

Synthetic and negative control evaluation framework for large-scale propensity score survival analysis

Yuxi Tian

M.D./Ph.D candidate

Department of Biomathematics, UCLA



Joint work with: Marc Suchard - University of California, Los Angeles
Martijn Schuemie - Janssen Research and Development

OHDSI Community Call; September 5, 2017

Propensity Score Adjustment

- PS = estimated probability of treatment assignment
address confounding in observational studies

Propensity Score Adjustment

- PS = estimated probability of treatment assignment
address confounding in observational studies

How is the PS Estimated?

Propensity Score Adjustment

- PS = estimated probability of treatment assignment
address confounding in observational studies

How is the PS Estimated?

Logistic Regression

Propensity Score Adjustment

- PS = estimated probability of treatment assignment
address confounding in observational studies

How is the PS Estimated?

Logistic Regression

How are Covariates Selected?

Propensity Score Adjustment

- PS = estimated probability of treatment assignment
address confounding in observational studies

How is the PS Estimated?

Logistic Regression

How are Covariates Selected?

Thousands of potential confounders

PS Model Selection

- Traditionally: Investigator Selection
- high-dimensional Propensity Score algorithm (hdPS)
univariate screen for significant covariates
based on exposure or outcome association

“exposure-based” : relative risk with treatment exposure

“bias-based” : relative risk with outcome of interest

- L1-regularization (LASSO)
multivariate model selection via penalized likelihood
coefficients of unimportant covariates shrunk to zero

Study Goals

- Detail framework to evaluate propensity score estimation method performance
 - simulations
 - negative control experiments
- Use evaluation to compare:
 - hdPS Algorithm : “exposure-based” and “bias-based”
 - L1-regularization (LASSO)

PS Details

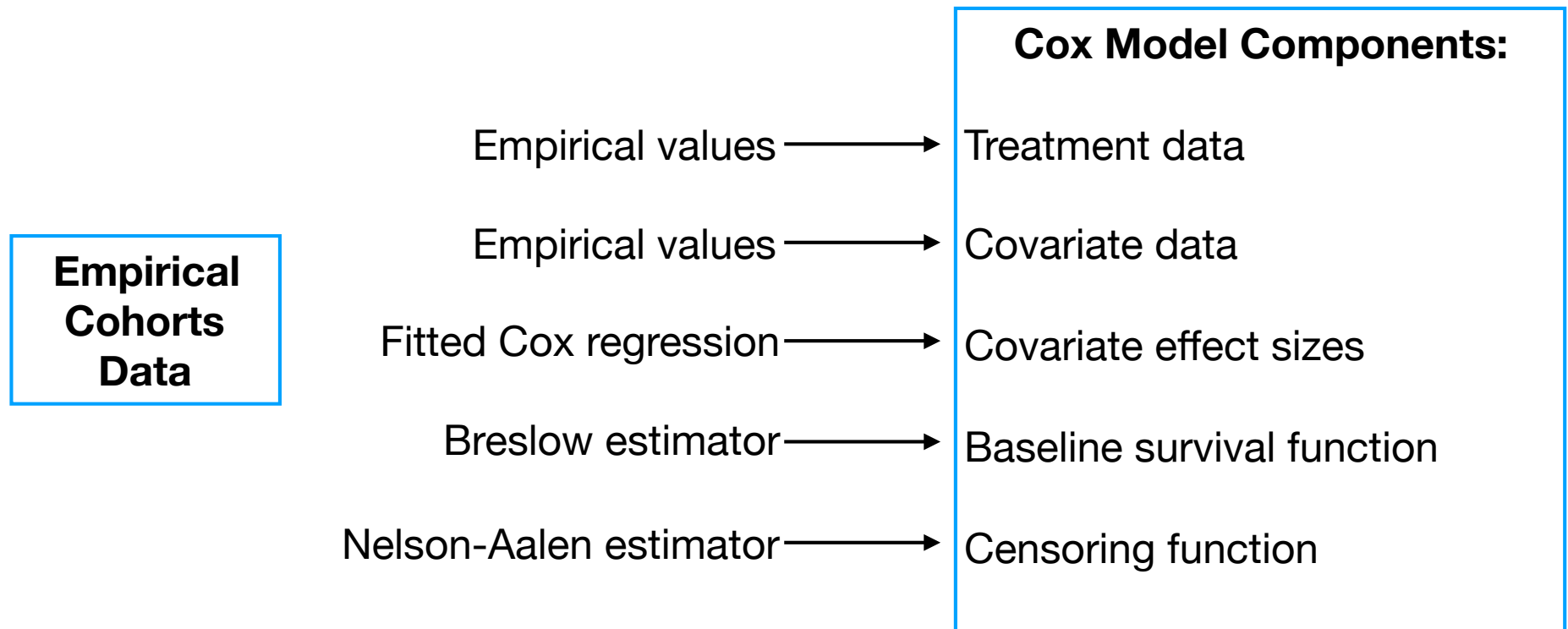
- hdPS Algorithm prescribes a certain set of data pre-processing:
 - aggregate covariates by coding
 - limit considered covariates to most prevalent
 - augment covariates by individual level frequency
 - 180 day lookback windows
- FeatureExtraction default uses more expansive set of covariates
 - eras, exposures, observations, measurements, scores
 - 30 day, 365 day, all day lookback windows
- We used L1-regularization on both (hdPS and CDM)

Simulations

- Keep treatment exposure and covariates from real-world data
- Simulate outcomes times under a survival model
- Simulate under **known hazard ratio** and with different outcome prevalences
- Extends the “plasmode” framework by Franklin et al. (2014)

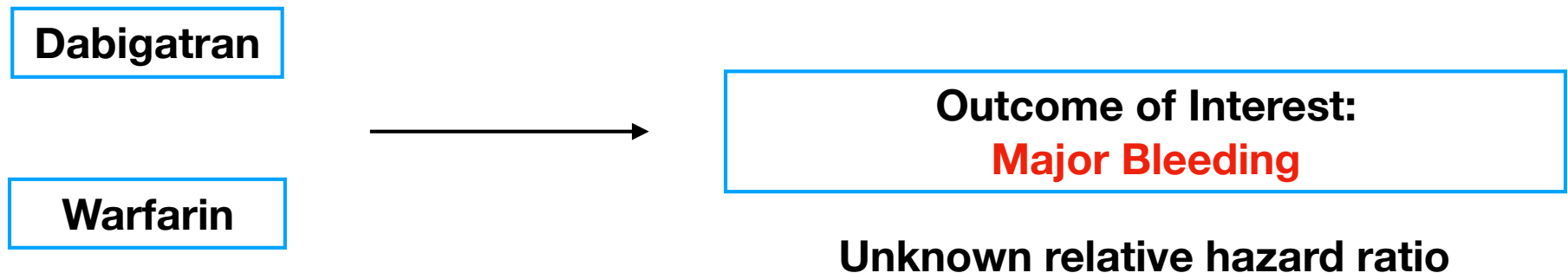
Simulations

- Simulate realistic survival data under a **known hazard ratio** in Cox proportional hazards model



