

---

# RealCause: Realistic Causal Inference Benchmarking

---

Brady Neal  
Mila, Université de Montréal

Chin-Wei Huang  
Mila, Université de Montréal

Sunand Raghupathi  
Mila, Université de Montréal

## Abstract

1        There are many different causal effect estimators in causal inference. However, it  
2        is unclear how to choose between these estimators because there is no ground-truth  
3        for causal effects. A commonly used option is to simulate synthetic data, where  
4        the ground-truth is known. However, the best causal estimators on synthetic data  
5        are unlikely to be the best causal estimators on real data. An ideal benchmark for  
6        causal estimators would both (a) yield ground-truth values of the causal effects and  
7        (b) be representative of real data. Using flexible generative models, we provide a  
8        benchmark that both yields ground-truth and is realistic. Using this benchmark,  
9        we evaluate over 1500 different causal estimators and provide evidence that it is  
10       rational to choose hyperparameters for causal estimators using predictive metrics.

## 11    1 Introduction

12    In causal inference, we want to measure causal effects of treatments on outcomes. Given some  
13    outcome  $Y$  and a binary treatment  $T$ , we are interested in the *potential outcomes*  $Y_i(1)$  and  $Y_i(0)$ .  
14    Respectively, these denote the outcome that unit  $i$  would have if they were to take the treatment  
15    ( $T = 1$ ) and the outcome they would have if they were to not take the treatment ( $T = 0$ ). We are  
16    often interested in causal estimands such as  $\mathbb{E}[Y(1)] - \mathbb{E}[Y(0)]$ , the *average treatment effect* (ATE).  
17    This is equivalent to the following expression using Pearl’s do-notation (Pearl, 1994, 2009, 2019):  
18     $\mathbb{E}[Y \mid \text{do}(T = 1)] - \mathbb{E}[Y \mid \text{do}(T = 0)]$ , where  $\text{do}(T = t)$  is a more mnemonic way of writing that  
19    we set the value of the treatment to  $t$ .

20    There are many different estimators for estimating causal estimands (see, e.g., Neal, 2020; Hernán &  
21    Robins, 2020; Morgan & Winship, 2014; Imbens & Rubin, 2015, and Appendix E). However, it is  
22    unclear how to choose between these estimators because the true values of the causal estimands are  
23    generally unknown. This is because we cannot observe both potential outcomes (Rubin, 1974), so  
24    we have no ground-truth. This is often referred to as the *fundamental problem of causal inference*  
25    (Holland, 1986). Supervised machine learning does not have this “no ground-truth” problem because  
26    it is only interested in estimating  $\mathbb{E}[Y \mid T]$ , which only requires samples from  $P(Y \mid T)$ , rather than  
27    samples from  $P(Y \mid \text{do}(T = 1))$  and  $P(Y \mid \text{do}(T = 0))$ . Yet, we must choose between causal  
28    estimators. How can we do that when faced with the fundamental problem of causal inference?

29    To evaluate causal estimators, people have created various benchmarks, each bringing different  
30    strengths and weaknesses that we will cover in Section 3. In this paper, we focus on how well causal  
31    estimators perform in the simplest setting, where there is no unobserved confounding, no selection  
32    bias, and no measurement error. It is straightforward to extend RealCause to these more complex  
33    settings. The ideal benchmark for choosing between causal estimators in this setting should have the  
34    following qualities: (1) yield ground-truth estimands, (2) be representative of a substantial subset of  
35    real data, (3) do not have unobserved confounders, and (4) yield many different data distributions of  
36    varying important characteristics (e.g. degree of overlap).

37    (1) is important in order to know which estimators yield estimates closer to the ground-truth. (2) is  
38    important so that we know that estimators that perform well on our benchmark will also perform well

on real datasets that we would apply them to. (3) is important so that we can rule out unobserved confounding as the explanation for an estimator performing poorly. (4) is important because it is unlikely that rankings of causal estimators on a single problem will generalize perfectly to all problems. Rather, we might expect that certain estimators perform better on distributions with certain properties and other estimators perform better on distributions with other specific properties. Existing benchmarks often have 1-3 of the above qualities (Section 3). Our benchmarking framework has all four.

We present a benchmark that simulates data from data generating processes (DGPs) that are statistically indistinguishable from observed real data. We first take the observed pretreatment covariates  $W$  as the only common causes of  $T$  and  $Y$ . Then, we fit generative models  $P_{\text{model}}(T | W)$  and  $P_{\text{model}}(Y | T, W)$  that closely match the real analogs  $P(T | W)$  and  $P(Y | T, W)$ . This allows us to simulate realistic data by first sampling  $W$  from the real data, then sampling  $T$  from  $P_{\text{model}}(T | W)$ , and finally sampling  $Y$  from  $P_{\text{model}}(Y | T, W)$ . Importantly, because we’ve fit generative models to the data, we can sample from *both* interventional distributions  $P_{\text{model}}(Y | \text{do}(T = 1), W)$  and  $P_{\text{model}}(Y | \text{do}(T = 0), W)$ , which means that we have access to ground-truth estimands for our realistic simulated data. That is, the fundamental problem of causal inference isn’t a problem in these DGPs. We then use this realistic simulated data for benchmarking.

## Main contributions

1. RealCause and corresponding realistic benchmarks
2. Application of RealCause to show evidence in favor of selecting hyperparameters based on predictive metrics (like in machine learning)
3. Open-source dataset for predicting causal performance of causal estimators from predictive performance

## 2 Preliminaries and notation

Let  $T$  be a binary scalar random variable denoting the treatment. Let  $W$  be a set of random variables that corresponds to the observed covariates. Let  $Y$  be a scalar random variable denoting the outcome of interest. Let  $e(w)$  denote the *propensity score*  $P(T = 1 | W = w)$ . We denote the treatment and outcome for unit  $i$  as  $T_i$  and  $Y_i$ .  $Y_i(1)$  (resp.  $Y_i(0)$ ) denotes the potential outcome that unit  $i$  would observe if  $T_i$  were 1, taking treatment (resp. if  $T_i$  were 0, not taking treatment).  $Y(t)$  is a random variable that is a function of all the relevant characteristics  $I$  (a set of random variables) that characterize the outcome of an individual (unit) under treatment  $t$ .

We define the *individual treatment effect* (ITE) for unit  $i$  as follows:  $\tau_i \triangleq Y_i(1) - Y_i(0)$ . We define the *average treatment effect* (ATE) as follows:  $\tau \triangleq \mathbb{E}[Y(1) - Y(0)]$ . Let  $C$  be a set of random variables, denoting all the common causes (confounders) of the causal effect of  $T$  on  $Y$ . We can identify the ATE from observational data if we observe  $C$ . This setting has many names: “no unobserved confounding,” “conditional ignorability,” “conditional exchangeability,” “selection on observables,” etc. In this setting, we can identify the ATE via the *adjustment formula* (Robins, 1986; Spirtes et al., 1993; Pearl et al., 2016; Pearl, 2009):

$$\tau = \mathbb{E}_C [\mathbb{E}[Y | T = 1, C] - \mathbb{E}[Y | T = 0, C]] \quad (1)$$

We define the *conditional average treatment effect* (CATE) similarly:

$$\tau(x) \triangleq \mathbb{E}[Y(1) - Y(0) | X = x] = \mathbb{E}_C [\mathbb{E}[Y | T = 1, x, C] - \mathbb{E}[Y | T = 0, x, C]] \quad (2)$$

Here,  $X$  is a set of random variables that corresponds to the characteristics that we are interested in measuring more specialized treatment effects with respect to ( $x$ -specific treatment effects). In this paper, we’ll only consider CATEs where  $X = W$ , so there is no further need for the variable  $X$ .

Similarly, we consider DGPs where  $W = C$ , for simplicity, so it suffices to use only the variable  $W$ . This means that we must adjust for all of  $W$  to get causal effects and that the CATEs reduce to

$$\tau(w) = \mathbb{E}[Y | T = 1, w] - \mathbb{E}[Y | T = 0, w] \triangleq \mu(1, w) - \mu(0, w), \quad (3)$$

where  $\mu$  is the *mean conditional outcome*. Our DGPs provide ground-truth CATEs by providing  $\mu$ . This allows our DGPs to capture unobserved causes of  $Y$  in the data.

### 85 3 Methods for evaluating causal estimators

#### 86 3.1 Simulated synthetic data

87 The simplest way to get ground truth ATEs is to simulate synthetic data that we construct so that  
88 the only confounders of the effect of  $T$  on  $Y$  are  $W$ . This gives us access to the true *outcome*  
89 *mechanism*  $P(Y | T, W)$ . Using the outcome mechanism, we have access to the ground-truth CATE  
90 via Equation 3 and the ground-truth ATE via Equation 1.

91 In these simulations, we additionally have access to the true *treatment selection mechanism*  $P(T |$   
92  $W)$  (or just “*selection mechanism*” for short). We must be able to sample from this to generate  
93 samples from  $P(W, T, Y)$  via ancestral sampling:  $P(W) \rightarrow P(T | W) \rightarrow P(Y | T, W)$ . Having  
94 access to  $P(T | W)$  gives us ground-truth for things like the propensity scores and the degree of  
95 positivity/overlap violations.

96 This is probably the most common method for evaluating causal estimators. However, it has several  
97 disadvantages. First, the data is completely synthetic, so we do not know if the rankings of estimators  
98 that we get will generalize to real data. Second, authors proposing new causal estimators are naturally  
99 interested in synthetic data with specific properties that their estimator was developed to perform  
100 well on. This means that different synthetic data used in different papers cannot be used for a fair  
101 comparison.

#### 102 3.2 Simulated semi-synthetic data with real covariates

103 One natural improvement on the completely synthetic data described in Section 3.1 is to make it  
104 more realistic by taking the covariates  $W$  from real data. This means that  $P(W)$  is realistic. Then,  
105 one can proceed with generating samples through ancestral sampling by simulating  $P(T | W)$  and  
106  $P(Y | T, W)$  as arbitrary stochastic functions. One of the main advantages of this is that these  
107 stochastic functions can be made to have any properties that its designers choose, such as degree  
108 of nonlinearity, positivity violation, treatment effect heterogeneity, etc. (Dorie et al., 2019). This is  
109 what many current benchmarks do (Dorie et al., 2019; Shimoni et al., 2018; Hahn et al., 2019). The  
110 main problem is that the selection mechanism  $P(T | W)$  and outcome mechanism  $P(Y | T, W)$  are  
111 unrealistic.

#### 112 3.3 Simulated data that is fit to real data

113 The way to fix the unrealistic selection and outcome mechanisms is to fit them to real data. This is  
114 what we do, and we are not the first. For example, there is work on this in economics (Knaus et al.,  
115 2018; Athey et al., 2019; Huber et al., 2013; Lechner & Wunsch, 2013), in healthcare (Wendling  
116 et al., 2018; Franklin et al., 2014), and in papers that are meant for a general audience (Abadie &  
117 Imbens, 2011; Schuler et al., 2017). Some fit relatively simple models (Franklin et al., 2014; Abadie  
118 & Imbens, 2011), whereas others fit more flexible models (Wendling et al., 2018; Athey et al., 2019;  
119 Schuler et al., 2017). Our work is distinguished from the above work in two key ways: we statistically  
120 test that our generative models are realistic using two samples tests and we provide knobs to vary  
121 important characteristics of the DGPs. See Appendix A.1 for more discussion on this.

122 **Using RCTs for ground-truth** Finally, there are several different ways to use RCTs for ground-  
123 truths, but they all have problems, which we discuss in Appendix A.2.

### 124 4 RealCause: a method for producing realistic benchmark datasets

125 The basic idea is to fit flexible generative models  $P_{\text{model}}(T | W)$  and  $P_{\text{model}}(Y | T, W)$  to the  
126 selection mechanism  $P(T | W)$  and the outcome mechanism  $P(Y | T, W)$ , respectively. For  
127  $P_{\text{model}}(W)$ , we simply sample from  $P(W)$ , just as is done in the semi-synthetic data simulations we  
128 described in Section 3.2. These three mechanisms give us a joint  $P_{\text{model}}(W, T, Y)$  that we would like  
129 to be the same as the true  $P(W, T, Y)$ . This is what makes our DGPs realistic.

130 **Architecture** We use neural networks to parameterize the conditioning of  $P_{\text{model}}(T | W)$  and  
131  $P_{\text{model}}(Y | T, W)$ ; that is, the input of the neural net is either  $W$  (to predict  $T$ ) or both  $W$  and  
132  $T$  (to predict  $Y$ ). A naive approach would be to concatenate  $W$  and  $T$  to predict the  $Y$ , but our

experiments on semi-synthetic data (where  $\tau$  is known) suggest that the resulting generative model tends to underestimate  $\tau$ . For example, this can happen from the network “ignoring”  $T$ , especially when  $W$  is high-dimensional. Therefore, we follow the TARNet structure (Shalit et al., 2017) to learn two separate conditionals  $P_{\text{model}}(Y \mid T = 0, W)$  and  $P_{\text{model}}(Y \mid T = 1, W)$ , encoding the importance of  $T$  into the structure of our network. Since all conditionals depend on  $W$ , we use a multi-layer perceptron (MLP) to extract common features  $h(W)$  of  $W$ . We then have three more MLPs to model  $T$ ,  $Y \mid T = 0$ , and  $Y \mid T = 1$  separately, taking in the features  $h(W)$  as input. These all use the same  $h(W)$ , which is also learned, like in Dragonnet (Shi et al., 2019). For simplicity, all four MLPs have the same architecture. The tunable hyperparameters are the number of layers, the number of hidden units, and the activation function.

**Distribution assumption** We use the output of the MLPs to parameterize the distributions of selection and outcome. For example, for binary data (such as treatment), we apply the logistic sigmoid activation function to the last layer to parameterize the mean parameter of the Bernoulli distribution. For real-valued data (such as the outcome variable), one option is to assume it follows a Gaussian distribution conditioned on the covariates, in which case we would have the neural net output the mean and log-variance parameters. The baseline model that we use is a linear model that outputs the parameters of a Gaussian distribution with a diagonal covariance matrix. The main (more flexible) generative model we use is the sigmoidal flow (Huang et al., 2018), which has been shown to be a universal density model capable of fitting arbitrary distributions.

For mixed random variables, we parameterize the likelihood as a mixture distribution:  $P(Y) = \pi_0 1_{Y \notin \mathcal{A}} P_c(Y) + \sum_{j=1}^K \pi_j 1_{Y=a_j}$  where  $\mathcal{A} = \{a_1, \dots, a_K\}$  is the set of (discrete) atoms,  $\pi_j$  for  $j = 0, \dots, K$  forms a convex sum, and  $P_c$  is the density function of the continuous component. We have dropped the conditioning to simplify the notation.

**Optimization** For all the datasets, we use a 50/10/40 split for the training set, validation set, and test set. To preprocess the covariate ( $W$ ) and the outcome ( $Y$ ), we either standardize the data to have zero mean and unit variance or normalize it so that the training data ranges from 0 to 1. We use the Adam optimizer to maximize the likelihood of the training data, and save the model with the best validation likelihood for evaluation and model selection. We perform grid search on the hyperparameters and select the model with the best (early-stopped) validation likelihood and with a p-value passing 0.05 on the validation set.

**Tunable knobs** After we fit a generative model to a dataset, we might like to get other models that are very similar but differ along important dimensions of interest. For example, this will allow us to test estimators in settings where there are positivity/overlap violations, where the causal effect is large/small, or where there is a lot of heterogeneity, no heterogeneity, etc. To do this, RealCause supports the following 3 knobs that we can turn to generate new but related distributions, after we’ve fit a model to a real dataset.

**Positivity/overlap knob** Let  $p_i$  be the probability of treatment for example  $i$  (i.e.  $p_i = P(T = 1 \mid W = w_i)$ ). The value of this knob  $\beta$  can be set to anywhere between 0 and 1 inclusive. We use  $\beta$  to linearly interpolate between  $p_i$  and the the extreme that  $p_i$  is closer to (0 or 1). Namely, we change  $p_i$  to  $p'_i$  according to the following equation:  $p'_i = \beta p_i + (1 - \beta) 1_{p_i \geq 0.5}$ . For example,  $\beta = 1$  corresponds to the regular data,  $\beta = 0$  corresponds to the setting where treatment selection is fully deterministic, and all other values of  $0 < \beta < 1$  correspond to somewhere in between.

**Heterogeneity knob** The value  $\gamma$  of the heterogeneity knob can be any real value between 0 and 1 inclusive. If  $\gamma$  is set to 1, the CATEs are the same as the regular dataset. If  $\gamma$  is set to 0, the CATEs are all equal to the ATE. If  $\gamma$  is somewhere between 0 and 1, the CATEs are the corresponding linear interpolation of the original CATE and the ATE.

**Causal effect scale knob** The value  $s$  of the causal effect scale knob can be any real number. This knob sets the scale of the causal effects by changing the potential outcomes according to the following equations:  $Y_i(1)' = s \frac{Y_i(1)}{\tau}$  and  $Y_i(0)' = s \frac{Y_i(0)}{\tau}$ .

