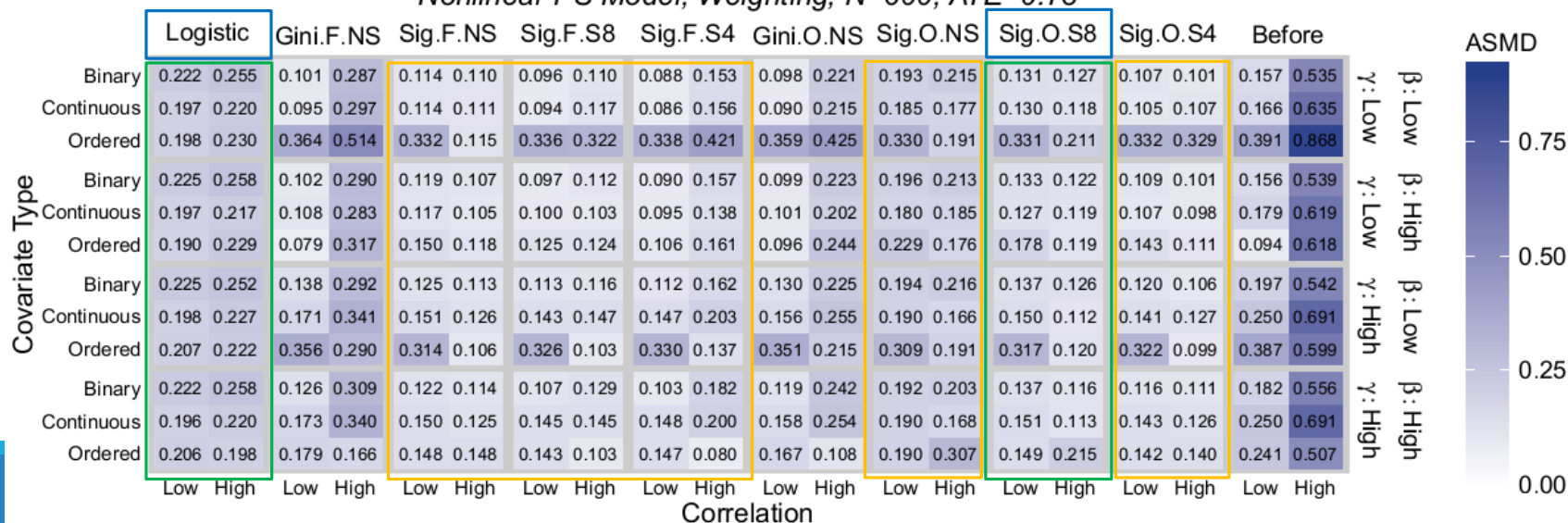


Mean ASMD

Nonlinear PS Model, Matching, N=600, ATE=0.73

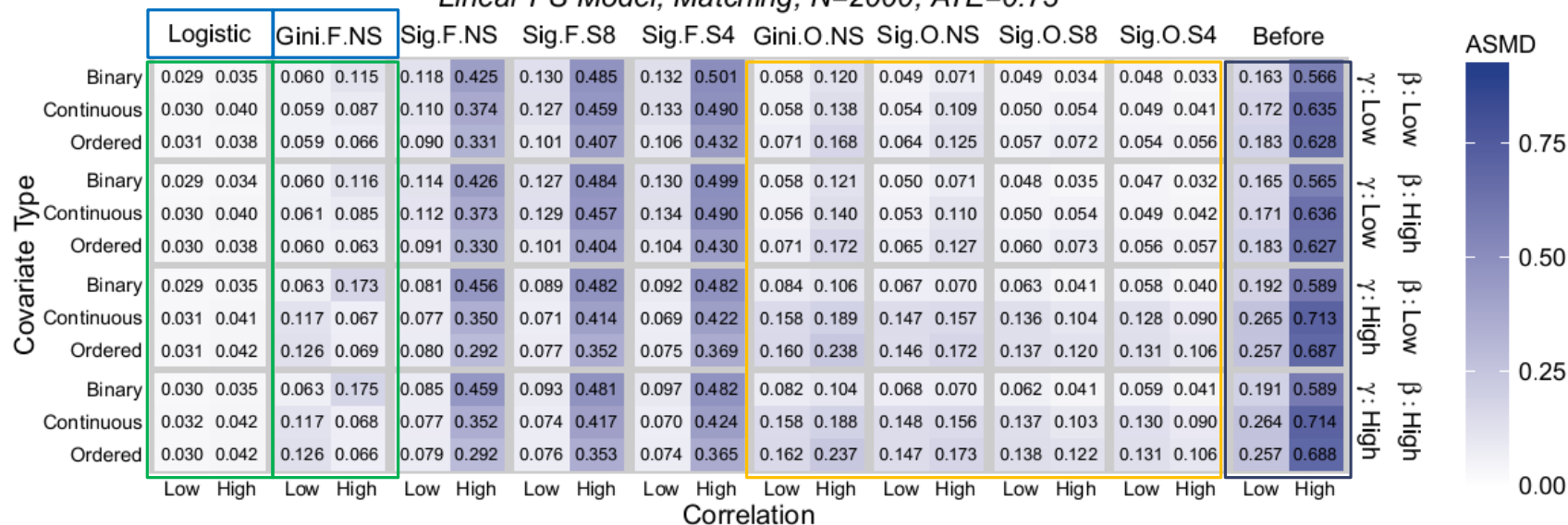


Nonlinear PS Model, Weighting, N=600, ATE=0.73

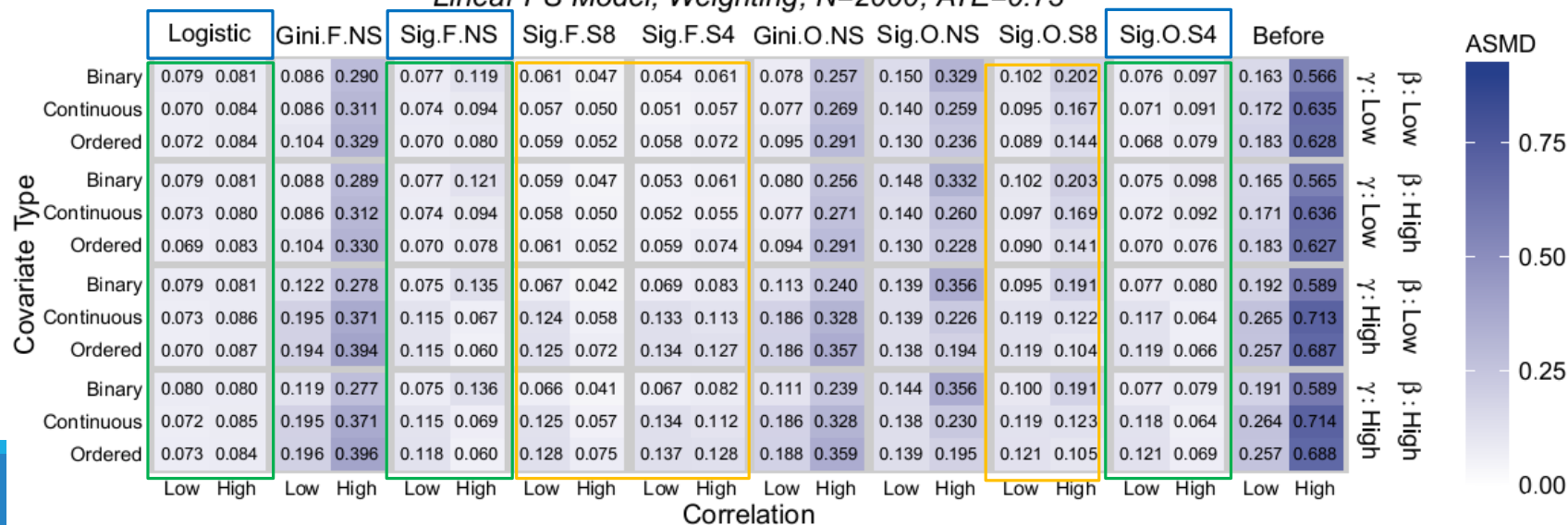


Mean ASMD