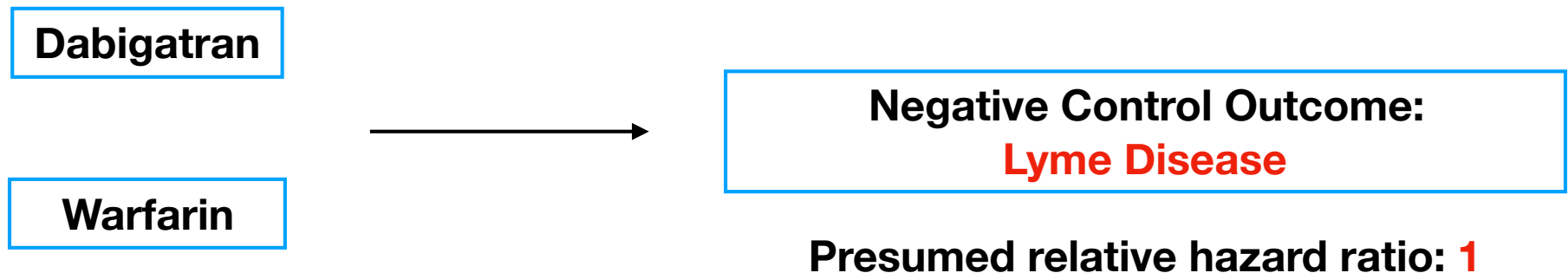
Negative Control Experiments

- Downside to simulations:
Do not capture full complexity of real-world data
- Negative controls:
Outcomes unaffected by the studied treatments



Negative Control Experiments

- Downside to simulations:
Do not capture full complexity of real-world data
- Negative controls:
Outcomes unaffected by the studied treatments

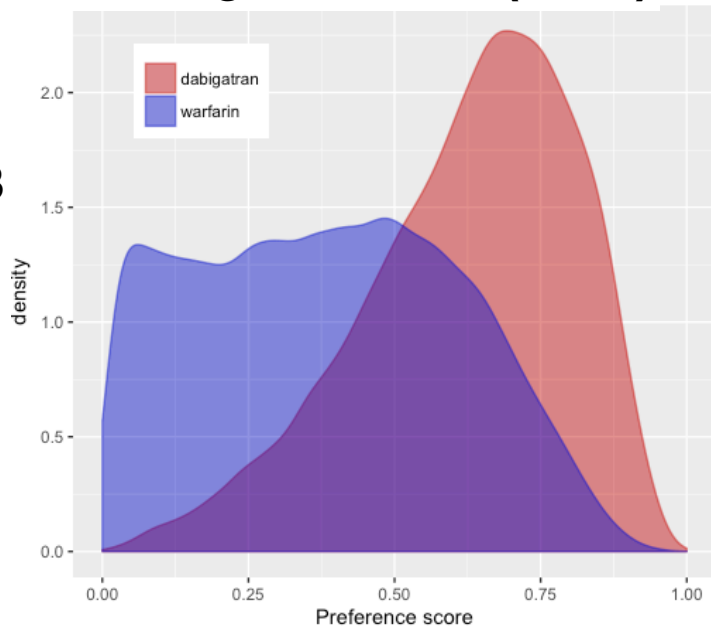


Empirical Data Used - Anticoagulants

- Replication of dabigatran vs warfarin observational study by Graham et al. (2014)
- Database: Truven Health Marketscan Medicare Supplemental and Coordination of Benefits Database
- Cohorts:
 - 19768 dabigatran users, 52721 warfarin users
 - 192 intracranial hemorrhage 0.26%
 - 98118 unique covariates

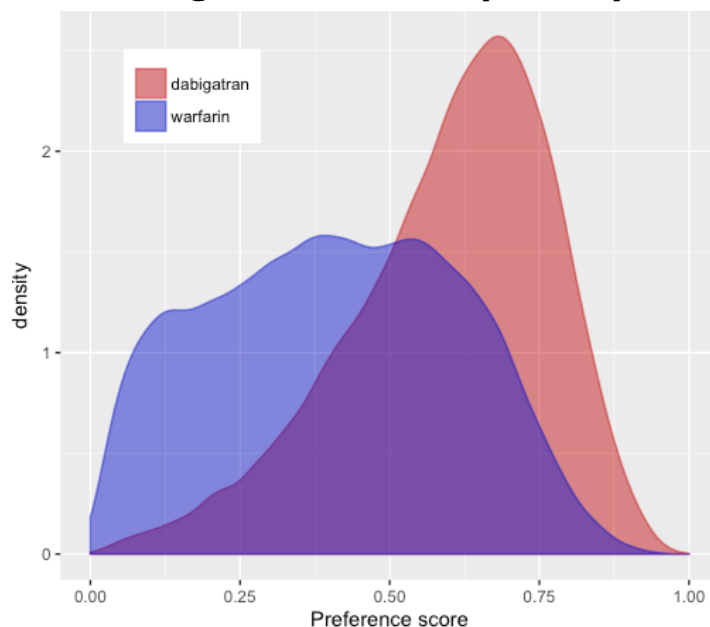
PS Distribution

L1 Regularization (CDM)



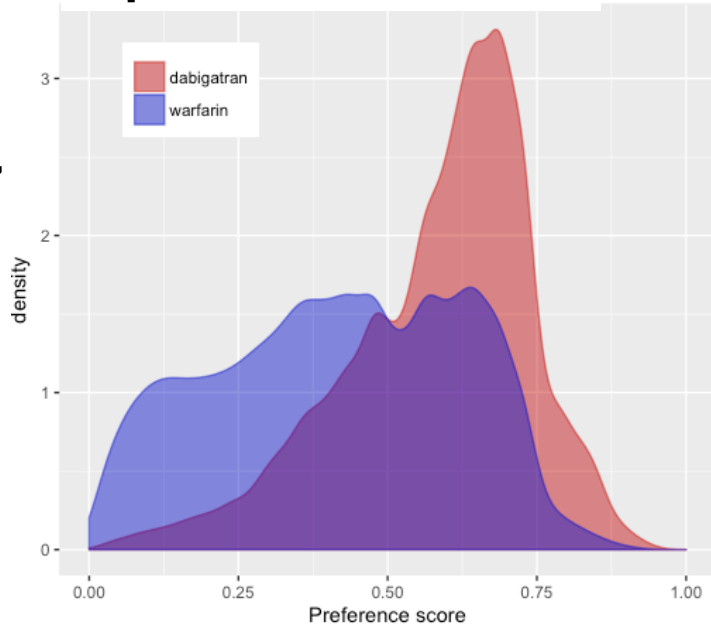
AUC: 0.793

L1 Regularization (hdPS)



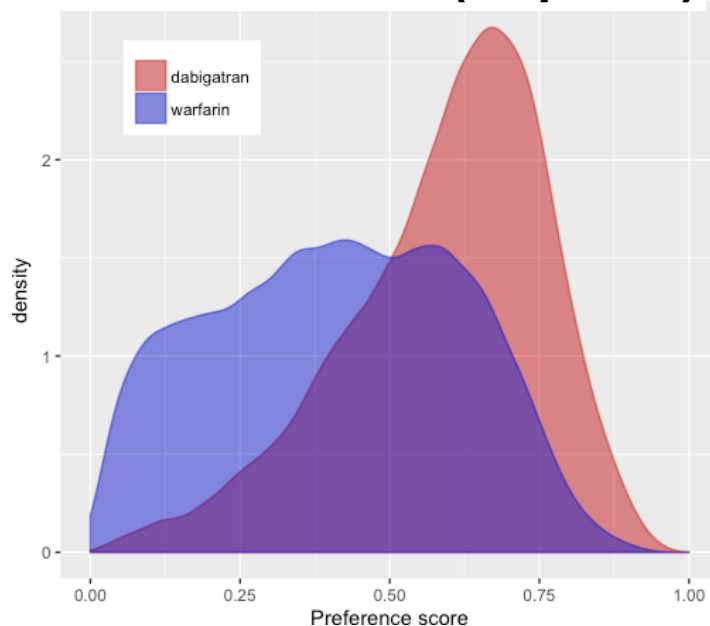
AUC: 0.760

exposure-based hdPS



AUC: 0.737

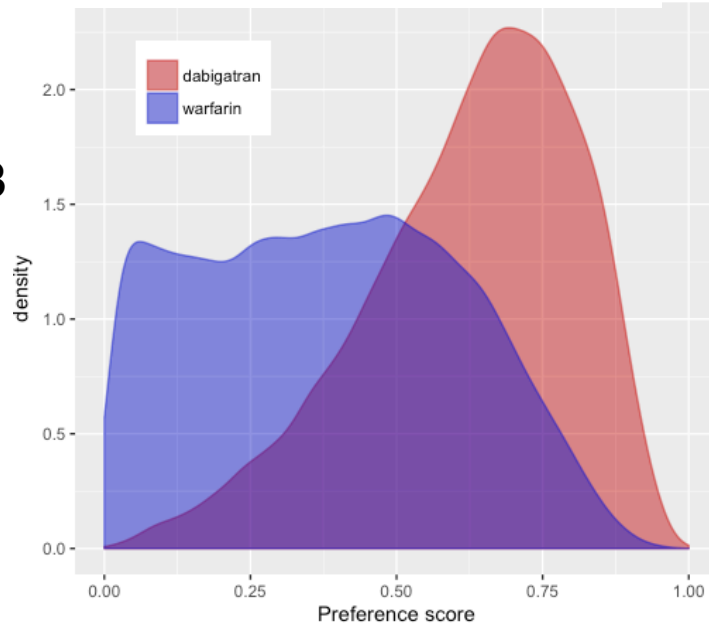
bias-based hdPS (empirical)