## 5 How realistic is RealCause?

In this section, we show that RealCause produces realistic datasets that are very close to the real ones. For all datasets, we show that the distribution of our generative model  $P_{\text{model}}(W, T, Y)$  is

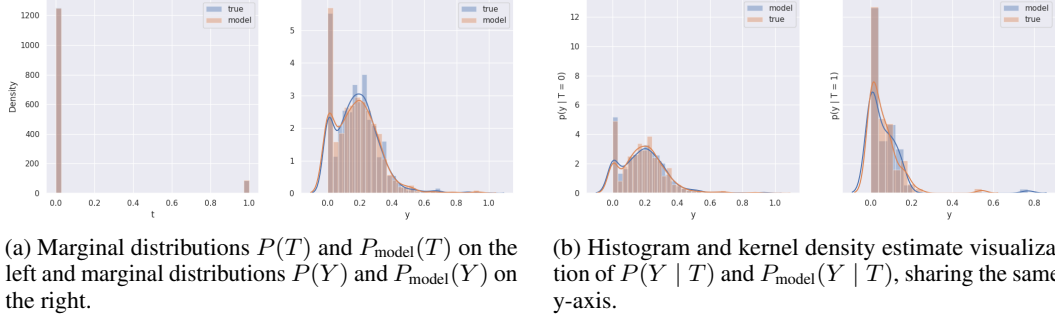


Figure 1: Visualizations of how well the generative model models the real LaLonde PSID data.

very close to the true distribution  $P(W, T, Y)$ . We show this by providing both visual comparisons and quantitative evaluations. We visually compare  $P_{\text{model}}(T, Y)$  and  $P(T, Y)$  using histograms and Gaussian kernel density estimation (see, e.g., Figure 1). We quantitatively compare  $P_{\text{model}}(W, T, Y)$  and  $P(W, T, Y)$  by running two-sample tests (Table 1).

Two-sample tests evaluate the probability that a sample from  $P_{\text{model}}(W, T, Y)$  and a sample from  $P(W, T, Y)$  came from the same distribution, under the null hypothesis that  $P_{\text{model}}(W, T, Y) = P(W, T, Y)$  (that the model distribution matches the true distribution). If that probability (p-value) is less than some small value  $\alpha$  such as 0.05, we say we have sufficient evidence to reject the null hypothesis that  $P_{\text{model}}(W, T, Y) = P(W, T, Y)$  (i.e. the generative model is not as realistic as we would like it to be). This is how we operationalize the hypothesis that our modeled distributions are “realistic.” Two-sample tests give us a way to falsify the hypothesis that our generative models are realistic.

However, two-sample tests do not work well in high-dimensions. Importantly, the power<sup>1</sup> of two-sample tests can decay with dimensionality (Ramdas et al., 2015) and  $W$  can have many dimensions in the datasets we consider. On the bright side, the treatment  $T$  and the outcome  $Y$  are each one-dimensional, so evaluating the statistical relationship between them is only a two-dimensional problem. This means that we might get more power from testing the hypothesis that  $P_{\text{model}}(T, Y) = P(T, Y)$  because it’s a lower-dimensional problem, even though this test will ignore  $W$  and its relationship to  $T$  and  $Y$ . Tests that use  $P(W, T, Y)$  could have more power because they use information about  $P(T, Y | W)$  (recall that  $P(W) = P_{\text{model}}(W)$ , by construction). Therefore, we run two-sample tests for both  $P(T, Y)$  and  $P(W, T, Y)$  (and the marginals). Finally, we stress that passing the marginal tests is not trivial, since we learn the conditional  $P(T, Y | W)$  and marginalize out  $P(W)$ , instead of learning  $P(T, Y)$ ,  $P(T)$ , or  $P(Y)$  directly.

**Datasets** We fit eight datasets in total. We fit generative models to three real datasets: LaLonde PSID, LaLonde CPS (LaLonde, 1986) (we use Dehejia & Wahba (1999)’s version), and Twins<sup>2</sup> (Louizos et al., 2017). We additionally fit generative models to five popular semi-synthetic datasets: IHDP (Hill, 2011) and four LBIDD datasets (Shimoni et al., 2018). On all of these datasets, we can fit generative models to model the observational distribution. Then, with the semi-synthetic datasets, we can also check that our generative models give roughly the same ground-truth causal effects as existing popular synthetic benchmarks.

**Visualization of modeled LaLonde PSID** Consider the LaLonde PSID dataset as our first example. We visualize  $P_{\text{model}}(T)$  vs.  $P(T)$  and  $P_{\text{model}}(Y)$  vs.  $P(Y)$  in Figure 1a.  $P_{\text{model}}(W)$  and  $P(W)$  are known to be the same distributions, by construction. We visualize  $P_{\text{model}}(T, Y)$  vs.  $P(T, Y)$  in Figure 1b. We provide similar visualizations of the other real datasets and corresponding similar models in Appendix B.

**Univariate statistical tests** The Kolmogorov-Smirnov (KS) test is the most popular way to test the hypothesis that two samples come from the same distribution. The Epps-Singleton (ES) test is more well-suited for discrete distributions and can have higher power than the KS test (Epps &

<sup>1</sup>For a fixed value of  $\alpha$ , power is the probability of rejecting the null hypothesis, given that the null hypothesis is false.

<sup>2</sup>The treatment selection mechanism for the Twins dataset is simulated. This is to ensure that there is some confounding, as the regular dataset might be unconfounded.



Singleton, 1986). We use the implementations of the KS and ES tests from *SciPy* (Virtanen et al., 2020). For all datasets, we report the p-values of the KS and ES tests for comparing the marginal distributions  $P_{\text{model}}(Y)$  and  $P(Y)$  and for comparing the marginal distributions  $P_{\text{model}}(T)$  and  $P(T)$  in the first section of Table 1. In all tests, the p-values are much larger than any reasonable value of  $\alpha$ , so we fail to reject the null hypothesis that the generated data and the true data come from the same distribution. This means that our generative models are reasonably realistic, at least if we only look at the marginals.

**Multivariate statistical tests** Extending the KS test to multiple dimensions is difficult. However, there are several multivariate tests such as the Friedman-Rafsky test (Friedman & Rafsky, 1979), k-nearest neighbor (kNN) test (Friedman & Rafsky, 1983), and energy test (Székely & Rizzo, 2013). We use the implementations of these tests in the *torch-two-sample* Python library (Djoulonga, 2017). These are just permutation tests and can be conducted with any statistic, so we additionally run permutation tests with the Wasserstein-1 and Wasserstein-2 distance metrics. We run each test with 1000 permutations. We display the corresponding p-values in the last two sections of Table 1. For all tests except the FR and kNN  $(T, Y)$  test on the LaLonde PSID dataset, the p-values are much larger than any reasonable value of  $\alpha$ . However, we might be worried that these multivariate two-sample tests don't have enough power when we include the higher-dimensional  $W$ .

**Demonstration of statistical power via linear baselines** We demonstrate that these tests do have a decent amount of statistical power (probability of rejecting the null when  $P_{\text{model}}$  and  $P$  differ) by fitting a linear Gaussian model to the data and displaying the corresponding p-values in Table 2. Even when  $W$  is high-dimensional, we are still able to reject the linear models as realistic. For example, we clearly have p-values that are below most reasonable values of  $\alpha$  for the LaLonde PSID, and all three nonlinear LBIDD datasets. As we might expect, for high-dimensional  $W$  such as in the LBIDD datasets, the  $(T, Y)$  tests have enough power to reject the null hypothesis because they operate in only two dimensions, whereas the  $(W, T, Y)$  tests do not because their power suffers from the high-dimensionality (179 dimensions). The LaLonde CPS dataset is an example where it can be useful to include  $W$  in the statistical test; all of the p-values for the  $(T, Y)$  tests are above  $\alpha = .075$ , whereas all but one of the p-values for the  $(W, T, Y)$  tests are below  $\alpha = .075$ . Our p-values for the Twins dataset are quite high, but this is not due to these tests not having enough power. Rather, it is because the Twins dataset is well modeled by a linear model:  $T$  and  $Y$  are both binary (two parameters) and  $W$  is 75-dimensional, so it makes sense that we can linearly predict these two parameters from 75 dimensions. We demonstrate how well the linear model fits Twins in

Table 1: Table of p-values for the various statistical hypothesis tests we run to test the null hypothesis that real data samples and samples from the generative model come from the same distribution. Large values (e.g.  $> 0.05$ ) mean that we don't have statistically significant evidence that the real and generated data come from different distributions, so we want to see large values. The first section is univariate tests. The second section is 2-dimensional tests to capture the dependence of  $Y$  on  $T$ . The third section can be much higher dimensional tests whose power may suffer from the high dimensionality, but these tests may be able to pick up on the dependence of  $T$  and  $Y$  on  $W$  that the 2-dimensional tests cannot pick up on.

TEST	LALONDE		TWINS	IHDP	LBIDD			
	PSID	CPS			QUAD	EXP	LOG	LINEAR
$T$ KS	0.9995	1.0	0.9837	0.9290	0.5935	0.9772	0.4781	0.3912
$T$ ES	0.6971	0.3325	0.7576	0.5587	0.8772	0.6975	0.4157	0.3815
$Y$ KS	0.4968	1.0	0.8914	0.3058	0.2204	0.9146	0.4855	0.4084
$Y$ ES	0.3069	0.1516	0.4466	0.3565	0.2264	0.7223	0.3971	0.1649
$(T, Y)$ Wass1	0.6914	0.435	0.5088	0.2894	0.3617	0.4391	0.3899	0.5046
$(T, Y)$ Wass2	0.6638	0.4356	0.4960	0.3365	0.4353	0.4709	0.4205	0.5063
$(T, Y)$ FR	0.0	0.4004	0.5549	0.4761	0.8610	0.5773	0.5132	0.8355
$(T, Y)$ kNN	0.0	0.4120	0.4318	0.5978	0.3166	0.3735	0.4902	0.4838
$(T, Y)$ Energy	0.6311	0.4396	0.5053	0.3186	0.2371	0.4453	0.3988	0.5086
$(W, T, Y)$ Wass1	0.4210	0.3854	0.4782	1.0	0.5191	0.4219	0.4866	0.5393
$(W, T, Y)$ Wass2	0.5347	0.3660	0.4728	1.0	0.5182	0.4160	0.4807	0.5381
$(W, T, Y)$ FR	0.2569	0.4033	0.5068	1.0	0.4829	0.4989	0.5027	0.4893
$(W, T, Y)$ kNN	0.2270	0.4343	0.4919	1.0	0.5104	0.5101	0.5223	0.4988
$(W, T, Y)$ Energy	0.5671	0.4177	0.5263	0.9409	0.5104	0.4423	0.5031	0.5421
$ W $ (n covariates)	8	8	75	25	177	177	177	177

Table 2: Table of p-values for the various statistical hypothesis tests we run to test the null hypothesis that real data samples and samples from a *linear* Gaussian generative model come from the same distribution. Small values (e.g.  $< 0.05$ ) mean that these tests have enough power to detect that the real data comes from a different distribution than the distribution generated by our linear Gaussian generative model.

TEST	LALONDE				LBIDD			
	PSID	CPS	TWINS	IHDP	QUAD	EXP	LOG	LINEAR
$(T, Y)$ Wass1	0.0304	0.1500	0.5004	0.2019	0.2009	0.0456	0.1510	0.2832
$(T, Y)$ Wass2	0.0123	0.0797	0.4924	0.1636	0.4277	0.1314	0.2380	0.3172
$(T, Y)$ FR	0.0	0.0776	0.5581	0.2825	0.0	0.0014	0.0140	0.7946
$(T, Y)$ kNN	0.0	0.1808	0.4541	0.4183	0.0	0.0023	0.0013	0.4070
$(T, Y)$ Energy	0.0482	0.1620	0.5094	0.2249	0.0002	0.0551	0.2020	0.3409
$(W, T, Y)$ Wass1	0.0470	0.0671	1.0	1.0	0.4917	0.5245	0.8230	0.6777
$(W, T, Y)$ Wass2	0.4001	0.0624	0.9966	1.0	0.4782	0.5204	0.7840	0.6257
$(W, T, Y)$ FR	0.1333	0.0525	0.9992	1.0	0.7655	0.6979	0.3651	0.7369
$(W, T, Y)$ kNN	0.5136	0.0711	1.0	1.0	0.8953	0.8416	0.4510	0.7968
$(W, T, Y)$ Energy	0.1080	0.2863	0.7389	0.8935	0.5099	0.5142	0.7429	0.7144
$ W $ (n covariates)	8	8	75	25	177	177	177	177

Figures 5c and 5d in Appendix B. Similarly, the p-values for IHDP are so high because the IHDP data is reasonably well fit by the linear model (see Figures 6d to 6f), and the IHDP tests have less power since the IHDP dataset is much smaller than the other datasets.

**Realistic causal effects** We also show that our generative model admits causal effect estimates that roughly match those of the popular semi-synthetic benchmarks IHDP and LBIDD. For each of these datasets, we report the true ATE, our generative model’s ATE estimate, the corresponding absolute bias, and the PEHE. We report these values in Table 3. The values in the table indicate that our model accurately models the causal effects. The one number that is relatively high relative to the others is the PEHE for IHDP; this is because the training sample for IHDP is only 374 examples.

**Limitations** Although we can statistically test how well RealCause fits the observed distribution  $P(W, T, Y)$ , we cannot test how well RealCause fits the interventional distributions  $P(Y \mid \text{do}(T = t), w)$  without making the no unobserved confounding assumption. Due to the fundamental problem of causal inference, there is no way of getting around this for arbitrary distributions. Fortunately, we can test the interventional distributions of synthetic data such as IHDP and LBIDD; this is why we include Table 3. That said, RealCause (or any realistic benchmark) could potentially not model the interventional distributions well on other datasets, resulting in suboptimal interventional distributions. Additionally, RealCause will be biased based on the specific architecture of the generative model it uses. Ideally, one would run RealCause benchmarks using many different generative model architectures.

## 6 Results

The reason we spent so much effort establishing that RealCause DGPs are realistic in Section 5 is that we can now trust the results that RealCause DGPs yield for important tasks such as the following: (a) benchmarking causal estimators and (b) evaluating whether *predictive* metrics can be used for model selection of *causal* estimators. We first apply RealCause to benchmarking causal estimators. We then use these results to analyze correlation between predictive performance and causal performance in Section 6.1.

Table 3: True causal effects, corresponding estimates from our generative model, and associated error.

	IHDP	LBIDD QUAD	LBIDD EXP	LBIDD LOG	LBIDD LINEAR
True ATE	4.0161	2.5437	-0.6613	0.0549	1.8592
ATE estimate	4.1908	2.4910	-0.6608	0.0555	1.7177
ATE abs bias	0.1747	0.0527	0.0004	0.0005	0.1415
PEHE	51.5279	0.1554	0.0225	0.0151	0.1367

**Datasets and estimators** In our evaluations in this section, we use 3 real datasets, 4 meta-estimators, 15 machine learning models for each of the meta-estimators, and roughly 10 different settings of the single most important hyperparameter for each of the machine learning models. Taking the Cartesian product over all of those yields over 1500 causal estimators. The 3 datasets we use are LaLonde PSID, LaLonde CPS, and Twins; we use RealCause to turn these into datasets where we know the ground-truth causal effects. The 4 meta-estimators from *causalib* (Shimoni et al., 2019) we use are standardization (or S-learner), stratified standardization (or T-learner), inverse probability weighting (IPW), and IPW with weight trimming. We use a variety of machine learning models from *scikit-learn* (Pedregosa et al., 2011) to plug in to these meta-estimators. For each model, we use a grid of values for the most important hyperparameter (according to van Rijn & Hutter (2018)). See Appendix E for more info on our estimators.

**Benchmarking causal estimators** As one would expect, different causal estimators perform better on different datasets. We choose causal estimators within a given model class according to the best cross-validated RMSE for standardization estimators and according to the best cross-validated average precision for IPW estimators. We divide the ATE RMSEs by each dataset’s ATE and show those weighted averages in Figure 2. Interestingly, most of our standardization estimators don’t perform very well, but then standardization pair with an RBF-SVM achieves the lowest ATE RMSE. While this estimator also achieves the lowest weighted averaged PEHE, it doesn’t have the lowest weighted averaged absolute bias. We provide the corresponding plots for ATE absolute bias and PEHE along with the more fine-grained full tables by dataset in Appendix C.

## 6.1 Predicting causal performance from predictive performance

The following is known and commonly stated: just because the model(s) used in a causal estimator are highly predictive does not mean that the causal estimator will perform well at estimating a causal parameter such as  $\tau$  or  $\tau(w)$ . Then, the following questions naturally arise: (1) How can I choose hyperparameters for causal estimators? (2) How can I inform model selection for causal problems? In machine learning, the answer is simple: run cross-validation using the relevant predictive metric for hyperparameter and model selection. However, we can’t do the analog in causal inference because we don’t have access to a corresponding causal metric, due to the fundamental problem of causal inference.

What if it turns out that the hyperparameters and models that yield the best predictive performance also yield the best causal performance? Then, hyperparameter and model selection for causal inference would be the same as it is for machine learning. We can measure if this is the case by measuring how correlated predictive metrics and causal metrics are.

**Correlation measures** While Pearson’s correlation coefficient is the most common method for measuring correlation, it only captures linear relationships. We are more interested in general monotonic relationships (e.g. if the prediction performance of model A is better than the predictive performance of model B, then will the causal performance of model A also be better than the causal performance of model B?). Therefore, we use Spearman’s rank correlation coefficient (equivalent to Pearson’s correlation coefficient on the *rank* of the random variables) and Kendall’s rank correlation coefficient. We also report a more intuitive measure: the probability that the causal performance of model A is at least as good as the causal performance of model B, given that the predictive performance of model A is at least as good as the predictive performance of model B.