Linear PS Model, Matching, N=2000, ATE=0.73

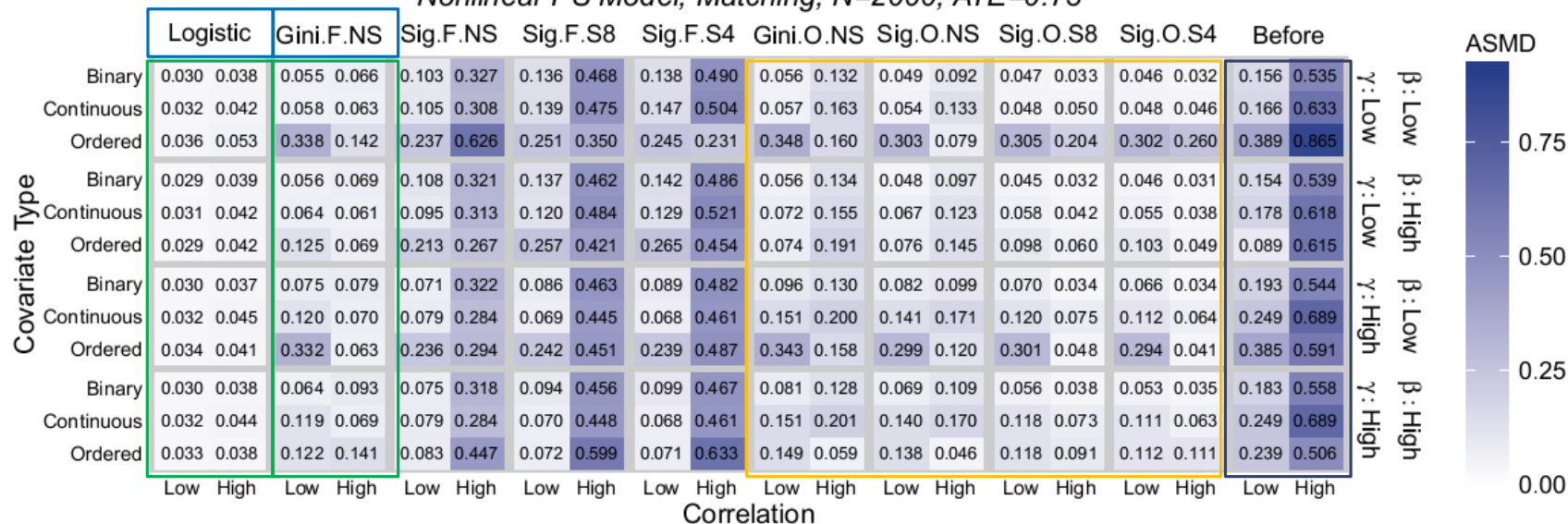


Linear PS Model, Weighting, N=2000, ATE=0.73



Mean ASMD

Nonlinear PS Model, Matching, N=2000, ATE=0.73



Nonlinear PS Model, Weighting, N=2000, ATE=0.73