**Empirical:
AUC: 0.747**

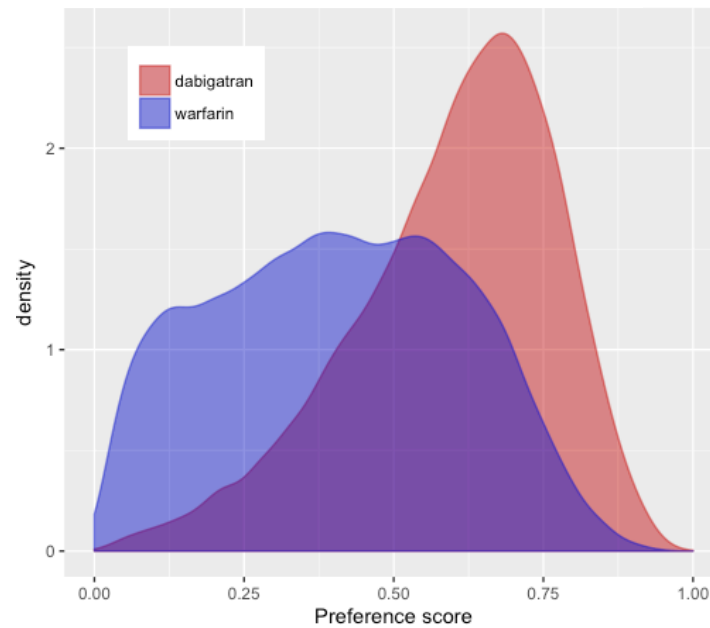
PS Distribution

L1 Regularization (CDM)



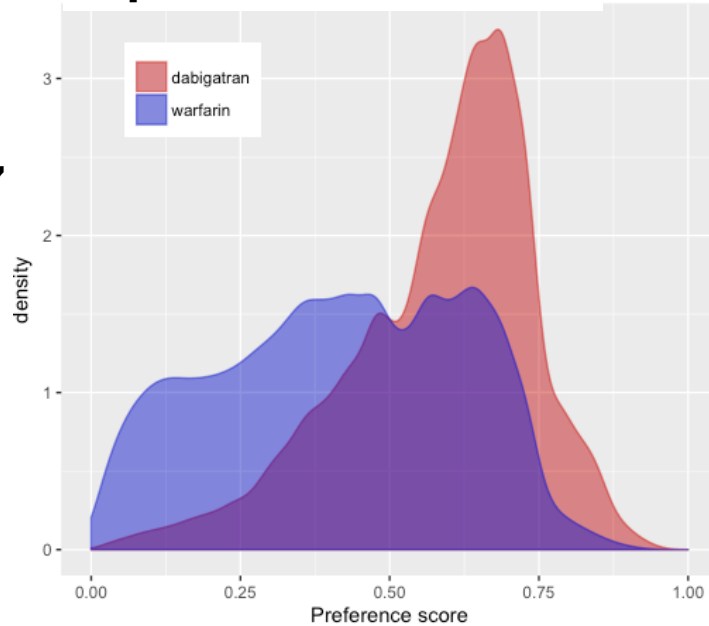
AUC: 0.793

L1 Regularization (hdPS)



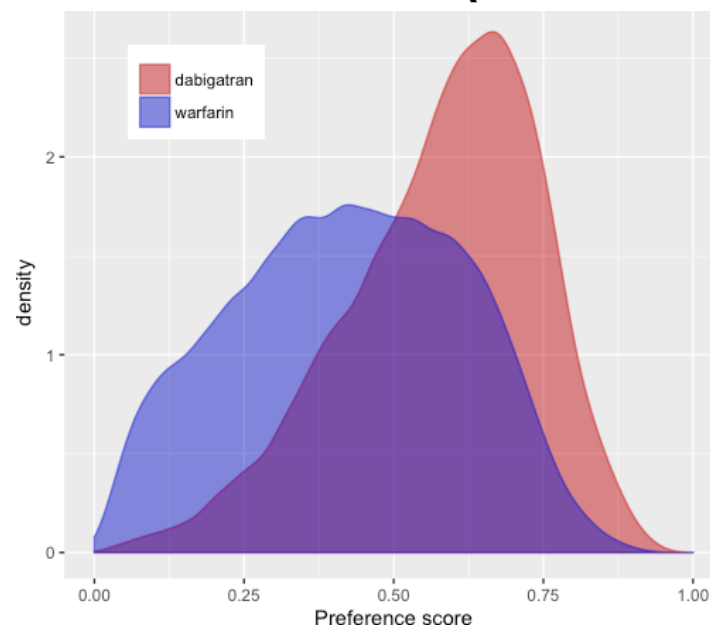
AUC: 0.760

exposure-based hdPS



AUC: 0.737

bias-based hdPS (simulation)