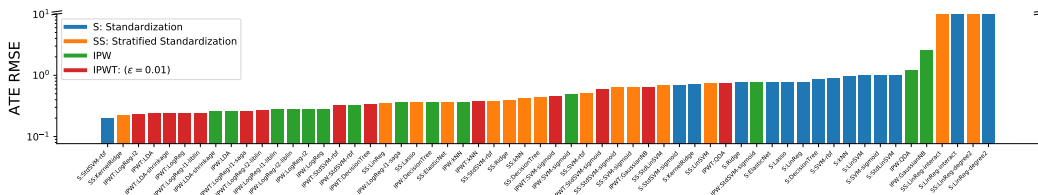


Figure 2: ATE RMSE of the different estimators, weighted averaged (by their inverse ATEs) over three datasets and color-coded by meta-estimator.



**Metrics** The main predictive metric we consider for outcome models (predict  $Y$  from  $T$  and  $W$ ) is the corresponding RMSE (root mean squared error). The main predictive metrics we consider for propensity score models (predict binary  $T$  from  $W$ ) are scikit-learn’s balanced F score, average precision, and balanced accuracy. The main causal metric is the PEHE (Hill, 2011).

**Selecting model hyperparameters** For a given dataset, meta-estimator, and machine learning model class, we must choose the hyperparameters for that specific model class. We show the full table of correlation coefficients for how predictive RMSE is of ATE RMSE and PEHE within every model class in Appendix D.1. We summarize this with just the median Spearman correlation coefficient and the median probability of better or equal causal performance given better or equal predictive performance in Table 4; these medians are taken over all models for standardization and stratified standardization estimators fit to a given dataset. **Importantly, these results show that, in this setting, it is a fairly good idea to select hyperparameters for causal estimators based on predictive performance.** For example, the median probabilities that a better predictive model corresponds to a better causal model hover around 80-95% in this summary table. We do the same for IPW and propensity score models in Appendix D.2.

Table 4: Median correlation of predictive RMSE with PEHE in standardization estimators.

DATASET	SPEARMAN	PROB BETTER
PSID	0.92	0.92
CPS	0.80	0.87
Twins	0.91	0.96

**Model selection** We just saw that predictive performance is indicative of causal performance when choosing hyperparameters within a model class, but what about selecting between model classes after choosing hyperparameters via predictive cross-validation? The results are much less positive and more dataset-specific. For standardization estimators, there isn’t much correlation on the LaLonde datasets, but there is a great deal of correlation on the Twins dataset. For IPW estimators, it is roughly the same, except for the fact that average precision has a modest correlation with ATE RMSE on the LaLonde CPS dataset. See Appendix D.3 for details.

**Open-source dataset for exploration** We created a dataset with 1568 rows (estimators) and 77 columns (predictive metrics, causal metrics, and estimator specification). Importantly, this dataset contains all the predictive metrics that scikit-learn provides and many different causal metrics that we compute using RealCause. In this section, we chose one line of analysis for this dataset, but there are many others. For example, one can use any machine learning model for predicting any subset of causal metrics from any subset of predictive metrics, one can cross-validate over different predictive metrics than the ones we used, one can group the data differently, etc. We already see that different predictive metrics correlate quite differently with ATE RMSE, depending on the model and dataset in Appendix D.2. This suggests that more value might be gained in doing more complex analyses on this dataset. We open-source our dataset at <https://github.com/bradyneal/causal-benchmark/blob/master/causal-predictive-analysis.csv>.

## 7 Conclusion and future work

Now that we’ve rigorously shown that RealCause produces realistic DGPs, we are hopeful that others will use it. We open-source default benchmark datasets, our trained RealCause generative models, and the code to train new generative models on other datasets at <https://github.com/bradyneal/causal-benchmark>.

There are many important extensions of RealCause that can be done. Adding even more causal estimators and more real datasets would be valuable to expand the open-source dataset of predictive and causal metrics that we started. Similarly, running the benchmarking suite with various non-default settings of RealCause’s knobs (e.g. zero overlap) could lead to useful empirical results about when to use various estimators. RealCause’s realism gives us confidence in our evidence that hyperparameters for causal estimators can be selected using cross-validation on a predictive metric. There is much potential for further analysis of our open-source dataset of predictive and causal metrics. For example, future papers or a Kaggle competition to predict causal metrics from predictive metrics would be valuable.

## Acknowledgements

We thank Uri Shalit, Yoshua Bengio, and Ioannis Mitliagkas for useful feedback on this paper.

## References

- Abadie, A. and Imbens, G. W. Bias-corrected matching estimators for average treatment effects. *Journal of Business & Economic Statistics*, 29(1):1–11, 2011.
- Arjovsky, M., Chintala, S., and Bottou, L. Wasserstein generative adversarial networks. In Precup, D. and Teh, Y. W. (eds.), *Proceedings of the 34th International Conference on Machine Learning*, volume 70 of *Proceedings of Machine Learning Research*, pp. 214–223, International Convention Centre, Sydney, Australia, 06–11 Aug 2017. PMLR.
- Athey, S., Imbens, G., Metzger, J., and Munro, E. Using wasserstein generative adversarial networks for the design of monte carlo simulations, 2019.
- Dehejia, R. H. and Wahba, S. Causal effects in nonexperimental studies: Reevaluating the evaluation of training programs. *Journal of the American Statistical Association*, 94(448):1053–1062, 1999.
- Djulonga, J. A pytorch library for differentiable two-sample tests. <https://github.com/josipd/torch-two-sample>, 2017.
- Dorie, V., Hill, J., Shalit, U., Scott, M., and Cervone, D. Automated versus do-it-yourself methods for causal inference: Lessons learned from a data analysis competition. *Statist. Sci.*, 34(1):43–68, 02 2019.
- Epps, T. and Singleton, K. J. An omnibus test for the two-sample problem using the empirical characteristic function. *Journal of Statistical Computation and Simulation*, 26(3-4):177–203, 1986.
- Franklin, J., Schneeweiss, S., Polinski, J., and Rassen, J. Plasmode simulation for the evaluation of pharmacoepidemiologic methods in complex healthcare databases. *Computational Statistics & Data Analysis*, 72:219–226, 04 2014.
- Friedman, J. H. and Rafsky, L. C. Multivariate generalizations of the wald-wolfowitz and smirnov two-sample tests. *Ann. Statist.*, 7(4):697–717, 07 1979.
- Friedman, J. H. and Rafsky, L. C. Graph-theoretic measures of multivariate association and prediction. *Ann. Statist.*, 11(2):377–391, 06 1983.
- Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., and Bengio, Y. Generative adversarial nets. In Ghahramani, Z., Welling, M., Cortes, C., Lawrence, N. D., and Weinberger, K. Q. (eds.), *Advances in Neural Information Processing Systems 27*, pp. 2672–2680. Curran Associates, Inc., 2014.
- Hahn, P. R., Dorie, V., and Murray, J. S. Atlantic causal inference conference (acic) data analysis challenge 2017, 2019.
- Hernán, M. A. and Robins, J. M. *Causal Inference: What If*. Boca Raton: Chapman & Hall/CRC, 2020.
- Hill, J. L. Bayesian nonparametric modeling for causal inference. *Journal of Computational and Graphical Statistics*, 20(1):217–240, 2011.
- Hill, J. L., Reiter, J. P., and Zanutto, E. L. *A Comparison of Experimental and Observational Data Analyses*, chapter 5, pp. 49–60. John Wiley & Sons, Ltd, 2004.
- Holland, P. W. Statistics and causal inference. *Journal of the American Statistical Association*, 81(396):945–960, 1986.
- Horvitz, D. G. and Thompson, D. J. A generalization of sampling without replacement from a finite universe. *Journal of the American Statistical Association*, 47(260):663–685, 1952.

- 413 Huang, C.-W., Krueger, D., Lacoste, A., and Courville, A. Neural autoregressive flows. In *International Conference on Machine Learning*, pp. 2078–2087, 2018.
- 414
- 415 Huber, M., Lechner, M., and Wunsch, C. The performance of estimators based on the propensity  
416 score. *Journal of Econometrics*, 175(1):1 – 21, 2013.
- 417 Imbens, G. W. and Rubin, D. B. *Causal Inference for Statistics, Social, and Biomedical Sciences: An  
418 Introduction*. Cambridge University Press, 2015.
- 419 Kallus, N., Puli, A. M., and Shalit, U. Removing hidden confounding by experimental grounding.  
420 In Bengio, S., Wallach, H., Larochelle, H., Grauman, K., Cesa-Bianchi, N., and Garnett, R.  
421 (eds.), *Advances in Neural Information Processing Systems*, volume 31, pp. 10888–10897. Curran  
422 Associates, Inc., 2018.
- 423 Knaus, M. C., Lechner, M., and Strittmatter, A. Machine learning estimation of heterogeneous causal  
424 effects: Empirical monte carlo evidence, 2018.
- 425 Künzel, S. R., Sekhon, J. S., Bickel, P. J., and Yu, B. Metalearners for estimating heterogeneous  
426 treatment effects using machine learning. *Proceedings of the National Academy of Sciences*, 116  
427 (10):4156–4165, 2019.
- 428 LaLonde, R. J. Evaluating the econometric evaluations of training programs with experimental data.  
429 *The American Economic Review*, 76(4):604–620, 1986.
- 430 Lechner, M. and Wunsch, C. Sensitivity of matching-based program evaluations to the availability of  
431 control variables. *Labour Economics*, 21:111 – 121, 2013.
- 432 Louizos, C., Shalit, U., Mooij, J. M., Sontag, D., Zemel, R., and Welling, M. Causal effect inference  
433 with deep latent-variable models. In Guyon, I., Luxburg, U. V., Bengio, S., Wallach, H., Fergus,  
434 R., Vishwanathan, S., and Garnett, R. (eds.), *Advances in Neural Information Processing Systems  
435 30*, pp. 6446–6456. Curran Associates, Inc., 2017.
- 436 Morgan, S. L. and Winship, C. *Counterfactuals and Causal Inference: Methods and Principles for  
437 Social Research*. Analytical Methods for Social Research. Cambridge University Press, 2 edition,  
438 2014.
- 439 Neal, B. *Introduction to Causal Inference*. 2020.
- 440 Pearl, J. A probabilistic calculus of actions. *ArXiv*, abs/1302.6835, 1994.
- 441 Pearl, J. *Causality*. Cambridge University Press, 2009.
- 442 Pearl, J. On the interpretation of do(x). *Journal of Causal Inference*, 7(1):20192002, 2019.
- 443 Pearl, J., Glymour, M., and Jewell, N. P. *Causal inference in statistics: A primer*. John Wiley & Sons,  
444 2016.
- 445 Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M.,  
446 Prettenhofer, P., Weiss, R., Dubourg, V., Vanderplas, J., Passos, A., Cournapeau, D., Brucher, M.,  
447 Perrot, M., and Duchesnay, E. Scikit-learn: Machine learning in Python. *Journal of Machine  
448 Learning Research*, 12:2825–2830, 2011.
- 449 Ramdas, A., Reddi, S. J., Póczos, B., Singh, A., and Wasserman, L. On the decreasing power of  
450 kernel and distance based nonparametric hypothesis tests in high dimensions. In *Proceedings of  
451 the Twenty-Ninth AAAI Conference on Artificial Intelligence, AAAI’15*, pp. 3571–3577. AAAI  
452 Press, 2015.
- 453 Robins, J. A new approach to causal inference in mortality studies with a sustained exposure  
454 period—application to control of the healthy worker survivor effect. *Mathematical Modelling*, 7  
455 (9):1393 – 1512, 1986.
- 456 Rubin, D. B. Estimating causal effects of treatments in randomized and nonrandomized studies.  
457 *Journal of educational Psychology*, 66(5):688, 1974.

458 Schuler, A., Jung, K., Tibshirani, R., Hastie, T., and Shah, N. Synth-validation: Selecting the best  
459 causal inference method for a given dataset, 2017.

460 Shadish, W. R., Clark, M. H., and Steiner, P. M. Can nonrandomized experiments yield accurate  
461 answers? a randomized experiment comparing random and nonrandom assignments. *Journal of*  
462 *the American Statistical Association*, 103(484):1334–1344, 2008.

463 Shalit, U., Johansson, F. D., and Sontag, D. Estimating individual treatment effect: generalization  
464 bounds and algorithms. In *International Conference on Machine Learning*, pp. 3076–3085. PMLR,  
465 2017.

466 Shi, C., Blei, D., and Veitch, V. Adapting neural networks for the estimation of treatment effects. In  
467 *Advances in Neural Information Processing Systems*, pp. 2507–2517, 2019.

468 Shimoni, Y., Yanover, C., Karavani, E., and Goldschmidt, Y. Benchmarking Framework for  
469 Performance-Evaluation of Causal Inference Analysis. *ArXiv preprint arXiv:1802.05046*, 2018.

470 Shimoni, Y., Karavani, E., Ravid, S., Bak, P., Ng, T. H., Alford, S. H., Meade, D., and Goldschmidt,  
471 Y. An evaluation toolkit to guide model selection and cohort definition in causal inference, 2019.

472 Snowden, J. M., Rose, S., and Mortimer, K. M. Implementation of G-Computation on a Simulated  
473 Data Set: Demonstration of a Causal Inference Technique. *American Journal of Epidemiology*,  
474 173(7):731–738, 03 2011.

475 Spirtes, P., Glymour, C., and Scheines, R. *Causation, Prediction, and Search*, volume 81. 01 1993.

476 Székely, G. J. and Rizzo, M. L. Energy statistics: A class of statistics based on distances. *Journal of*  
477 *Statistical Planning and Inference*, 143(8):1249 – 1272, 2013.

478 Turner, R. and Neal, B. How well does your sampler really work? In *Uncertainty in Artificial*  
479 *Intelligence*. AUAI Press, 2018.

480 van Rijn, J. N. and Hutter, F. *Hyperparameter Importance Across Datasets*, pp. 2367–2376. Associa-  
481 tion for Computing Machinery, New York, NY, USA, 2018.

482 Virtanen, P., Gommers, R., Oliphant, T. E., Haberland, M., Reddy, T., Cournapeau, D., Burovski, E.,  
483 Peterson, P., Weckesser, W., Bright, J., van der Walt, S. J., Brett, M., Wilson, J., Millman, K. J.,  
484 Mayorov, N., Nelson, A. R. J., Jones, E., Kern, R., Larson, E., Carey, C. J., Polat, İ., Feng, Y.,  
485 Moore, E. W., VanderPlas, J., Laxalde, D., Perktold, J., Cimrman, R., Henriksen, I., Quintero,  
486 E. A., Harris, C. R., Archibald, A. M., Ribeiro, A. H., Pedregosa, F., van Mulbregt, P., and SciPy  
487 1.0 Contributors. SciPy 1.0: Fundamental Algorithms for Scientific Computing in Python. *Nature*  
488 *Methods*, 17:261–272, 2020.

489 Wendling, T., Jung, K., Callahan, A., Schuler, A., Shah, N. H., and Gallego, B. Comparing  
490 methods for estimation of heterogeneous treatment effects using observational data from health  
491 care databases. *Statistics in Medicine*, 37(23):3309–3324, 2018.

## Checklist

### 1. For all authors...

- (a) Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope? [Yes]
- (b) Did you describe the limitations of your work? [Yes]
- (c) Did you discuss any potential negative societal impacts of your work? [No]
- (d) Have you read the ethics review guidelines and ensured that your paper conforms to them? [Yes]

### 2. If you are including theoretical results...

- (a) Did you state the full set of assumptions of all theoretical results? [N/A]
- (b) Did you include complete proofs of all theoretical results? [N/A]

### 3. If you ran experiments (e.g. for benchmarks)...

- (a) Did you include the code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL)? [Yes]
- (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? [Yes]
- (c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? [No]
- (d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? [No]

### 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...

- (a) If your work uses existing assets, did you cite the creators? [Yes]
- (b) Did you mention the license of the assets? [Yes]
- (c) Did you include any new assets either in the supplemental material or as a URL? [Yes]
- (d) Did you discuss whether and how consent was obtained from people whose data you're using/curating? [No] All datasets we use are well-known datasets in causal inference.
- (e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? [No] All datasets we use are well-known datasets in causal inference.

### 5. If you used crowdsourcing or conducted research with human subjects...

- (a) Did you include the full text of instructions given to participants and screenshots, if applicable? [N/A]
- (b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [N/A]
- (c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [N/A]



# Appendices

## Appendix A Further Related Work Discussion

### A.1 More on Simulated Data That is Fit to Real Data

Our work is distinguished from the work in Section 3.3 in key two ways: we statistically test that our generative models are realistic using two samples tests and we provide knobs to vary important characteristics of that data. See Appendix A.1 for more discussion on this.

First, most of the above work does not use modern generative models that can fit most distributions so well that they pass two-sample tests.<sup>3</sup> Athey et al. (2019) is the exception as they use conditional Wasserstein Generative Adversarial Networks (WGANs) (Goodfellow et al., 2014; Arjovsky et al., 2017) for both the outcome mechanism and the reverse of the selection mechanism:  $P(W | T)$ . However, they do not report two-sample tests to rigorously test the hypothesis that their generative model is statistically similar to the distributions they are fit to. We fit several datasets well and run two-sample tests to test this claim in Section 5.

Second, our method allows us to have “knobs” to vary important aspects of the data distributions, just as Dorie et al. (2019) are able to maintain in their semi-synthetic study where they specify random functions for the outcome mechanism and selection mechanism. Wendling et al. (2018) illustrate the nontriviality of this when they wrote “The design of a simulation study is usually a trade-off between realism and control.” We are able to get both realism *and* control.