**Empirical:
AUC: 0.747**

**Simulation:
AUC: 0.742**

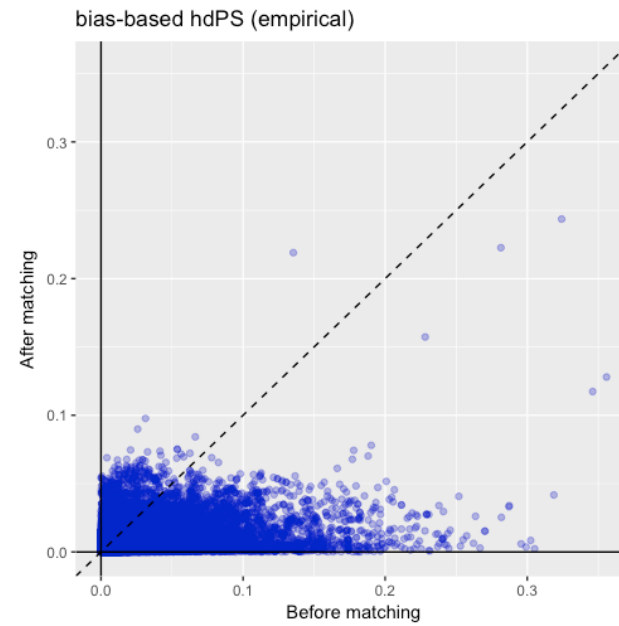
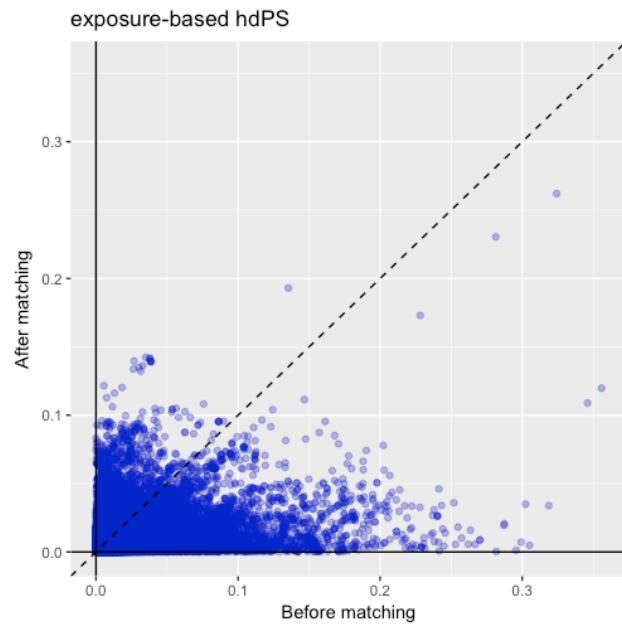
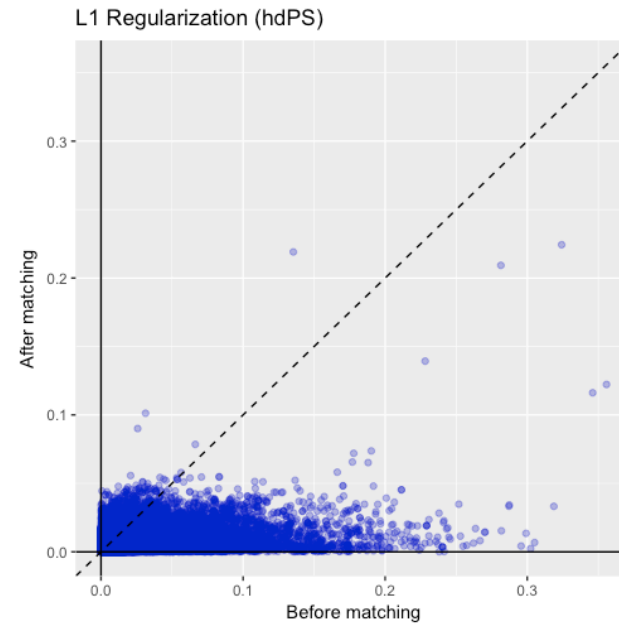
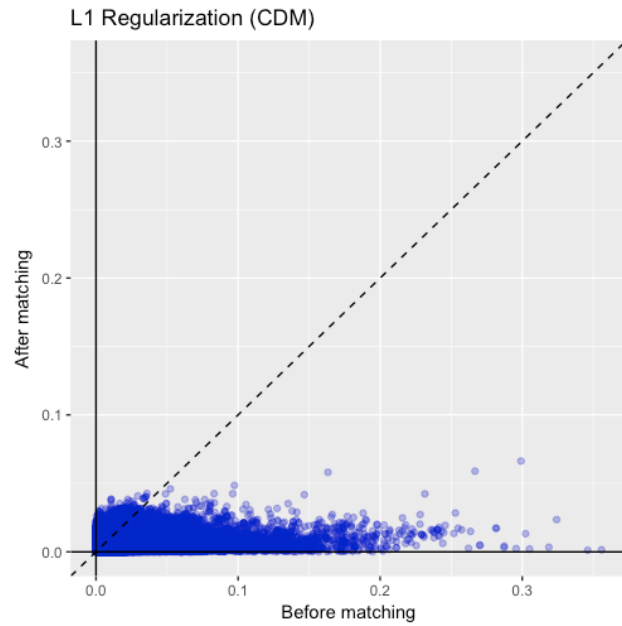
Covariate Balance

- standardized difference of covariates before and after propensity score matching

Which covariates to consider?

- All covariates
- “true confounders”
 - approximated by simulation model covariates
 - note: these include “hdPS Algorithm Covariates” and “CDM Covariates”

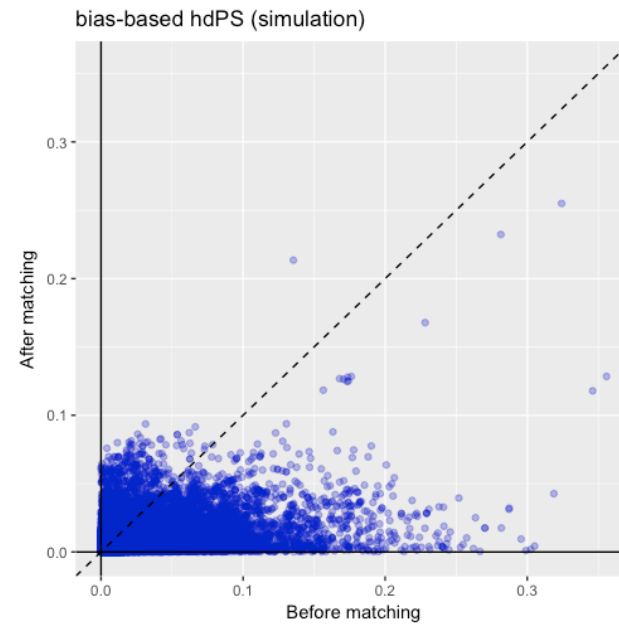
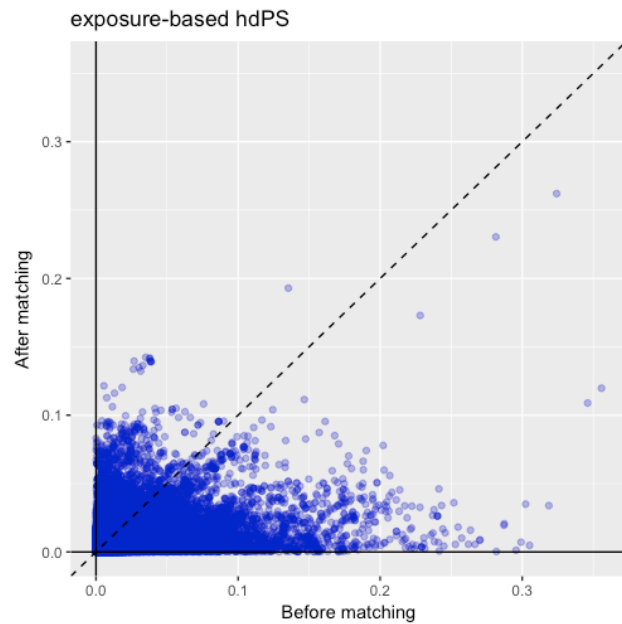
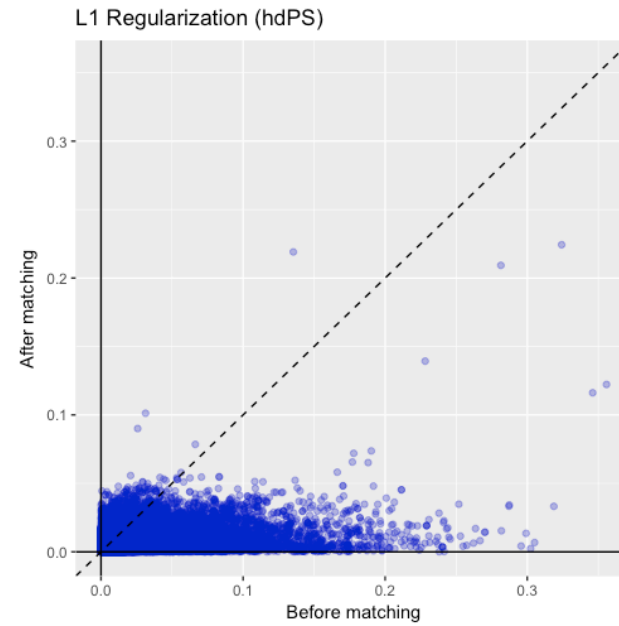
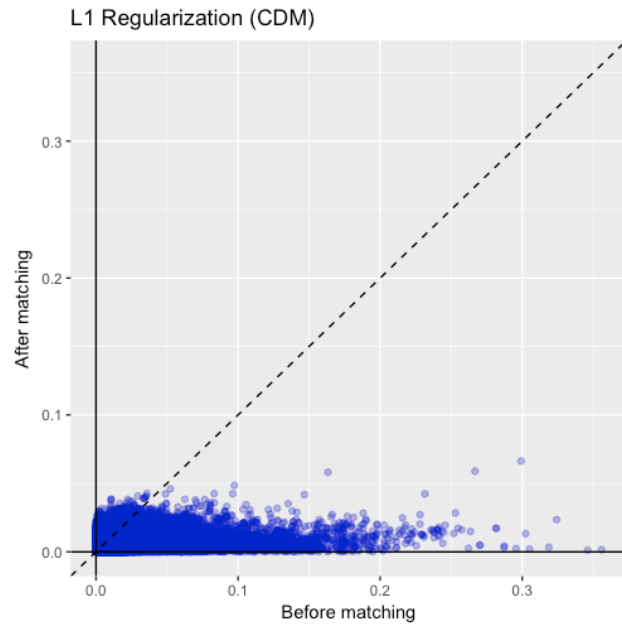
All Covariates



**Outcome
Independent
Metric**

10:1 variable ratio matching

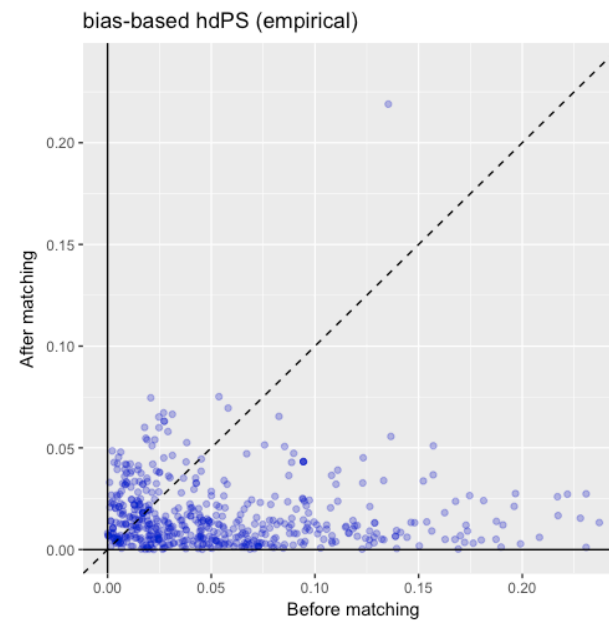
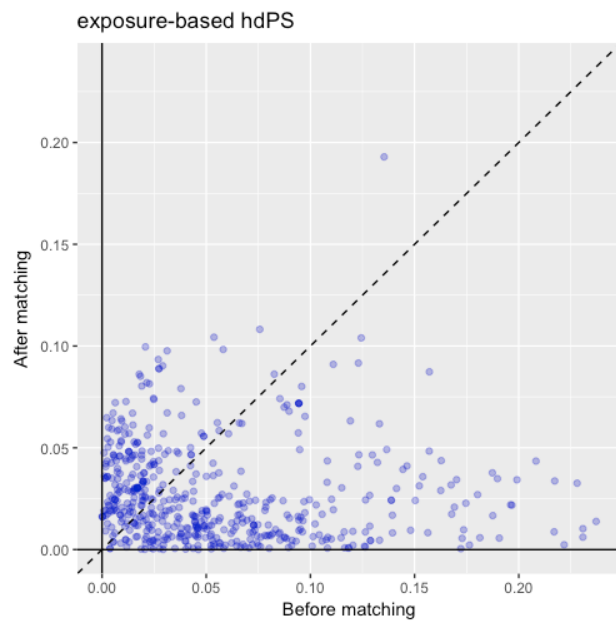
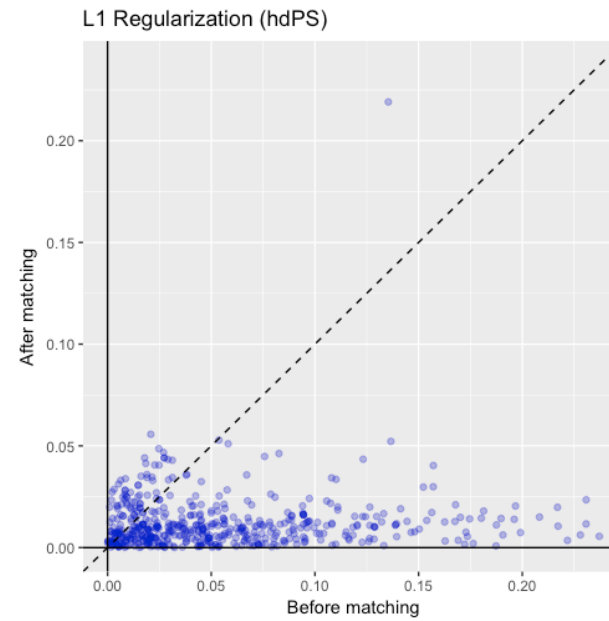
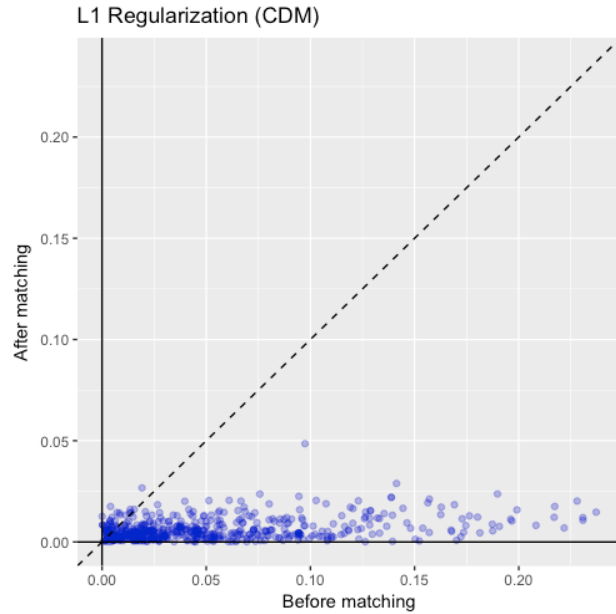
All Covariates