### A.2 Using RCTs for Ground-Truth

**Constructed observational studies** One can first take a randomized control trial (RCT), and get an unbiased estimate of the ground-truth ATE from that. Then, one can construct a corresponding observational study by swapping out the RCT control group with observational data for people who were not part of the RCT control group. LaLonde (1986) was the first to do this. We refer to this type of study as a *constructed observational study* (a term coined by Hill et al. (2004)). There are two problems with this type of study: (1) We do not know if we have observed all of the confounders for the observational data, so we do not know if an estimate that differs from the RCT ATE is due to unobserved confounding or due to the estimator doing poorly regardless of unobserved confounding. (2) The population that the observational data comes from is often not the same population as the population that the RCT data comes from, so it is not clear if the RCT ATE is the same as the true ATE of the constructed observational data.

**Doubly randomized preference trials (DRPTs)** In a *double randomized preference trial* (DRPT), one runs an RCT and an observational study in parallel on the same population. This can be done by first randomizing units into the RCT or observational study. The units in the RCT are then randomized into treatment groups. The units in the observational study are allowed to select which treatment they take. Shadish et al. (2008) were the first to run a DRPT to evaluate observational methods. The main problem with DRPTs (which is also a problem for constructed observational studies) is that you only get a single DGP, and it is prohibitively expensive to run many different DRPTs, which is important because we do not expect that the rankings of estimators will be the same across all DGPs. Additionally, if a causal estimator performs poorly, we do not know if it is simply because of unobserved confounding.

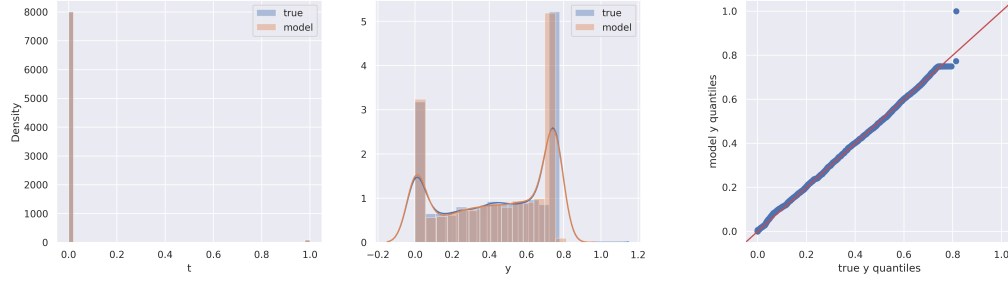
**Introducing selection bias via selective subsampling** While  $\hat{\mathbb{E}}[Y | T = 1] - \hat{\mathbb{E}}[Y | T = 0]$  is an unbiased estimate of the ATE in an RCT, we can turn this into a biased estimate if we introduce selection bias (see, e.g., Kallus et al. (2018)). We can introduce selection bias by selectively subsampling the data based on  $T$  and  $Y$  and giving that subsampled data to the causal estimator. Graphically, this creates a collider  $C$  that is a child of both  $T$  and  $Y$  in the causal graph. And we’re conditioning on this collider by giving the estimator access to only the subsampled data. This introduces selection bias (see, e.g., Hernán & Robins (2020, Chapter 8)), which means that  $\hat{\mathbb{E}}[Y | T = 1] - \hat{\mathbb{E}}[Y | T = 0]$  is a biased estimate in the subsampled data. The two problems with

<sup>3</sup>In related work outside of causal inference, Turner & Neal (2018) applied modern generative models and two-sample tests for benchmarking Markov chain Monte Carlo (MCMC) samplers.

576 this approach are (a) the selection mechanism is chosen by humans, so it may not be realistic, and  
577 (b) the graphical structural of selection bias is different from the graphical structure of confounding  
578 (common effect of  $T$  and  $Y$  vs. common cause of  $T$  and  $Y$ ).

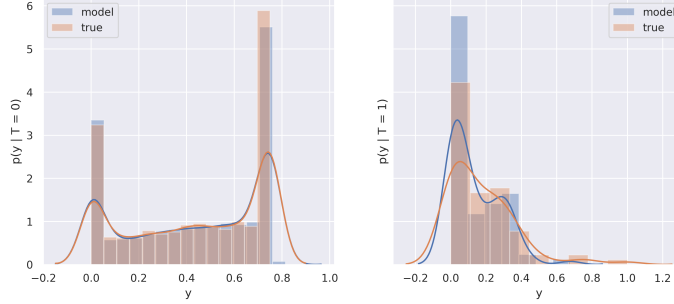
## 579 **Appendix B Visual Comparisons of Generated Distribution vs. Real** 580 **Distributions**

581 In this appendix, we provide the visualizations of the sigmoidal flows and the baseline linear generative  
582 models for each dataset. Each figure takes up a single page and corresponds to a single dataset. For  
583 each figure, the first half of the plots are for the sigmoidal flow and the second half of the plots are  
584 for the linear model. Because each figure takes up a whole page, the figures begin on the next page.

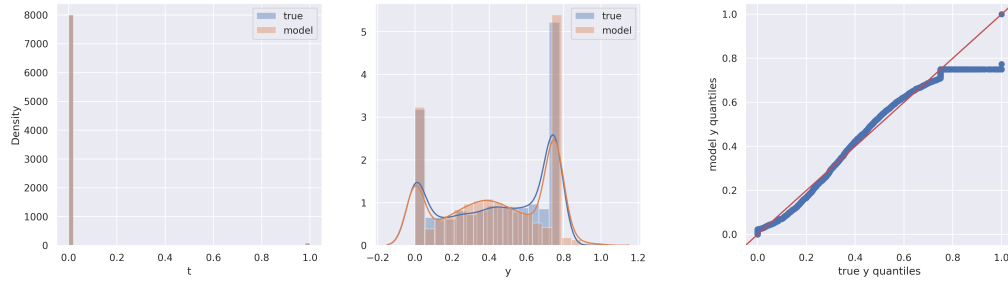


(a) Marginal distributions  $P(T)$  and  $P_{\text{model}}(T)$  on the left and marginal distributions  $P(Y)$  and  $P_{\text{model}}(Y)$  on the right.

(b) Q-Q plot of  $P_{\text{model}}(Y)$  and  $P(Y)$ .

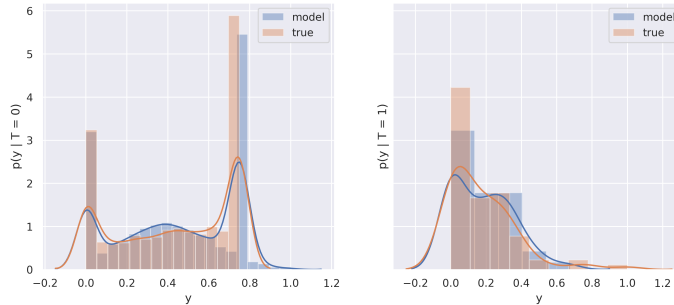


(c) Histogram and kernel density estimate visualization of  $P(Y | T)$  and  $P_{\text{model}}(Y | T)$ . Both graphs share the same y-axis.



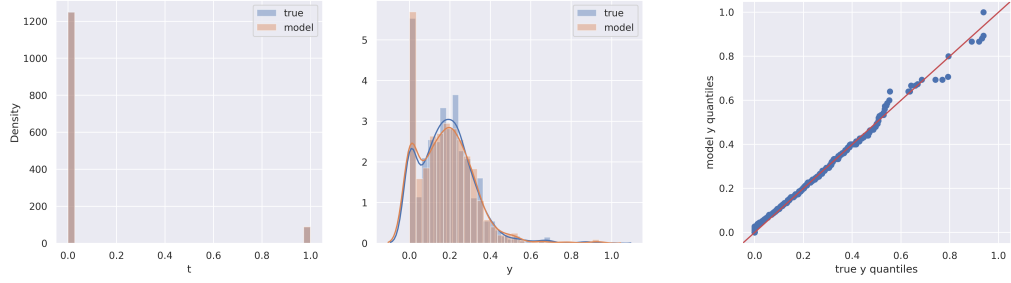
(d) Marginal distributions  $P(T)$  and  $P_{\text{model}}(T)$  on the left and marginal distributions  $P(Y)$  and  $P_{\text{model}}(Y)$  on the right.

(e) Q-Q plot of  $P_{\text{model}}(Y)$  and  $P(Y)$ .



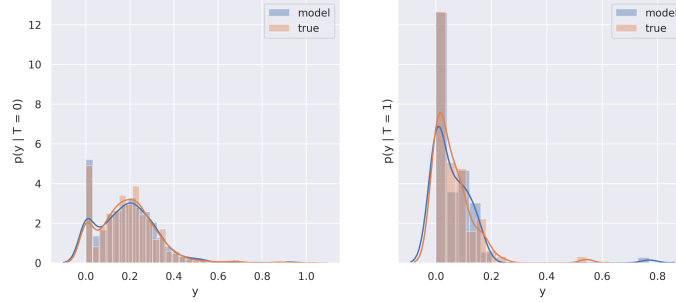
(f) Histogram and kernel density estimate visualization of  $P(Y | T)$  and  $P_{\text{model}}(Y | T)$ . Both graphs share the same y-axis.

Figure 3: LaLonde CPS – Visualizations of how well the generative model models the real data. Figures (a) - (c) are visualizations of the sigmoidal flow model. Figures (d) - (f) are visualizations of the baseline linear Gaussian model.

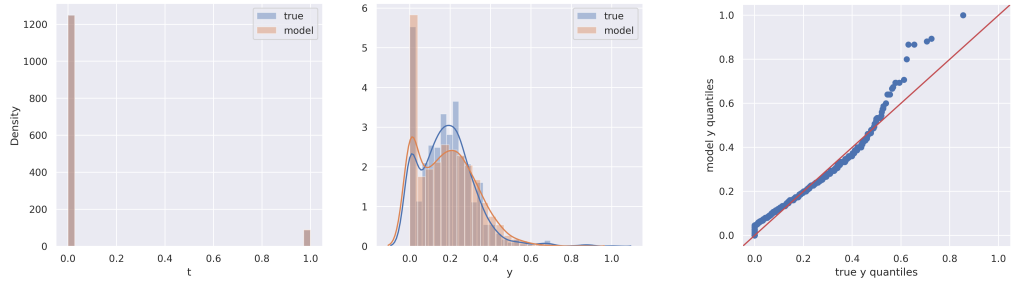


(a) Marginal distributions  $P(T)$  and  $P_{\text{model}}(T)$  on the left and marginal distributions  $P(Y)$  and  $P_{\text{model}}(Y)$  on the right.

(b) Q-Q plot of  $P_{\text{model}}(Y)$  and  $P(Y)$ .

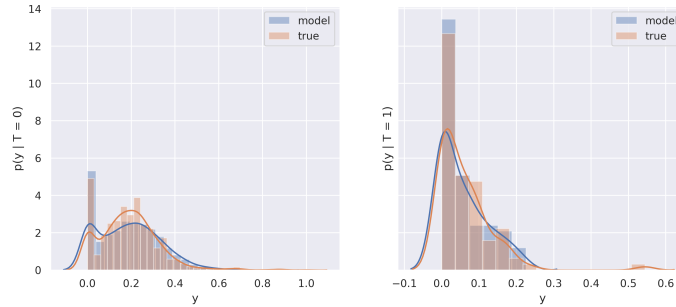


(c) Histogram and kernel density estimate visualization of  $P(Y | T)$  and  $P_{\text{model}}(Y | T)$ . Both graphs share the same y-axis.



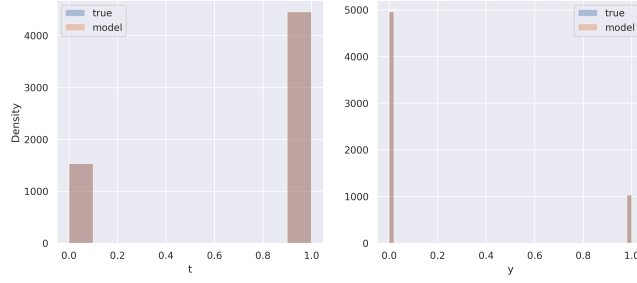
(d) Marginal distributions  $P(T)$  and  $P_{\text{model}}(T)$  on the left and marginal distributions  $P(Y)$  and  $P_{\text{model}}(Y)$  on the right.

(e) Q-Q plot of  $P_{\text{model}}(Y)$  and  $P(Y)$ .

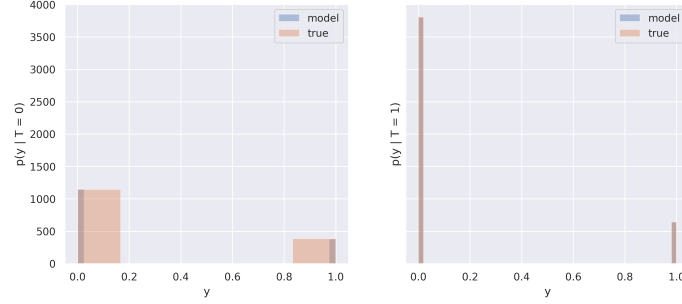


(f) Histogram and kernel density estimate visualization of  $P(Y | T)$  and  $P_{\text{model}}(Y | T)$ . Both graphs share the same y-axis.

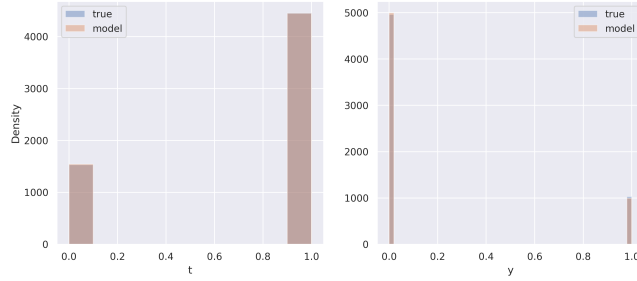
Figure 4: LaLonde PSID – Visualizations of how well the generative model models the real data. Figures (a) - (c) are visualizations of the sigmoidal flow model. Figures (d) - (f) are visualizations of the baseline linear Gaussian model.



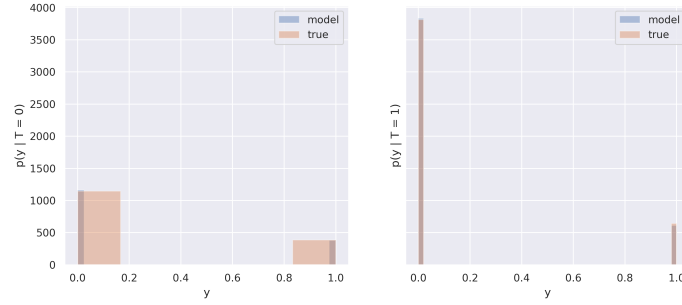
(a) Marginal distributions  $P(T)$  and  $P_{\text{model}}(T)$  on the left and marginal distributions  $P(Y)$  and  $P_{\text{model}}(Y)$  on the right.



(b) Histogram and kernel density estimate visualization of  $P(Y | T)$  and  $P_{\text{model}}(Y | T)$ . Both graphs share the same y-axis.



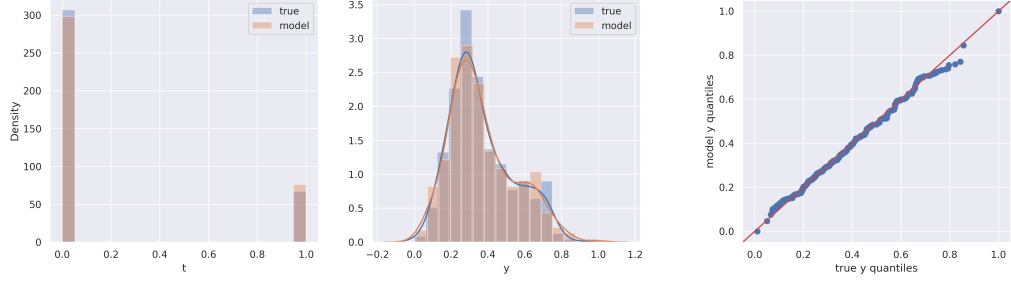
(c) Marginal distributions  $P(T)$  and  $P_{\text{model}}(T)$  on the left and marginal distributions  $P(Y)$  and  $P_{\text{model}}(Y)$  on the right.



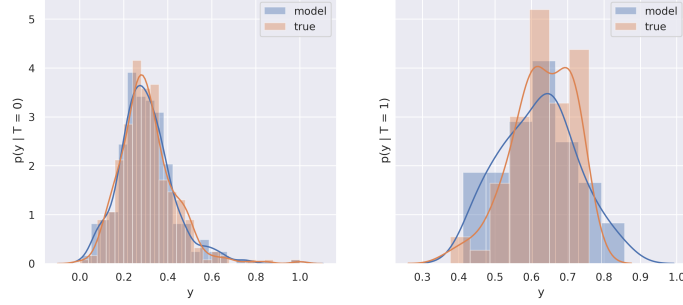
(d) Histogram and kernel density estimate visualization of  $P(Y | T)$  and  $P_{\text{model}}(Y | T)$ . Both graphs share the same y-axis.

Figure 5: Twins – Visualizations of how well the generative model models the dataset. Figures (a) - (c) are visualizations of the sigmoidal flow model. Figures (d) - (f) are visualizations of the baseline linear Gaussian model.