**Outcome
Independent
Metric**

10:1 variable ratio matching

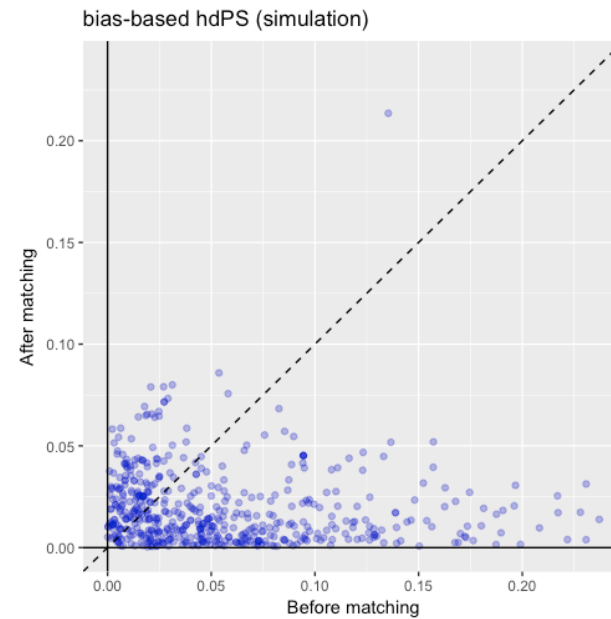
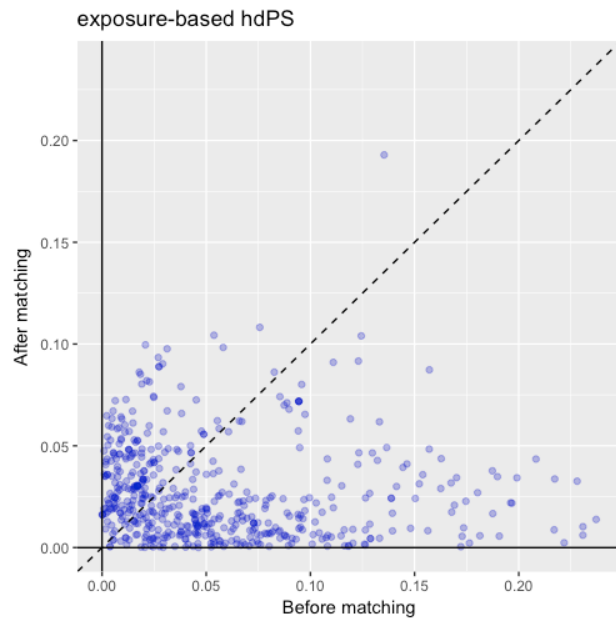
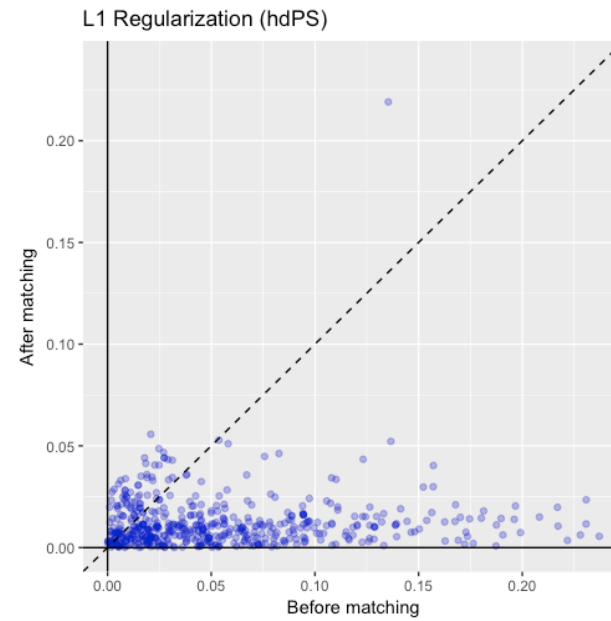
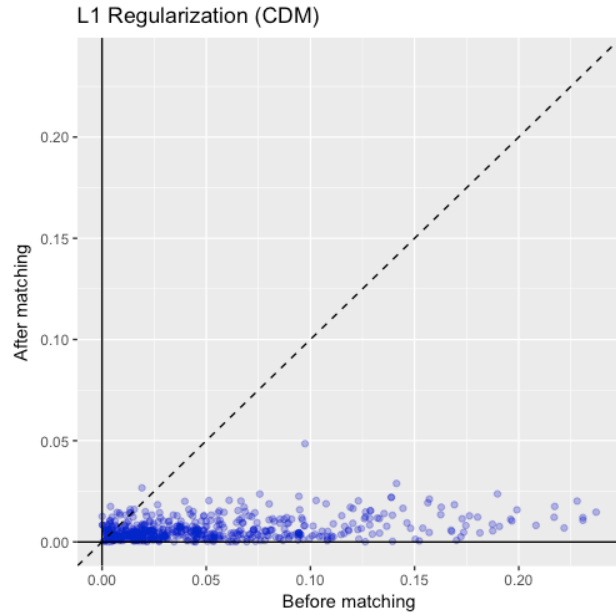
Simulation Model Covariates



**Outcome
Independent
Metric**

10:1 variable ratio matching

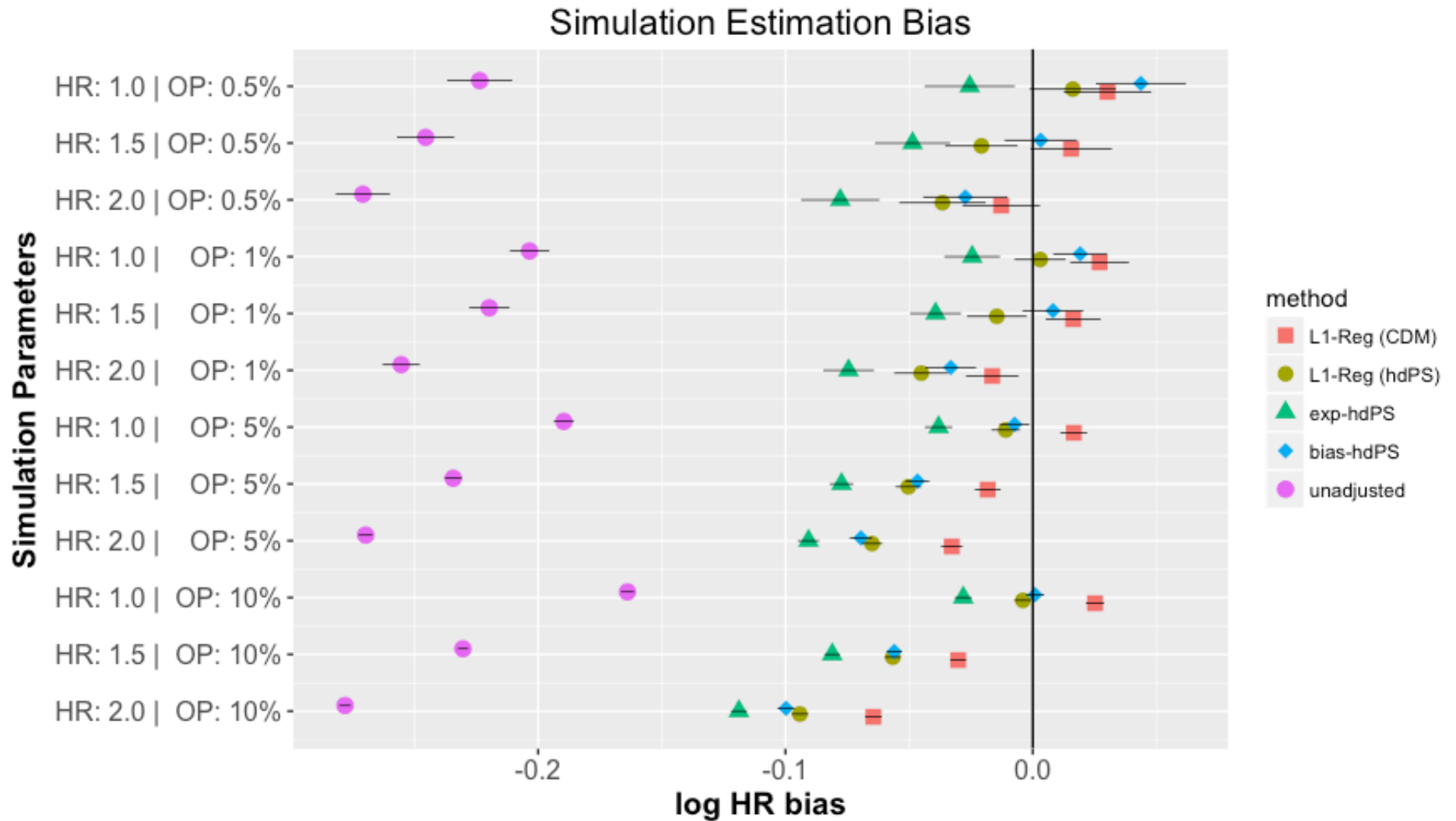
Simulation Model Covariates



**Outcome
Independent
Metric**

10:1 variable ratio matching

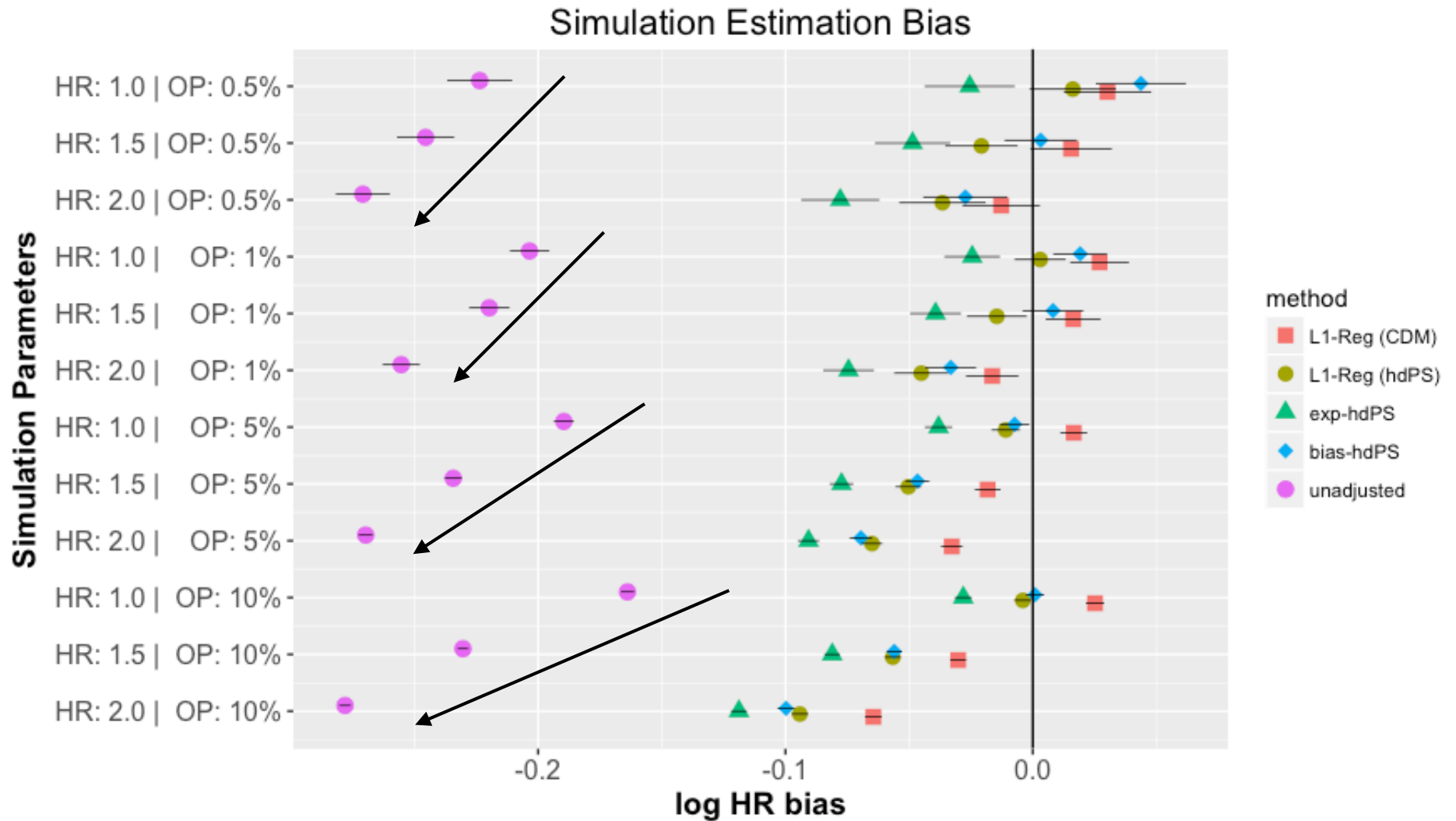
Bias Reduction: Simulations



**Outcome
Dependent
Metric**

10:1 variable ratio matching

Bias Reduction: Simulations



**Outcome
Dependent
Metric**

10:1 variable ratio matching

Simulation Bias

Survival Simulation; consider 1:1 matching

$$\hat{\eta} = \log N_1 - \log N_0$$

N_1 : exposed has event, time before unexposed

N_0 : unexposed has event, time before exposed

$$\Pr(\text{set in } N_1) = \int_0^\infty \left(\frac{\partial}{\partial t} S(t)^{\exp\{\theta_{1,k}\}} \right) S(t)^{\exp\{\theta_{0,k}\}} C(t) C(t) dt$$

The diagram consists of four labels at the bottom with arrows pointing upwards to the corresponding terms in the equation above:

- exposed hazard** (with **contains true effect size** in red text below it) has an arrow pointing to $\frac{\partial}{\partial t} S(t)^{\exp\{\theta_{1,k}\}}$.
- unexposed hazard** has an arrow pointing to $S(t)^{\exp\{\theta_{0,k}\}}$.
- survival function** has an arrow pointing to $S(t)^{\exp\{\theta_{1,k}\}}$.
- censoring function** has an arrow pointing to $C(t)$.

Simulation Bias

Survival Simulation; consider 1:1 matching

$$\hat{\eta} = \log N_1 - \log N_0$$

N_1 : exposed has event, time before unexposed

N_0 : unexposed has event, time before exposed

$$\Pr(\text{set in } N_1) = \int_0^\infty \left(\frac{\partial}{\partial t} S(t)^{\exp\{\theta_{1,k}\}} \right) S(t)^{\exp\{\theta_{0,k}\}} C(t) C(t) dt$$

↑
survival function

↑
exposed hazard
contains true effect size

↑
censoring function

↑
unexposed hazard

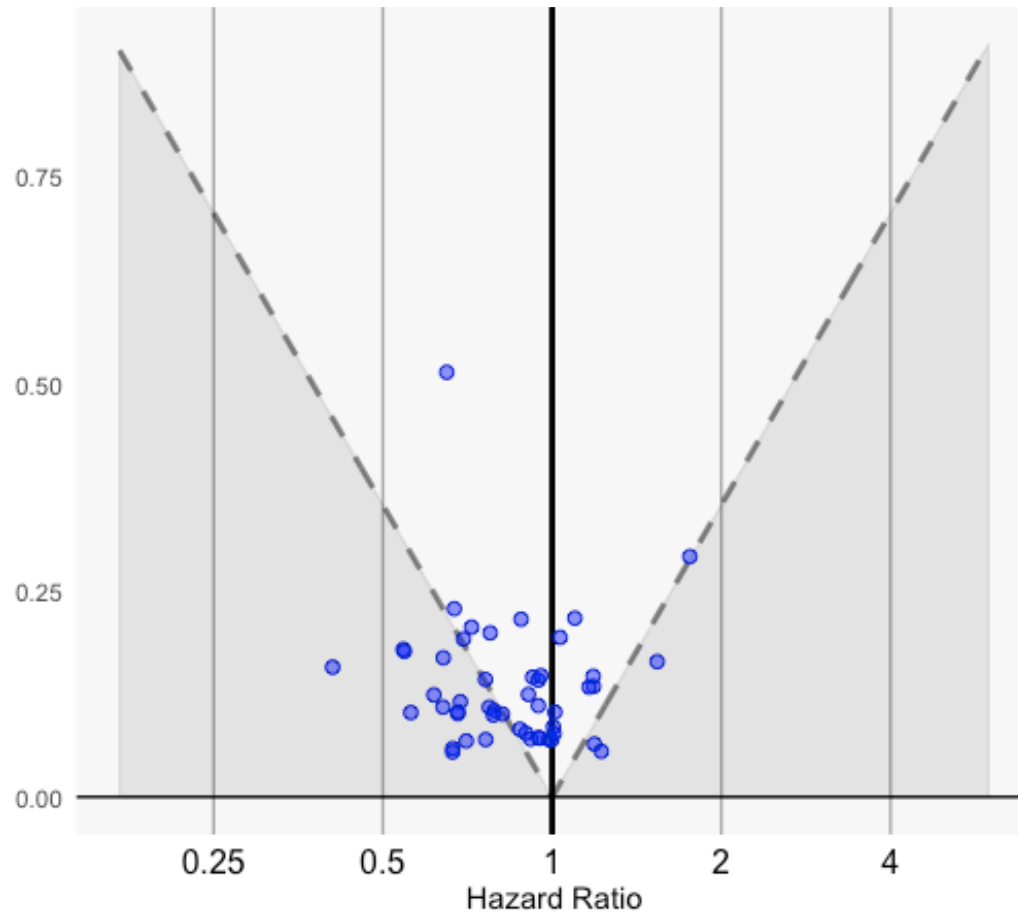
**Not unbiased
when there is
variance in**