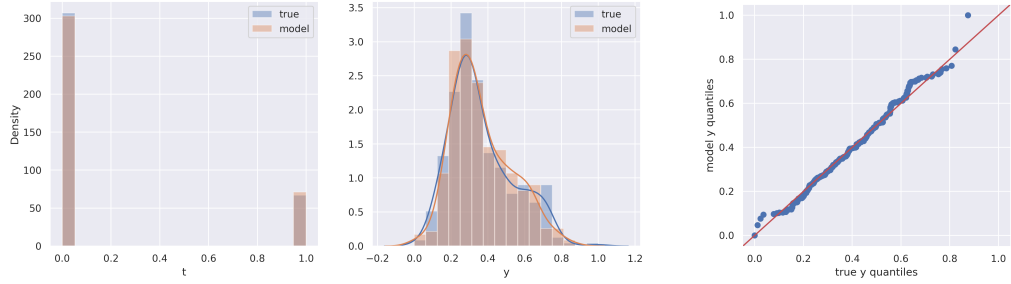




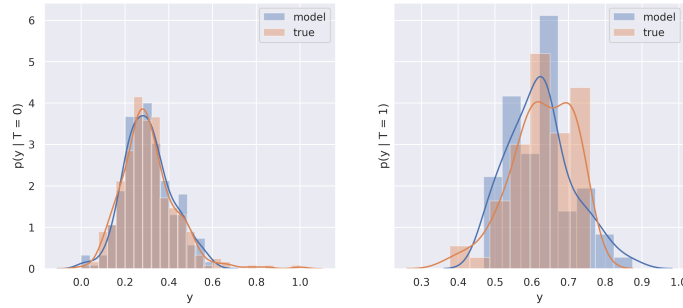
(a) Marginal distributions  $P(T)$  and  $P_{\text{model}}(T)$  on the left and marginal distributions  $P(Y)$  and  $P_{\text{model}}(Y)$  on the right. (b) Q-Q plot of  $P_{\text{model}}(Y)$  and  $P(Y)$ .



(c) Histogram and kernel density estimate visualization of  $P(Y | T)$  and  $P_{\text{model}}(Y | T)$ . Both graphs share the same y-axis.

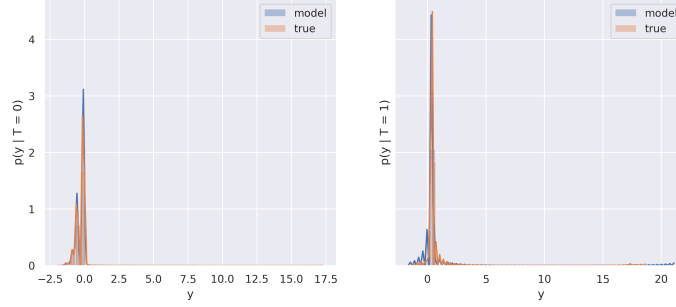
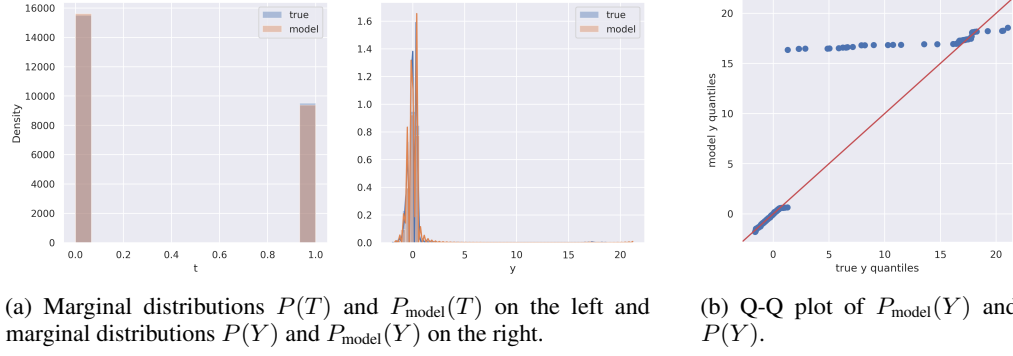


(d) Marginal distributions  $P(T)$  and  $P_{\text{model}}(T)$  on the left and marginal distributions  $P(Y)$  and  $P_{\text{model}}(Y)$  on the right. (e) Q-Q plot of  $P_{\text{model}}(Y)$  and  $P(Y)$ .

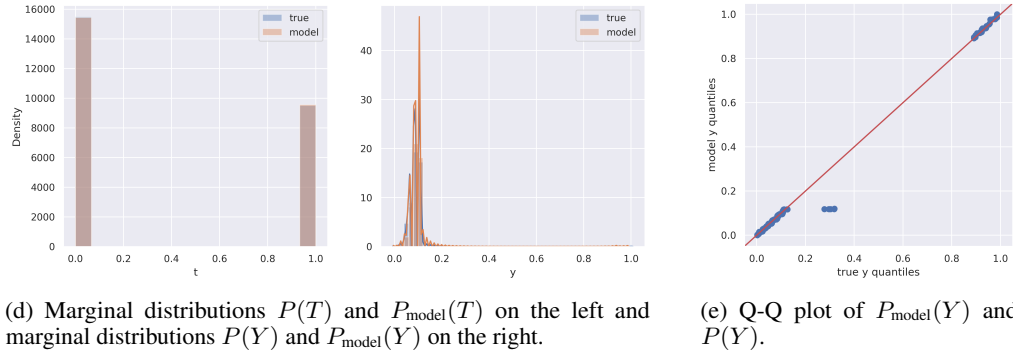


(f) Histogram and kernel density estimate visualization of  $P(Y | T)$  and  $P_{\text{model}}(Y | T)$ . Both graphs share the same y-axis.

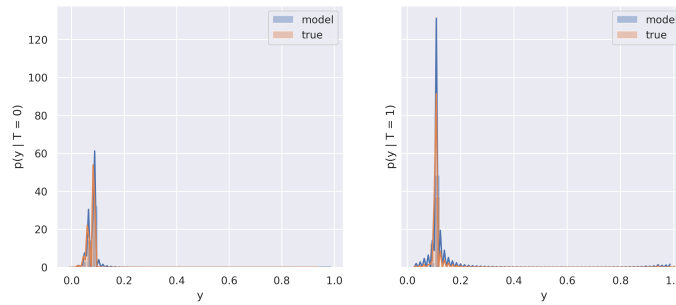
Figure 6: IHDP – Visualizations of how well the generative model models the dataset. Figures (a) - (c) are visualizations of the sigmoidal flow model. Figures (d) - (f) are visualizations of the baseline linear Gaussian model.



(c) Histogram and kernel density estimate visualization of  $P(Y | T)$  and  $P_{\text{model}}(Y | T)$ . Both graphs share the same y-axis.

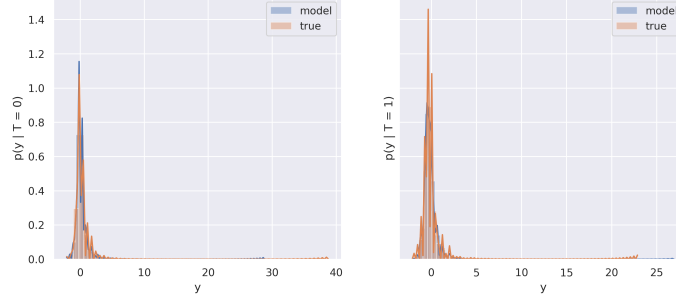
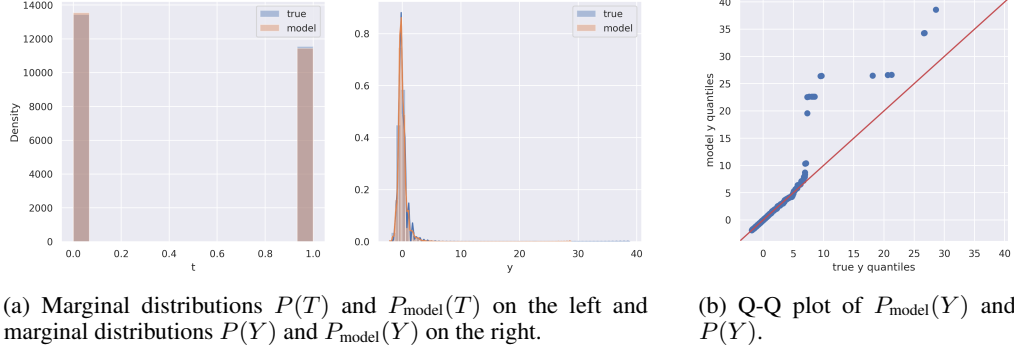


(d) Marginal distributions  $P(T)$  and  $P_{\text{model}}(T)$  on the left and marginal distributions  $P(Y)$  and  $P_{\text{model}}(Y)$  on the right.

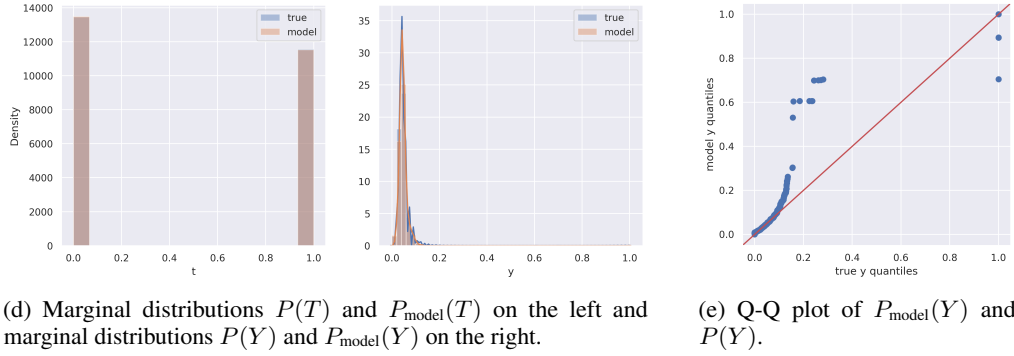


(f) Histogram and kernel density estimate visualization of  $P(Y | T)$  and  $P_{\text{model}}(Y | T)$ . Both graphs share the same y-axis.

Figure 7: LBIDD-Quadratic – Visualizations of how well the generative model models the dataset. Figures (a) - (c) are visualizations of the sigmoidal flow model. Figures (d) - (f) are visualizations of the baseline linear Gaussian model.



(c) Histogram and kernel density estimate visualization of  $P(Y | T)$  and  $P_{\text{model}}(Y | T)$ . Both graphs share the same y-axis.

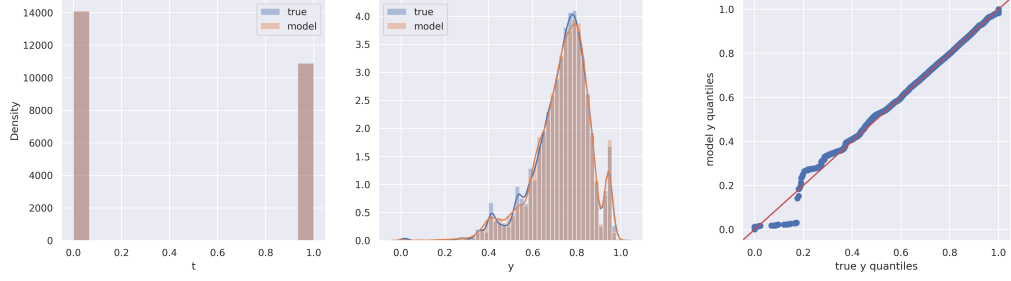


(d) Marginal distributions  $P(T)$  and  $P_{\text{model}}(T)$  on the left and marginal distributions  $P(Y)$  and  $P_{\text{model}}(Y)$  on the right.

(e) Q-Q plot of  $P_{\text{model}}(Y)$  and  $P(Y)$ .

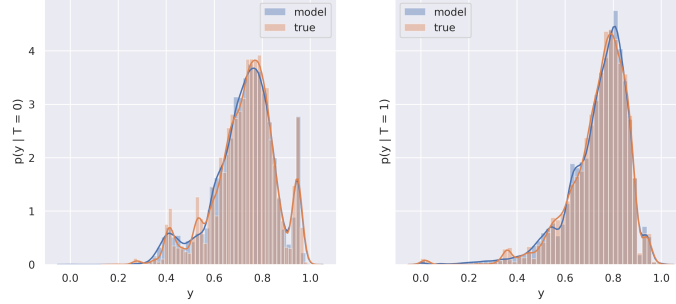
(f) Histogram and kernel density estimate visualization of  $P(Y | T)$  and  $P_{\text{model}}(Y | T)$ . Both graphs share the same y-axis.

Figure 8: LBIDD-Exponential – Visualizations of how well the generative model models the dataset. Figures (a) - (c) are visualizations of the sigmoidal flow model. Figures (d) - (f) are visualizations of the baseline linear Gaussian model.

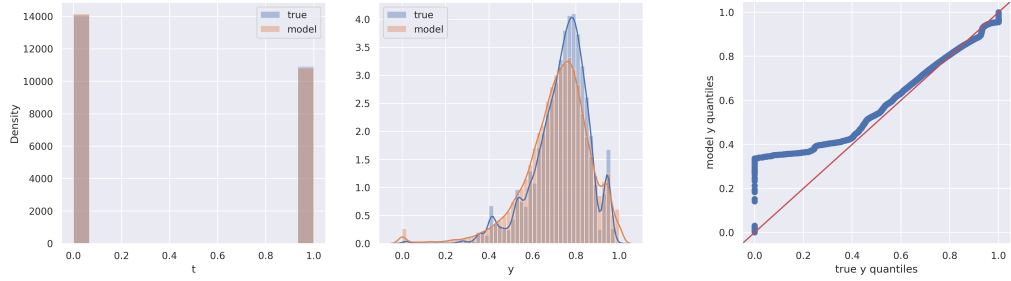


(a) Marginal distributions  $P(T)$  and  $P_{\text{model}}(T)$  on the left and marginal distributions  $P(Y)$  and  $P_{\text{model}}(Y)$  on the right.

(b) Q-Q plot of  $P_{\text{model}}(Y)$  and  $P(Y)$ .

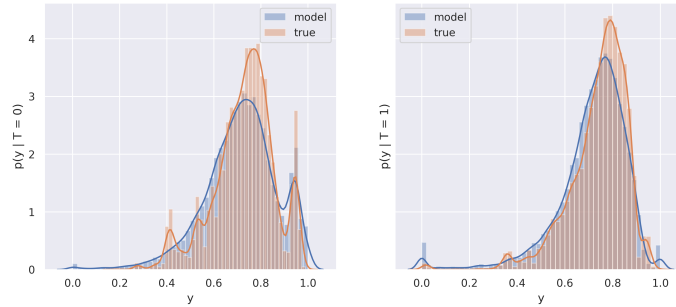


(c) Histogram and kernel density estimate visualization of  $P(Y | T)$  and  $P_{\text{model}}(Y | T)$ . Both graphs share the same y-axis.



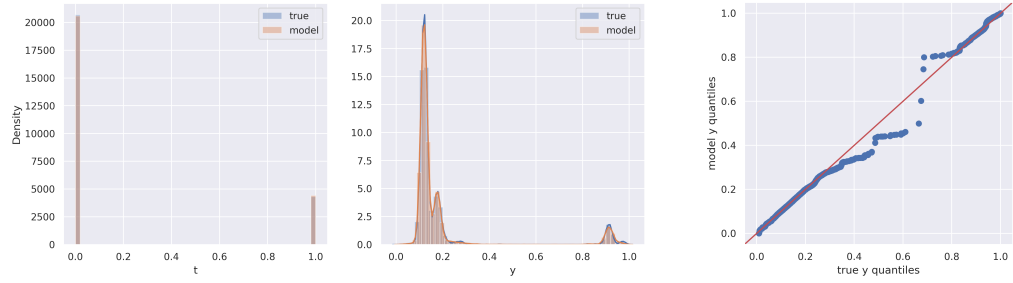
(d) Marginal distributions  $P(T)$  and  $P_{\text{model}}(T)$  on the left and marginal distributions  $P(Y)$  and  $P_{\text{model}}(Y)$  on the right.

(e) Q-Q plot of  $P_{\text{model}}(Y)$  and  $P(Y)$ .



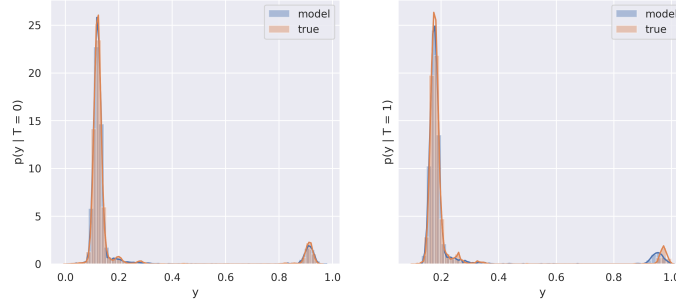
(f) Histogram and kernel density estimate visualization of  $P(Y | T)$  and  $P_{\text{model}}(Y | T)$ . Both graphs share the same y-axis.

Figure 9: LBIDD-Log – Visualizations of how well the generative model models the dataset. Figures (a) - (c) are visualizations of the sigmoidal flow model. Figures (d) - (f) are visualizations of the baseline linear Gaussian model.

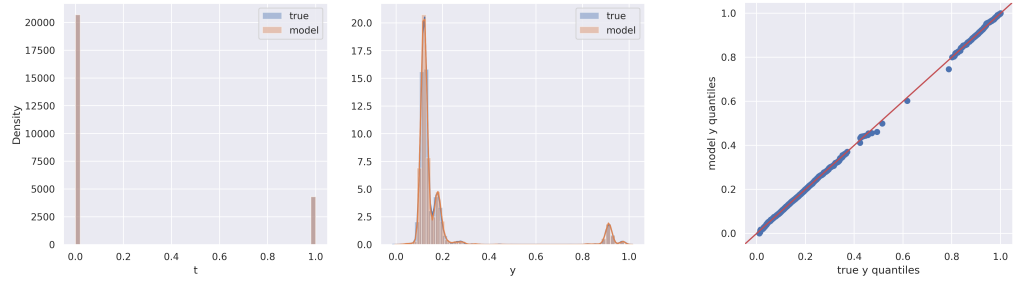


(a) Marginal distributions  $P(T)$  and  $P_{\text{model}}(T)$  on the left and marginal distributions  $P(Y)$  and  $P_{\text{model}}(Y)$  on the right.

(b) Q-Q plot of  $P_{\text{model}}(Y)$  and  $P(Y)$ .

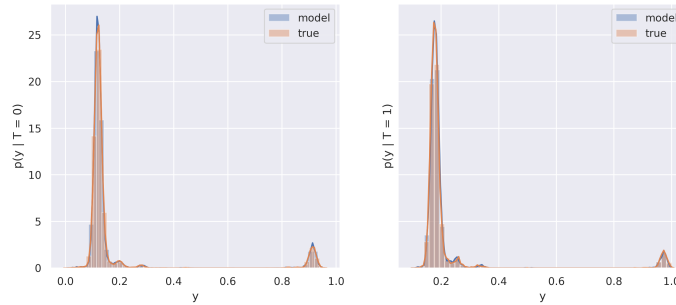


(c) Histogram and kernel density estimate visualization of  $P(Y | T)$  and  $P_{\text{model}}(Y | T)$ . Both graphs share the same y-axis.



(d) Marginal distributions  $P(T)$  and  $P_{\text{model}}(T)$  on the left and marginal distributions  $P(Y)$  and  $P_{\text{model}}(Y)$  on the right.

(e) Q-Q plot of  $P_{\text{model}}(Y)$  and  $P(Y)$ .



(f) Histogram and kernel density estimate visualization of  $P(Y | T)$  and  $P_{\text{model}}(Y | T)$ . Both graphs share the same y-axis.

Figure 10: LBIDD-Linear – Visualizations of how well the generative model models the dataset. Figures (a) - (c) are visualizations of the sigmoidal flow model. Figures (d) - (f) are visualizations of the baseline linear Gaussian model.



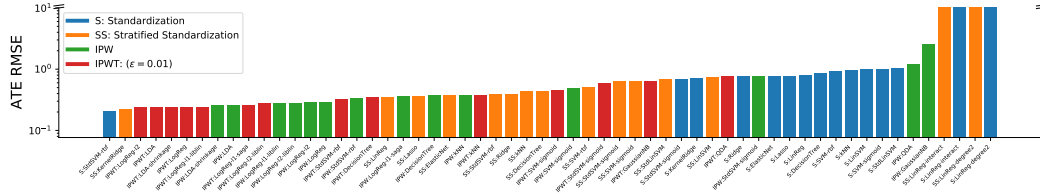


Figure 11: ATE RMSE of the different estimators, weighted averaged (by their inverse ATEs) over three datasets and color-coded by meta-estimator. This plot is in the main paper, but we include it here for comparisons next to the below plots.

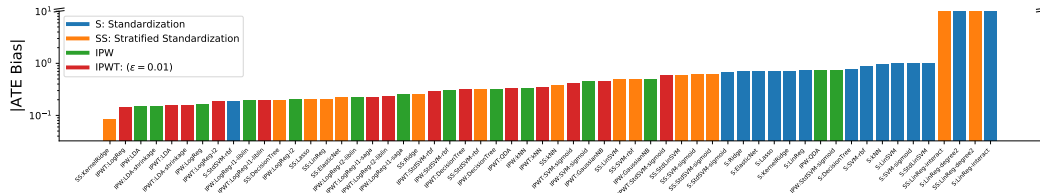


Figure 12: ATE absolute bias of the different estimators, weighted averaged (by their inverse ATEs) over three datasets and color-coded by meta-estimator.

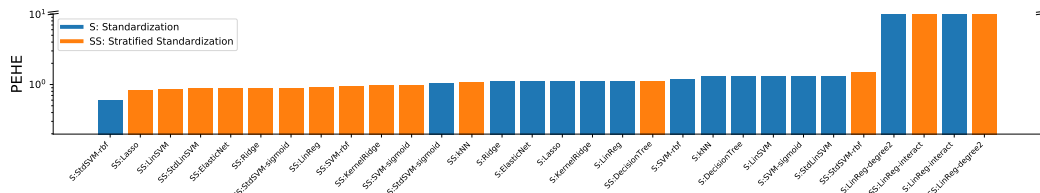


Figure 13: Average PEHE of the different standardization estimators, weighted averaged (by their inverse ATEs) over three datasets and color-coded by meta-estimator.

Table 5: Lalonde CPS Estimators sorted by ATE RMSE

meta-estimator	outcome_model	prop_score_model	ate_rmse	mean_pehe	ate_abs_bias	ate_std_error
ipw_trimeps.01		DecisionTree	1600.61		1292.06	944.74
ipw_trimeps.01		LogisticRegression_l2	1603.91		950.68	1291.79
ipw_trimeps.01		LogisticRegression	1612.56		973.41	1285.62
ipw_trimeps.01		LogisticRegression_l1_liblinear	1674.75		1089.23	1272.16
standardization	Standardized_SVM_sigmoid		1678.65	5633.50	1478.32	795.26
ipw		DecisionTree	1691.40		1063.80	1314.97
ipw_trimeps.01		Standardized_SVM_rbf	1826.27		1299.72	1282.97
ipw		Standardized_SVM_rbf	1875.67		1474.07	1159.84
standardization	Standardized_SVM_rbf		2052.88	3058.74	1882.01	819.98
ipw_trimeps.01		LogisticRegression_l1_saga	2214.21		1881.65	1167.10
ipw_trimeps.01		LogisticRegression_l2_liblinear	2420.16		1909.48	1486.96
ipw		kNN	2461.28		2320.38	820.81
ipw_trimeps.01		kNN	2472.01		2331.80	820.70
ipw_trimeps.01		GaussianNB	2487.79		2021.32	1450.29
ipw		LogisticRegression_l1_liblinear	2509.14		982.32	2308.86
ipw		LogisticRegression_l2_liblinear	2607.33		1579.18	2074.70
ipw		LogisticRegression	2607.86		1308.94	2255.57
ipw		LogisticRegression_l2	2625.46		1166.96	2351.86
ipw_trimeps.01		LDA	2627.57		1684.08	2016.93
ipw_trimeps.01		LDA_shrinkage	2643.51		1701.77	2022.90
ipw		Standardized_SVM_sigmoid	2664.59		2486.91	956.71
ipw_trimeps.01		Standardized_SVM_sigmoid	2806.76		2726.40	666.81
ipw_trimeps.01		SVM_sigmoid	2853.99		2757.42	736.14
ipw		LDA_shrinkage	3073.85		1503.30	2681.16
ipw		LDA	3076.71		1492.77	2690.31
stratified_standardization	Lasso		3360.80	5064.64	1826.03	2821.45
ipw		SVM_sigmoid	3362.98		3194.03	1052.52
stratified_standardization	LinearRegression		3383.14	5124.72	1813.37	2856.10
stratified_standardization	ElasticNet		3489.11	4889.81	2159.56	2740.47
standardization	LinearRegression_degree2		3593.31	5488.03	1201.96	3386.32
standardization	LinearRegression_interact		3789.72	5637.26	1560.70	3453.43
stratified_standardization	Ridge		3936.74	4927.52	2904.55	2657.35
stratified_standardization	Standardized_LinearSVM		3980.05	5995.41	2405.59	3170.80
ipw		LogisticRegression_l1_saga	4287.54		2293.69	3622.43
ipw_trimeps.01		QDA	4554.85		4173.15	1825.23
stratified_standardization	Standardized_SVM_rbf		4639.97	5326.08	4406.92	1452.04
stratified_standardization	Standardized_SVM_sigmoid		4708.39	7895.10	4613.74	939.35
stratified_standardization	DecisionTree		4726.51	6953.37	372.14	4711.83
stratified_standardization	LinearSVM		4902.21	5209.82	2173.07	4394.25
stratified_standardization	SVM_rbf		5281.01	6240.25	5212.29	849.16
stratified_standardization	SVM_sigmoid		5510.08	6883.21	5451.39	802.03
ipw		GaussianNB	5831.03		4051.81	4193.29
standardization	Standardized_LinearSVM		6530.70	9157.57	6488.87	737.98
standardization	LinearRegression		6971.02	9467.81	6910.01	920.23
standardization	Ridge		6998.69	9508.04	6987.78	390.59
standardization	ElasticNet		7050.64	9550.25	7049.30	137.34
standardization	Lasso		7053.98	9552.70	7052.56	141.86
standardization	KernelRidge		7092.93	9577.83	7082.44	385.52
standardization	DecisionTree		7135.67	9613.87	7135.67	0.00
standardization	LinearSVM		7135.67	9613.87	7135.67	0.00
standardization	SVM_rbf		7135.67	9613.87	7135.67	0.00
standardization	SVM_sigmoid		7135.67	9613.87	7135.67	0.00
standardization	kNN		7135.82	9613.85	7135.82	0.41
ipw		QDA	8259.79		6544.48	5039.24
stratified_standardization	LinearRegression_interact		44091.19	105675.94	11020.68	42691.66
stratified_standardization	LinearRegression_degree2		57156.93	100705.35	15734.67	54948.48
stratified_standardization	kNN					

Table 6: Lalonde PSID Estimators sorted by ATE RMSE

meta-estimator	outcome_model	prop_score_model	ate_rmse	mean_pehe	ate_abs_bias	ate_std_error
stratified_standardization	SVM_sigmoid		1374.50	10744.12	435.80	1303.58
stratified_standardization	kNN		1683.85	3065.36	406.78	1633.97
ipw		Standardized_SVM_rbf	1901.10		1327.76	1360.59
ipw_trimeps.01		Standardized_SVM_rbf	1905.76		1332.77	1362.23
standardization	Standardized_SVM_rbf		2298.04	6375.94	2098.06	937.60
ipw_trimeps.01		LogisticRegression_l1_saga	2348.67		1520.06	1790.45
ipw		kNN	2388.52		1351.69	1969.26
ipw		LogisticRegression_l1_saga	2403.84		1850.84	1533.89
ipw_trimeps.01		kNN	2412.94		1394.47	1969.19
ipw		LogisticRegression_l2_liblinear	2475.93		1933.00	1547.18
ipw_trimeps.01		LogisticRegression_l2_liblinear	2481.44		1785.19	1723.55
ipw		LogisticRegression_l2	2483.25		1952.25	1534.68
ipw		LogisticRegression_l1_liblinear	2484.97		1950.50	1539.68
ipw_trimeps.01		LogisticRegression_l2	2488.65		1803.00	1715.39
ipw_trimeps.01		LogisticRegression_l1_liblinear	2494.65		1798.03	1729.27
ipw		LDA_shrinkage	2548.83		1101.51	2298.53
ipw		LDA	2549.87		1097.66	2301.52
ipw_trimeps.01		LDA_shrinkage	2560.36		1127.30	2298.83
ipw_trimeps.01		LDA	2561.22		1123.23	2301.78
ipw_trimeps.01		LogisticRegression	2582.76		1969.73	1670.57
ipw		LogisticRegression	2590.52		2127.49	1478.04
stratified_standardization	SVM_rbf		2674.25	4416.33	2320.72	1328.86
stratified_standardization	Standardized_SVM_rbf		2703.05	4312.85	2438.89	1165.47
stratified_standardization	Standardized_SVM_sigmoid		2787.82	5333.67	2398.72	1420.60
ipw_trimeps.01		DecisionTree	2851.76		2257.44	1742.56
stratified_standardization	DecisionTree		2894.01	8570.48	1429.13	2516.52
ipw_trimeps.01		SVM_sigmoid	2905.19		1892.72	2204.02
ipw		GaussianNB	3049.13		2636.48	1531.73
ipw		SVM_sigmoid	3423.57		2462.19	2378.76
stratified_standardization	LinearRegression		3694.49	5379.93	2855.89	2343.74
ipw		DecisionTree	3731.41		3152.58	1996.15
standardization	LinearRegression_degree2		3814.48	5961.98	3122.71	2190.65
standardization	LinearRegression_interact		3877.79	6094.42	3157.29	2251.40
stratified_standardization	KernelRidge		3965.28	4584.82	258.97	3956.81
stratified_standardization	Lasso		4123.09	4614.18	2797.95	3028.43
ipw_trimeps.01		QDA	4128.07		3521.42	2154.20
stratified_standardization	ElasticNet		4240.31	4767.89	2716.97	3255.51
stratified_standardization	Ridge		4287.98	4825.75	2706.11	3326.22
ipw_trimeps.01		GaussianNB	5184.90		4643.51	2306.73
stratified_standardization	Standardized_LinearSVM		6060.79	8109.11	5676.84	2122.89
stratified_standardization	LinearSVM		6703.79	8962.51	2077.53	6373.75
ipw		QDA	9170.70		7210.86	5666.14
ipw_trimeps.01		Standardized_SVM_sigmoid	10482.60		10209.46	2377.36
standardization	Standardized_SVM_sigmoid		10574.70	15968.21	10567.02	402.89
standardization	DecisionTree		13124.53	18066.16	13124.53	0.00
standardization	ElasticNet		13124.53	18066.16	13124.53	0.00
standardization	Lasso		13124.53	18066.16	13124.53	0.00
standardization	LinearSVM		13124.53	18066.16	13124.53	0.00
standardization	SVM_rbf		13124.53	18066.16	13124.53	0.00
standardization	SVM_sigmoid		13124.53	18066.16	13124.53	0.00
standardization	Ridge		13125.30	18066.73	13125.30	1.18
standardization	kNN		13125.92	18065.93	13125.92	1.29
standardization	KernelRidge		13198.56	18120.34	13198.06	113.96
standardization	LinearRegression		14189.19	18818.59	14097.20	1613.08
standardization	Standardized_LinearSVM		14556.98	19110.22	14486.66	1429.18
ipw		Standardized_SVM_sigmoid	17323.20		16823.85	4129.33
stratified_standardization	LinearRegression_interact		53972.80	105866.18	37749.80	38574.81
stratified_standardization	LinearRegression_degree2		103618.34	201053.48	93842.31	43936.11

Table 7: Twins Estimators sorted by ATE RMSE

meta-estimator	outcome_model	prop_score_model	ate_rmse	mean_pehe	ate_abs_bias	ate_std_error
standardization	KernelRidge		0.01	0.04	0.01	0.01
standardization	Standardized_SVM_rbf		0.01	0.06	0.01	0.01
stratified_standardization	KernelRidge		0.01	0.11	0.01	0.01
ipw		LDA	0.01		0.01	0.01
ipw		LDA_shrinkage	0.01		0.01	0.01
ipw_trimeps.01		LDA	0.01		0.01	0.01
ipw_trimeps.01		LDA_shrinkage	0.01		0.01	0.01
standardization	ElasticNet		0.02	0.04	0.01	0.01
standardization	Lasso		0.02	0.04	0.01	0.01
standardization	LinearRegression		0.02	0.04	0.01	0.01
standardization	Ridge		0.02	0.04	0.01	0.01
stratified_standardization	Lasso		0.02	0.10	0.01	0.01
stratified_standardization	ElasticNet		0.02	0.11	0.01	0.01
stratified_standardization	LinearRegression		0.02	0.11	0.01	0.01
stratified_standardization	Ridge		0.02	0.11	0.01	0.01
stratified_standardization	Standardized_SVM_rbf		0.02	0.23	0.01	0.01
ipw		LogisticRegression	0.02		0.01	0.01
ipw_trimeps.01		LogisticRegression	0.02		0.01	0.01
ipw		LogisticRegression_l1_liblinear	0.02		0.02	0.01
ipw		LogisticRegression_l1_saga	0.02		0.02	0.01
ipw		LogisticRegression_l2	0.02		0.02	0.01
ipw		LogisticRegression_l2_liblinear	0.02		0.02	0.01
ipw_trimeps.01		LogisticRegression_l1_liblinear	0.02		0.02	0.01
ipw_trimeps.01		LogisticRegression_l1_saga	0.02		0.02	0.01
ipw_trimeps.01		LogisticRegression_l2	0.02		0.02	0.01
ipw_trimeps.01		LogisticRegression_l2_liblinear	0.02		0.02	0.01
stratified_standardization	DecisionTree		0.03	0.12	0.03	0.01
standardization	DecisionTree		0.04	0.08	0.02	0.04
stratified_standardization	SVM_rbf		0.04	0.11	0.04	0.01
ipw		DecisionTree	0.04		0.04	0.01
ipw		Standardized_SVM_rbf	0.04		0.04	0.02
ipw		Standardized_SVM_sigmoid	0.04		0.04	0.01
ipw		kNN	0.04		0.04	0.02
ipw_trimeps.01		DecisionTree	0.04		0.04	0.01
ipw_trimeps.01		Standardized_SVM_rbf	0.04		0.04	0.02
ipw_trimeps.01		Standardized_SVM_sigmoid	0.04		0.04	0.01
ipw_trimeps.01		kNN	0.04		0.04	0.02
standardization	SVM_rbf		0.05	0.06	0.04	0.01
stratified_standardization	kNN		0.05	0.13	0.05	0.02
ipw		SVM_sigmoid	0.05		0.05	0.01
ipw_trimeps.01		SVM_sigmoid	0.05		0.05	0.01
standardization	kNN		0.06	0.08	0.06	0.00
standardization	LinearSVM		0.07	0.08	0.07	0.00
standardization	SVM_sigmoid		0.07	0.08	0.07	0.00
standardization	Standardized_LinearSVM		0.07	0.08	0.07	0.00
standardization	Standardized_SVM_sigmoid		0.07	0.08	0.07	0.00
stratified_standardization	LinearSVM		0.07	0.08	0.07	0.00
stratified_standardization	SVM_sigmoid		0.07	0.08	0.07	0.00
stratified_standardization	Standardized_LinearSVM		0.07	0.08	0.07	0.00
stratified_standardization	Standardized_SVM_sigmoid		0.07	0.08	0.07	0.00
ipw_trimeps.01		GaussianNB	0.08		0.05	0.06
ipw_trimeps.01		QDA	0.09		0.01	0.09
ipw		QDA	0.12		0.05	0.10
ipw		GaussianNB	0.46		0.05	0.46
stratified_standardization	LinearRegression_interact		3006941.41	152626640.11	646047.20	2936719.20
standardization	LinearRegression_interact		3341204.71	168756300.74	1965661.75	2701818.43
stratified_standardization	LinearRegression_degree2		3928023.92	403823525.32	1505613.68	3628015.90
standardization	LinearRegression_degree2		10639314.01	122823195.44	1478780.90	10536043.36