baseline hazards

Negative Controls

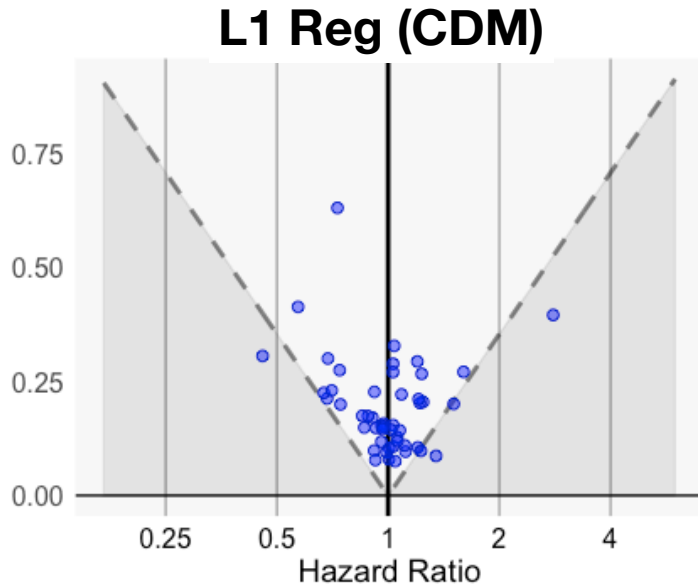
Unadjusted

Coverage:
 0.53 ± 0.07

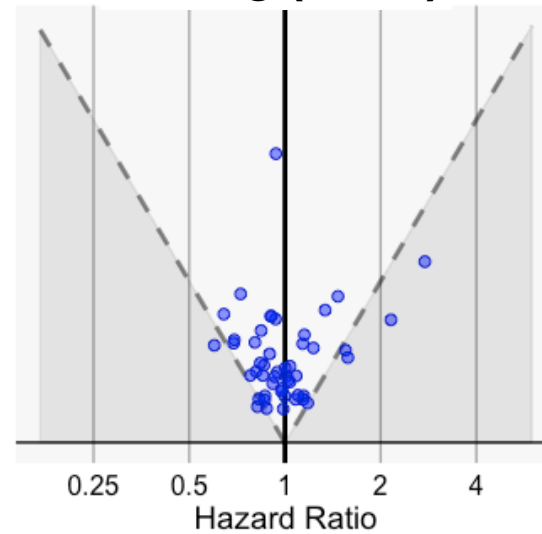


Bias Reduction: Negative Outcomes

Coverage:
 0.90 ± 0.04



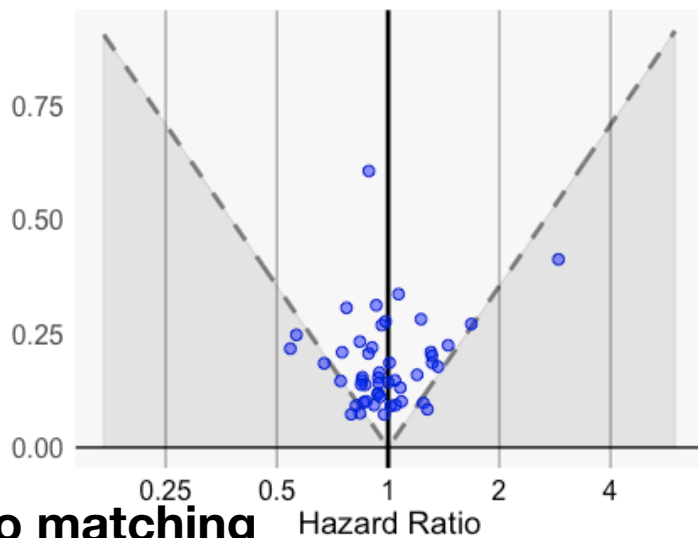
L1 Reg (hdPS)



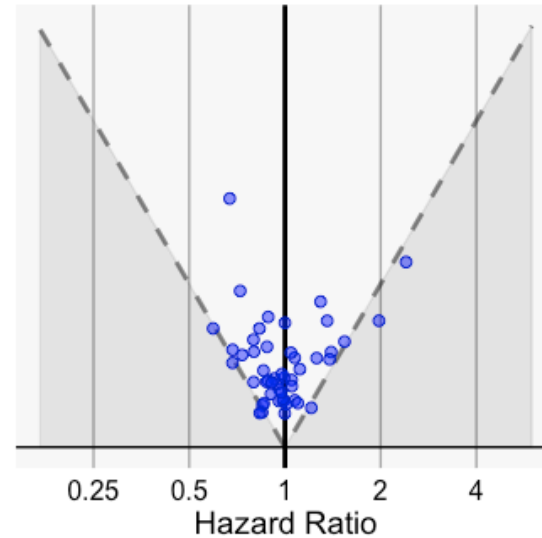
Coverage:
 0.86 ± 0.05

Coverage:
 0.80 ± 0.06

exposure-based hdPS



bias-based hdPS



Coverage:
 0.86 ± 0.05

**Outcome
Dependent
Metric**

10:1 variable ratio matching

Outcome Dependent Metrics

Outcome Dependent Metrics

- Susceptible to bias:
 - PS adjustment techniques
 - simulation design choices
 - negative control misspecification

Outcome Dependent Metrics

- Susceptible to bias:
 - PS adjustment techniques
 - simulation design choices
 - negative control misspecification
- Different outcomes can yield different results

Outcome Dependent Metrics

- Susceptible to bias:
 - PS adjustment techniques
 - simulation design choices
 - negative control misspecification
- Different outcomes can yield different results
- Outcome independent metrics more generalizable

Instrument Variables

- Variables that predict treatment exposure but has no effect on outcome (or correlation with any confounder)
- Inclusion in PS can increase bias and variance of estimate

Suppose:

- eye color perfectly separates treatment groups (all blue eyed receive A, all brown eyed receive B)
- eye color does not influence outcome
- no power in experiment

Instrument Variables

- Variables that predict treatment exposure but has no effect on outcome (or correlation with any confounder)
- Inclusion in PS can increase bias and variance of estimate

Suppose:

- absent of IV, PS correlated with outcome hazard, PS matches patients with similar baseline outcome hazard
- add in IV, PS of many exposed people increases
- exposed people now matched with higher hazard
- negative bias results

Instrument Variables

- True IV are rare, impact on real-world data unproven
- IV only problematic if uncorrelated with any confounders - unlikely situation in real-world data
- Identifying IV's is difficult
- bias-based hdPS uses outcome information in PS to avoid IV's, but breaks Rubin's unconfoundedness assumption

Instrument Variables - Solution?

- If certain IV's are suspected, stratify on them in the PS logistic regression -> conditional logistic regression (CLR)
- CLR avoids estimating any effect size from IVs
- Keeps unconfoundedness while eliminates effects on PS
- Issue:
CLR computationally expensive for large strata
CLR approximations can be very inaccurate
- Future direction:
Efficient CLR implementation, apply to PS

Take Away Points

- L1 Regularization favorable over hdPS Algorithm
- Simulations and negative controls provided useful evidence
- Regularization solves PS “convergence” problem (no MLE for regression exists)