Table 8: Lalonde CPS Estimators sorted by Absolute Bias

meta-estimator	outcome_model	prop_score_model	ate_abs_bias	ate_rmse	mean_pehe	ate_std_error
stratified_standardization	DecisionTree		372.14	4726.51	6953.37	4711.83
ipw_trimeps.01		LogisticRegression_l2	950.68	1603.91		1291.79
ipw_trimeps.01		LogisticRegression	973.41	1612.56		1285.62
ipw		LogisticRegression_l1_liblinear	982.32	2509.14		2308.86
ipw		DecisionTree	1063.80	1691.40		1314.97
ipw_trimeps.01	LinearRegression_degree2	LogisticRegression_l1_liblinear	1089.23	1674.75		1272.16
ipw		LogisticRegression_l2	1166.96	2625.46		2351.86
standardization			1201.96	3593.31	5488.03	3386.32
ipw_trimeps.01		DecisionTree	1292.06	1600.61		944.74
ipw_trimeps.01		Standardized_SVM_rbf	1299.72	1826.27		1282.97
ipw	Standardized_SVM_sigmoid	LogisticRegression	1308.94	2607.86		2255.57
ipw		Standardized_SVM_rbf	1474.07	1875.67		1159.84
standardization			1478.32	1678.65	5633.50	795.26
ipw		LDA	1492.77	3076.71		2690.31
ipw		LDA_shrinkage	1503.30	3073.85		2681.16
standardization	LinearRegression_interact		1560.70	3789.72	5637.26	3453.43
ipw		LogisticRegression_l2_liblinear	1579.18	2607.33		2740.70
ipw_trimeps.01		LDA	1684.08	2627.57		2016.93
ipw_trimeps.01		LDA_shrinkage	1701.77	2643.51		2022.90
stratified_standardization			1813.37	3383.14	5124.72	2856.10
stratified_standardization	Lasso		1826.03	3360.80	5064.64	2821.45
ipw_trimeps.01		LogisticRegression_l1_saga	1881.65	2214.21		1167.10
standardization	Standardized_SVM_rbf		1882.01	2052.88	3058.74	819.98
ipw_trimeps.01		LogisticRegression_l2_liblinear	1909.48	2420.16		1486.96
ipw_trimeps.01		GaussianNB	2021.32	2487.79		1450.29
stratified_standardization	ElasticNet		2159.56	3489.11	4889.81	2740.47
stratified_standardization	LinearSVM		2173.07	4902.21	5209.82	4394.25
ipw		LogisticRegression_l1_saga	2293.69	4287.54		3622.43
ipw	Standardized_LinearSVM	kNN	2320.38	2461.28		820.81
ipw_trimeps.01		kNN	2331.80	2472.01		820.70
stratified_standardization			2405.59	3980.05	5995.41	3170.80
ipw	Standardized_LinearSVM	Standardized_SVM_sigmoid	2486.91	2664.59		956.71
ipw_trimeps.01		Standardized_SVM_sigmoid	2726.40	2806.76		666.81
ipw_trimeps.01		SVM_sigmoid	2757.42	2853.99		736.14
stratified_standardization	Ridge		2904.55	3936.74	4927.52	2657.35
ipw		SVM_sigmoid	3194.03	3362.98		1052.52
ipw		GaussianNB	4051.81	5831.03		4193.29
ipw_trimeps.01	Standardized_LinearSVM	QDA	4173.15	4554.85		1825.23
stratified_standardization			4406.92	4639.97	5326.08	1452.04
stratified_standardization		Standardized_SVM_sigmoid	4613.74	4708.39	7895.10	939.35
stratified_standardization	SVM_rbf		5212.29	5281.01	6240.25	849.16
stratified_standardization		SVM_sigmoid	5451.39	5510.08	6883.21	802.03
standardization	Standardized_LinearSVM		6488.87	6530.70	9157.57	737.98
ipw		QDA	6544.48	8259.79		5039.24
standardization	LinearRegression		6910.01	6971.02	9467.81	920.23
standardization		Ridge	6987.78	6998.69	9508.04	390.59
standardization	ElasticNet		7049.30	7050.64	9550.25	137.34
standardization		Lasso	7052.56	7053.98	9552.70	141.86
standardization	KernelRidge		7082.44	7092.93	9577.83	385.52
standardization		DecisionTree	7135.67	7135.67	9613.87	0.00
standardization	LinearSVM		7135.67	7135.67	9613.87	0.00
standardization		SVM_rbf	7135.67	7135.67	9613.87	0.00
standardization	SVM_sigmoid		7135.67	7135.67	9613.87	0.00
standardization		kNN	7135.82	7135.82	9613.85	0.41
stratified_standardization	LinearRegression_interact		11020.68	44091.19	105675.94	42691.66
stratified_standardization		LinearRegression_degree2	15734.67	57156.93	100705.35	54948.48
stratified_standardization	kNN					
stratified_standardization						



Table 9: Lalonde PSID Estimators sorted by Absolute Bias

meta-estimator	outcome_model	prop_score_model	ate_abs_bias	ate_rmse	mean_pehe	ate_std_error
stratified_standardization	KernelRidge		258.97	3965.28	4584.82	3956.81
stratified_standardization	kNN		406.78	1683.85	3065.36	1633.97
stratified_standardization	SVM_sigmoid		435.80	1374.50	10744.12	1303.58
ipw		LDA	1097.66	2549.87		2301.52
ipw		LDA_shrinkage	1101.51	2548.83		2298.53
ipw_trimeps.01		LDA	1123.23	2561.22		2301.78
ipw_trimeps.01		LDA_shrinkage	1127.30	2560.36		2298.83
ipw		Standardized_SVM_rbf	1327.76	1901.10		1360.59
ipw_trimeps.01		Standardized_SVM_rbf	1332.77	1905.76		1362.23
ipw		kNN	1351.69	2388.52		1969.26
ipw_trimeps.01		kNN	1394.47	2412.94		1969.19
stratified_standardization	DecisionTree		1429.13	2894.01	8570.48	2516.52
ipw_trimeps.01		LogisticRegression_l1_saga	1520.06	2348.67		1790.45
ipw_trimeps.01		LogisticRegression_l2_liblinear	1785.19	2481.44		1723.55
ipw_trimeps.01		LogisticRegression_l1_liblinear	1798.03	2494.65		1729.27
ipw_trimeps.01		LogisticRegression_l2	1803.00	2488.65		1715.39
ipw		LogisticRegression_l1_saga	1850.84	2403.84		1533.89
ipw_trimeps.01		SVM_sigmoid	1892.72	2905.19		2204.02
ipw		LogisticRegression_l2_liblinear	1933.00	2475.93		1547.18
ipw		LogisticRegression_l1_liblinear	1950.50	2484.97		1539.68
ipw		LogisticRegression_l2	1952.25	2483.25		1534.68
ipw_trimeps.01		LogisticRegression	1969.73	2582.76		1670.57
stratified_standardization	LinearSVM		2077.53	6703.79	8962.51	6373.75
standardization	Standardized_SVM_rbf		2098.06	2298.04	6375.94	937.60
ipw		LogisticRegression	2127.49	2590.52		1478.04
ipw_trimeps.01		DecisionTree	2257.44	2851.76		1742.56
stratified_standardization	SVM_rbf		2320.72	2674.25	4416.33	1328.86
stratified_standardization	Standardized_SVM_sigmoid		2398.72	2787.82	5333.67	1420.60
stratified_standardization	Standardized_SVM_rbf		2438.89	2703.05	4312.85	1165.47
ipw		SVM_sigmoid	2462.19	3423.57		2378.76
ipw		GaussianNB	2636.48	3049.13		1531.73
stratified_standardization	Ridge		2706.11	4287.98	4825.75	3326.22
stratified_standardization	ElasticNet		2716.97	4240.31	4767.89	3255.51
stratified_standardization	Lasso		2797.95	4123.09	4614.18	3028.43
stratified_standardization	LinearRegression		2855.89	3694.49	5379.93	2343.74
standardization	LinearRegression_degree2		3122.71	3814.48	5961.98	2190.65
ipw		DecisionTree	3152.58	3731.41		1996.15
standardization	LinearRegression_interact		3157.29	3877.79	6094.42	2251.40
ipw_trimeps.01		QDA	3521.42	4128.07		2154.20
ipw_trimeps.01		GaussianNB	4643.51	5184.90		2306.73
stratified_standardization	Standardized_LinearSVM		5676.84	6060.79	8109.11	2122.89
ipw		QDA	7210.86	9170.70		5666.14
ipw_trimeps.01		Standardized_SVM_sigmoid	10209.46	10482.60		2377.36
standardization	Standardized_SVM_sigmoid		10567.02	10574.70	15968.21	402.89
standardization	DecisionTree		13124.53	13124.53	18066.16	0.00
standardization	ElasticNet		13124.53	13124.53	18066.16	0.00
standardization	Lasso		13124.53	13124.53	18066.16	0.00
standardization	LinearSVM		13124.53	13124.53	18066.16	0.00
standardization	SVM_rbf		13124.53	13124.53	18066.16	0.00
standardization	SVM_sigmoid		13124.53	13124.53	18066.16	0.00
standardization	Ridge		13125.30	13125.30	18066.73	1.18
standardization	kNN		13125.92	13125.92	18065.93	1.29
standardization	KernelRidge		13198.06	13198.56	18120.34	113.96
standardization	LinearRegression		14097.20	14189.19	18818.59	1613.08
standardization	Standardized_LinearSVM		14486.66	14556.98	19110.22	1429.18
ipw		Standardized_SVM_sigmoid	16823.85	17323.20		4129.33
stratified_standardization	LinearRegression_interact		37749.80	53972.80	105866.18	38574.81
stratified_standardization	LinearRegression_degree2		93842.31	103618.34	201053.48	43936.11

Table 10: Twins Estimators sorted by Absolute Bias

meta-estimator	outcome_model	prop_score_model	ate_abs_bias	ate_rmse	mean_pehe	ate_std_error
standardization	KernelRidge		0.01	0.01	0.04	0.01
standardization	Standardized_SVM_rbf		0.01	0.01	0.06	0.01
stratified_standardization	KernelRidge		0.01	0.01	0.11	0.01
ipw		LDA	0.01	0.01		0.01
ipw		LDA_shrinkage	0.01	0.01		0.01
ipw_trimeps.01		LDA	0.01	0.01		0.01
ipw_trimeps.01		LDA_shrinkage	0.01	0.01		0.01
standardization	ElasticNet		0.01	0.02	0.04	0.01
standardization	Lasso		0.01	0.02	0.04	0.01
standardization	LinearRegression		0.01	0.02	0.04	0.01
standardization	Ridge		0.01	0.02	0.04	0.01
stratified_standardization	Lasso		0.01	0.02	0.10	0.01
stratified_standardization	ElasticNet		0.01	0.02	0.11	0.01
stratified_standardization	LinearRegression		0.01	0.02	0.11	0.01
stratified_standardization	Ridge		0.01	0.02	0.11	0.01
stratified_standardization	Standardized_SVM_rbf		0.01	0.02	0.23	0.01
ipw		LogisticRegression	0.01	0.02		0.01
ipw_trimeps.01		LogisticRegression	0.01	0.02		0.01
ipw_trimeps.01		QDA	0.01	0.09		0.09
ipw		LogisticRegression_l1_liblinear	0.02	0.02		0.01
ipw		LogisticRegression_l1_saga	0.02	0.02		0.01
ipw		LogisticRegression_l2	0.02	0.02		0.01
ipw		LogisticRegression_l2_liblinear	0.02	0.02		0.01
ipw_trimeps.01		LogisticRegression_l1_liblinear	0.02	0.02		0.01
ipw_trimeps.01		LogisticRegression_l1_saga	0.02	0.02		0.01
ipw_trimeps.01		LogisticRegression_l2	0.02	0.02		0.01
ipw_trimeps.01		LogisticRegression_l2_liblinear	0.02	0.02		0.01
standardization	DecisionTree		0.02	0.04	0.08	0.04
stratified_standardization	DecisionTree		0.03	0.03	0.12	0.01
stratified_standardization	SVM_rbf		0.04	0.04	0.11	0.01
ipw		DecisionTree	0.04	0.04		0.01
ipw		Standardized_SVM_rbf	0.04	0.04		0.02
ipw		Standardized_SVM_sigmoid	0.04	0.04		0.01
ipw		kNN	0.04	0.04		0.02
ipw_trimeps.01		DecisionTree	0.04	0.04		0.01
ipw_trimeps.01		Standardized_SVM_rbf	0.04	0.04		0.02
ipw_trimeps.01		Standardized_SVM_sigmoid	0.04	0.04		0.01
ipw_trimeps.01		kNN	0.04	0.04		0.02
standardization	SVM_rbf		0.04	0.05	0.06	0.01
stratified_standardization	kNN		0.05	0.05	0.13	0.02
ipw		SVM_sigmoid	0.05	0.05		0.01
ipw_trimeps.01		SVM_sigmoid	0.05	0.05		0.01
ipw_trimeps.01		GaussianNB	0.05	0.08		0.06
ipw		QDA	0.05	0.12		0.10
ipw		GaussianNB	0.05	0.46		0.46
standardization	kNN		0.06	0.06	0.08	0.00
standardization	LinearSVM		0.07	0.07	0.08	0.00
standardization	SVM_sigmoid		0.07	0.07	0.08	0.00
standardization	Standardized_LinearSVM		0.07	0.07	0.08	0.00
standardization	Standardized_SVM_sigmoid		0.07	0.07	0.08	0.00
stratified_standardization	LinearSVM		0.07	0.07	0.08	0.00
stratified_standardization	SVM_sigmoid		0.07	0.07	0.08	0.00
stratified_standardization	Standardized_LinearSVM		0.07	0.07	0.08	0.00
stratified_standardization	Standardized_SVM_sigmoid		0.07	0.07	0.08	0.00
stratified_standardization	LinearRegression_interact		646047.20	3006941.41	152626640.11	2936719.20
standardization	LinearRegression_degree2		1478780.90	10639314.01	122823195.44	10536043.36
stratified_standardization	LinearRegression_degree2		1505613.68	3928023.92	403823525.32	3628015.90
standardization	LinearRegression_interact		1965661.75	3341204.71	168756300.74	2701818.43

Table 11: Lalonde CPS Estimators sorted by PEHE

meta-estimator	outcome_model	mean_pehe	ate_rmse	ate_abs_bias	ate_std_error
standardization	Standardized_SVM_rbf	3058.74	2052.88	1882.01	819.98
stratified_standardization	ElasticNet	4889.81	3489.11	2159.56	2740.47
stratified_standardization	Ridge	4927.52	3936.74	2904.55	2657.35
stratified_standardization	Lasso	5064.64	3360.80	1826.03	2821.45
stratified_standardization	LinearRegression	5124.72	3383.14	1813.37	2856.10
stratified_standardization	LinearSVM	5209.82	4902.21	2173.07	4394.25
stratified_standardization	Standardized_SVM_rbf	5326.08	4639.97	4406.92	1452.04
standardization	LinearRegression_degree2	5488.03	3593.31	1201.96	3386.32
standardization	Standardized_SVM_sigmoid	5633.50	1678.65	1478.32	795.26
standardization	LinearRegression_interact	5637.26	3789.72	1560.70	3453.43
stratified_standardization	Standardized_LinearSVM	5995.41	3980.05	2405.59	3170.80
stratified_standardization	SVM_rbf	6240.25	5281.01	5212.29	849.16
stratified_standardization	SVM_sigmoid	6883.21	5510.08	5451.39	802.03
stratified_standardization	DecisionTree	6953.37	4726.51	372.14	4711.83
stratified_standardization	Standardized_SVM_sigmoid	7895.10	4708.39	4613.74	939.35
standardization	Standardized_LinearSVM	9157.57	6530.70	6488.87	737.98
standardization	LinearRegression	9467.81	6971.02	6910.01	920.23
standardization	Ridge	9508.04	6998.69	6987.78	390.59
standardization	ElasticNet	9550.25	7050.64	7049.30	137.34
standardization	Lasso	9552.70	7053.98	7052.56	141.86
standardization	KernelRidge	9577.83	7092.93	7082.44	385.52
standardization	kNN	9613.85	7135.82	7135.82	0.41
standardization	DecisionTree	9613.87	7135.67	7135.67	0.00
standardization	LinearSVM	9613.87	7135.67	7135.67	0.00
standardization	SVM_rbf	9613.87	7135.67	7135.67	0.00
standardization	SVM_sigmoid	9613.87	7135.67	7135.67	0.00
stratified_standardization	LinearRegression_degree2	100705.35	57156.93	15734.67	54948.48
stratified_standardization	LinearRegression_interact	105675.94	44091.19	11020.68	42691.66
stratified_standardization	kNN				

Table 12: Lalonde PSID Estimators sorted by PEHE

meta-estimator	outcome_model	mean_pehe	ate_rmse	ate_abs_bias	ate_std_error
stratified_standardization	kNN	3065.36	1683.85	406.78	1633.97
stratified_standardization	Standardized_SVM_rbf	4312.85	2703.05	2438.89	1165.47
stratified_standardization	SVM_rbf	4416.33	2674.25	2320.72	1328.86
stratified_standardization	KernelRidge	4584.82	3965.28	258.97	3956.81
stratified_standardization	Lasso	4614.18	4123.09	2797.95	3028.43
stratified_standardization	ElasticNet	4767.89	4240.31	2716.97	3255.51
stratified_standardization	Ridge	4825.75	4287.98	2706.11	3326.22
stratified_standardization	Standardized_SVM_sigmoid	5333.67	2787.82	2398.72	1420.60
stratified_standardization	LinearRegression	5379.93	3694.49	2855.89	2343.74
standardization	LinearRegression_degree2	5961.98	3814.48	3122.71	2190.65
standardization	LinearRegression_interact	6094.42	3877.79	3157.29	2251.40
standardization	Standardized_SVM_rbf	6375.94	2298.04	2098.06	937.60
stratified_standardization	Standardized_LinearSVM	8109.11	6060.79	5676.84	2122.89
stratified_standardization	DecisionTree	8570.48	2894.01	1429.13	2516.52
stratified_standardization	LinearSVM	8962.51	6703.79	2077.53	6373.75
stratified_standardization	SVM_sigmoid	10744.12	1374.50	435.80	1303.58
standardization	Standardized_SVM_sigmoid	15968.21	10574.70	10567.02	402.89
standardization	kNN	18065.93	13125.92	13125.92	1.29
standardization	DecisionTree	18066.16	13124.53	13124.53	0.00
standardization	ElasticNet	18066.16	13124.53	13124.53	0.00
standardization	Lasso	18066.16	13124.53	13124.53	0.00
standardization	LinearSVM	18066.16	13124.53	13124.53	0.00
standardization	SVM_rbf	18066.16	13124.53	13124.53	0.00
standardization	SVM_sigmoid	18066.16	13124.53	13124.53	0.00
standardization	Ridge	18066.73	13125.30	13125.30	1.18
standardization	KernelRidge	18120.34	13198.56	13198.06	113.96
standardization	LinearRegression	18818.59	14189.19	14097.20	1613.08
standardization	Standardized_LinearSVM	19110.22	14556.98	14486.66	1429.18
stratified_standardization	LinearRegression_interact	105866.18	53972.80	37749.80	38574.81
stratified_standardization	LinearRegression_degree2	201053.48	103618.34	93842.31	43936.11

Table 13: Twins Estimators sorted by PEHE

meta-estimator	outcome_model	mean_pehe	ate_rmse	ate_abs_bias	ate_std_error
standardization	KernelRidge	0.04	0.01	0.01	0.01
standardization	ElasticNet	0.04	0.02	0.01	0.01
standardization	Lasso	0.04	0.02	0.01	0.01
standardization	LinearRegression	0.04	0.02	0.01	0.01
standardization	Ridge	0.04	0.02	0.01	0.01
standardization	Standardized_SVM_rbf	0.06	0.01	0.01	0.01
standardization	SVM_rbf	0.06	0.05	0.04	0.01
standardization	DecisionTree	0.08	0.04	0.02	0.04
standardization	kNN	0.08	0.06	0.06	0.00
standardization	LinearSVM	0.08	0.07	0.07	0.00
standardization	SVM_sigmoid	0.08	0.07	0.07	0.00
standardization	Standardized_LinearSVM	0.08	0.07	0.07	0.00
standardization	Standardized_SVM_sigmoid	0.08	0.07	0.07	0.00
stratified_standardization	LinearSVM	0.08	0.07	0.07	0.00
stratified_standardization	SVM_sigmoid	0.08	0.07	0.07	0.00
stratified_standardization	Standardized_LinearSVM	0.08	0.07	0.07	0.00
stratified_standardization	Standardized_SVM_sigmoid	0.08	0.07	0.07	0.00
stratified_standardization	Lasso	0.10	0.02	0.01	0.01
stratified_standardization	KernelRidge	0.11	0.01	0.01	0.01
stratified_standardization	ElasticNet	0.11	0.02	0.01	0.01
stratified_standardization	LinearRegression	0.11	0.02	0.01	0.01
stratified_standardization	Ridge	0.11	0.02	0.01	0.01
stratified_standardization	SVM_rbf	0.11	0.04	0.04	0.01
stratified_standardization	DecisionTree	0.12	0.03	0.03	0.01
stratified_standardization	kNN	0.13	0.05	0.05	0.02
stratified_standardization	Standardized_SVM_rbf	0.23	0.02	0.01	0.01
standardization	LinearRegression_degree2	122823195.44	10639314.01	1478780.90	10536043.36
stratified_standardization	LinearRegression_interact	152626640.11	3006941.41	646047.20	2936719.20
standardization	LinearRegression_interact	168756300.74	3341204.71	1965661.75	2701818.43
stratified_standardization	LinearRegression_degree2	403823525.32	3928023.92	1505613.68	3628015.90

587 **D.1 Hyperparameter Selection in Standardization Estimators**

Table 14: Lalonde CPS Outcome Model Correlations

meta-estimator	outcome_model	spearman	kendall	pearson	prob_better_better	prob_better_or_equal	causal_score
standardization	DecisionTree	0.81	0.73	0.89	0.75	0.92	ate_rmse
standardization	DecisionTree	0.83	0.73	0.91	0.75	0.92	mean_pehe
standardization	ElasticNet	0.86	0.69	0.59	0.76	0.89	ate_rmse
standardization	ElasticNet	0.84	0.64	0.58	0.73	0.87	mean_pehe
standardization	KernelRidge	0.81	0.60	0.74	0.80	0.80	ate_rmse
standardization	KernelRidge	0.49	0.29	0.70	0.64	0.64	mean_pehe
standardization	Lasso	0.77	0.64	0.63	0.69	0.89	ate_rmse
standardization	Lasso	0.57	0.30	0.62	0.53	0.73	mean_pehe
standardization	LinearSVM	0.96	0.91	0.86	0.96	0.96	ate_rmse
standardization	LinearSVM	0.96	0.91	0.85	0.96	0.96	mean_pehe
standardization	Ridge	0.70	0.62	0.75	0.71	0.84	ate_rmse
standardization	Ridge	0.40	0.22	0.72	0.53	0.67	mean_pehe
standardization	SVM_rbf				0.00	1.00	ate_rmse
standardization	SVM_rbf				0.00	1.00	mean_pehe
standardization	SVM_sigmoid	0.52	0.45	1.00	0.20	1.00	ate_rmse
standardization	SVM_sigmoid	0.52	0.45	1.00	0.20	1.00	mean_pehe
standardization	Standardized_LinearSVM	0.37	0.38	0.41	0.69	0.69	ate_rmse
standardization	Standardized_LinearSVM	0.62	0.51	0.43	0.76	0.76	mean_pehe
standardization	Standardized_SVM_rbf	0.98	0.91	0.96	0.96	0.96	ate_rmse
standardization	Standardized_SVM_rbf	1.00	1.00	0.98	1.00	1.00	mean_pehe
standardization	Standardized_SVM_sigmoid	1.00	1.00	1.00	1.00	1.00	ate_rmse
standardization	Standardized_SVM_sigmoid	1.00	1.00	1.00	1.00	1.00	mean_pehe
standardization	kNN	0.85	0.73	0.87	0.87	0.87	ate_rmse
standardization	kNN	0.50	0.38	0.90	0.69	0.69	mean_pehe
stratified_standardization	DecisionTree	0.70	0.61	0.59	0.81	0.81	ate_rmse
stratified_standardization	DecisionTree	0.73	0.67	0.80	0.83	0.83	mean_pehe
stratified_standardization	ElasticNet	0.88	0.73	0.83	0.87	0.87	ate_rmse
stratified_standardization	ElasticNet	-0.02	0.02	0.28	0.51	0.51	mean_pehe
stratified_standardization	Lasso	0.86	0.74	0.86	0.82	0.89	ate_rmse
stratified_standardization	Lasso	-0.50	-0.19	-0.21	0.38	0.44	mean_pehe
stratified_standardization	LinearSVM	0.87	0.78	0.52	0.89	0.89	ate_rmse
stratified_standardization	LinearSVM	0.87	0.78	0.66	0.89	0.89	mean_pehe
stratified_standardization	Ridge	0.31	0.17	0.68	0.51	0.64	ate_rmse
stratified_standardization	Ridge	-0.36	-0.12	-0.32	0.38	0.51	mean_pehe
stratified_standardization	SVM_rbf	0.32	0.11	0.78	0.56	0.56	ate_rmse
stratified_standardization	SVM_rbf	0.84	0.73	0.99	0.87	0.87	mean_pehe
stratified_standardization	SVM_sigmoid	1.00	1.00	1.00	1.00	1.00	ate_rmse
stratified_standardization	SVM_sigmoid	1.00	1.00	1.00	1.00	1.00	mean_pehe
stratified_standardization	Standardized_LinearSVM	0.78	0.69	0.61	0.84	0.84	ate_rmse
stratified_standardization	Standardized_LinearSVM	0.77	0.64	0.77	0.82	0.82	mean_pehe
stratified_standardization	Standardized_SVM_rbf	0.49	0.16	0.77	0.58	0.58	ate_rmse
stratified_standardization	Standardized_SVM_rbf	0.94	0.87	0.97	0.93	0.93	mean_pehe
stratified_standardization	Standardized_SVM_sigmoid	0.33	0.11	1.00	0.56	0.56	ate_rmse
stratified_standardization	Standardized_SVM_sigmoid	0.99	0.96	1.00	0.98	0.98	mean_pehe
stratified_standardization	kNN	0.90	0.79	0.89	0.89	0.89	ate_rmse
stratified_standardization	kNN	1.00	1.00	0.99	1.00	1.00	mean_pehe

Table 15: Lalonde PSID Outcome Model Correlations

meta-estimator	outcome_model	spearman	kendall	pearson	prob_better_better	prob_better_or_equal	causal_score
standardization	DecisionTree	-0.13	-0.08	-0.01	0.44	0.47	ate_rmse
standardization	DecisionTree	0.87	0.70	0.93	0.83	0.86	mean_pehe
standardization	ElasticNet	0.99	0.97	0.99	0.93	1.00	ate_rmse
standardization	ElasticNet	0.99	0.97	0.99	0.93	1.00	mean_pehe
standardization	KernelRidge	0.84	0.69	0.28	0.84	0.84	ate_rmse
standardization	KernelRidge	0.81	0.64	0.28	0.82	0.82	mean_pehe
standardization	Lasso	0.98	0.93	0.99	0.87	1.00	ate_rmse
standardization	Lasso	0.98	0.94	0.99	0.87	1.00	mean_pehe
standardization	LinearSVM	0.94	0.82	0.87	0.91	0.91	ate_rmse
standardization	LinearSVM	0.94	0.82	0.87	0.91	0.91	mean_pehe
standardization	Ridge	0.99	0.97	0.99	0.93	1.00	ate_rmse
standardization	Ridge	0.99	0.97	1.00	0.93	1.00	mean_pehe
standardization	SVM_rbf				0.00	1.00	ate_rmse
standardization	SVM_rbf				0.00	1.00	mean_pehe
standardization	SVM_sigmoid	-0.52	-0.45	-1.00	0.00	0.80	ate_rmse
standardization	SVM_sigmoid				0.00	1.00	mean_pehe
standardization	Standardized_LinearSVM	-0.56	-0.29	-0.72	0.36	0.36	ate_rmse
standardization	Standardized_LinearSVM	-0.56	-0.29	-0.72	0.36	0.36	mean_pehe
standardization	Standardized_SVM_rbf	1.00	1.00	0.99	1.00	1.00	ate_rmse
standardization	Standardized_SVM_rbf	1.00	1.00	1.00	1.00	1.00	mean_pehe
standardization	Standardized_SVM_sigmoid	1.00	1.00	1.00	1.00	1.00	ate_rmse
standardization	Standardized_SVM_sigmoid	0.75	0.73	1.00	0.87	0.87	mean_pehe
standardization	kNN	0.71	0.56	0.74	0.78	0.78	ate_rmse
standardization	kNN	0.95	0.87	0.97	0.93	0.93	mean_pehe
stratified_standardization	DecisionTree	-0.67	-0.50	-0.66	0.25	0.25	ate_rmse
stratified_standardization	DecisionTree	0.87	0.72	0.94	0.86	0.86	mean_pehe
stratified_standardization	ElasticNet	-0.95	-0.87	-0.89	0.07	0.07	ate_rmse
stratified_standardization	ElasticNet	0.88	0.78	0.97	0.89	0.89	mean_pehe
stratified_standardization	KernelRidge	-0.36	-0.24	-0.05	0.38	0.38	ate_rmse
stratified_standardization	KernelRidge	0.50	0.33	-0.03	0.67	0.67	mean_pehe
stratified_standardization	Lasso	-0.97	-0.88	-0.95	0.04	0.11	ate_rmse
stratified_standardization	Lasso	0.99	0.98	0.89	0.93	1.00	mean_pehe
stratified_standardization	LinearSVM	0.77	0.60	0.46	0.80	0.80	ate_rmse
stratified_standardization	LinearSVM	0.77	0.60	0.57	0.80	0.80	mean_pehe
stratified_standardization	Ridge	-0.98	-0.92	-0.91	0.02	0.09	ate_rmse
stratified_standardization	Ridge	0.88	0.74	0.97	0.82	0.89	mean_pehe
stratified_standardization	SVM_rbf	0.48	0.42	0.50	0.71	0.71	ate_rmse
stratified_standardization	SVM_rbf	0.99	0.96	1.00	0.98	0.98	mean_pehe
stratified_standardization	SVM_sigmoid	1.00	1.00	1.00	1.00	1.00	ate_rmse
stratified_standardization	SVM_sigmoid	1.00	1.00	1.00	1.00	1.00	mean_pehe
stratified_standardization	Standardized_LinearSVM	0.82	0.69	0.95	0.84	0.84	ate_rmse
stratified_standardization	Standardized_LinearSVM	0.89	0.78	0.96	0.89	0.89	mean_pehe
stratified_standardization	Standardized_SVM_rbf	0.53	0.47	0.35	0.73	0.73	ate_rmse
stratified_standardization	Standardized_SVM_rbf	0.99	0.96	1.00	0.98	0.98	mean_pehe
stratified_standardization	Standardized_SVM_sigmoid	0.88	0.82	1.00	0.91	0.91	ate_rmse
stratified_standardization	Standardized_SVM_sigmoid	0.75	0.73	1.00	0.87	0.87	mean_pehe
stratified_standardization	kNN	0.86	0.79	0.66	0.89	0.89	ate_rmse
stratified_standardization	kNN	1.00	1.00	0.96	1.00	1.00	mean_pehe

Table 16: Twins Outcome Model Correlations

meta-estimator	outcome_model	spearman	kendall	pearson	prob_better_better	prob_better_or_equal	causal_score
standardization	DecisionTree	-0.74	-0.59	-0.66	0.08	0.47	ate_rmse
standardization	DecisionTree	0.85	0.77	0.96	0.81	0.92	mean_pehe
standardization	ElasticNet	0.90	0.84	0.97	0.47	1.00	ate_rmse
standardization	ElasticNet	0.91	0.88	0.98	0.47	1.00	mean_pehe
standardization	KernelRidge	0.92	0.89	0.98	0.53	1.00	ate_rmse
standardization	KernelRidge	1.00	1.00	0.99	0.53	1.00	mean_pehe
standardization	Lasso	1.00	1.00	1.00	0.51	1.00	ate_rmse
standardization	Lasso	0.99	0.96	0.98	0.47	1.00	mean_pehe
standardization	LinearSVM	-0.93	-0.85	-0.89	0.00	0.44	ate_rmse
standardization	LinearSVM	-0.93	-0.85	-0.89	0.00	0.44	mean_pehe
standardization	Ridge	0.86	0.84	0.98	0.38	1.00	ate_rmse
standardization	Ridge	1.00	1.00	1.00	0.38	1.00	mean_pehe
standardization	SVM_rbf	-0.67	-0.56	-0.78	0.09	0.53	ate_rmse
standardization	SVM_rbf	0.93	0.89	0.96	0.69	1.00	mean_pehe
standardization	SVM_sigmoid	0.99	0.97	1.00	0.93	1.00	ate_rmse
standardization	SVM_sigmoid	0.99	0.97	1.00	0.93	1.00	mean_pehe
standardization	Standardized_LinearSVM	0.76	0.62	0.97	0.53	0.89	ate_rmse
standardization	Standardized_LinearSVM	0.76	0.62	0.94	0.53	0.89	mean_pehe
standardization	Standardized_SVM_rbf	-0.28	-0.25	-0.23	0.27	0.56	ate_rmse
standardization	Standardized_SVM_rbf	1.00	1.00	0.98	0.71	1.00	mean_pehe
standardization	Standardized_SVM_sigmoid	0.79	0.65	1.00	0.78	0.84	ate_rmse
standardization	Standardized_SVM_sigmoid	0.92	0.84	1.00	0.87	0.93	mean_pehe
standardization	kNN	0.14	0.12	0.43	0.18	0.89	ate_rmse
standardization	kNN	0.76	0.70	0.43	0.69	0.91	mean_pehe
stratified_standardization	DecisionTree				0.00	1.00	ate_rmse
stratified_standardization	DecisionTree	0.85	0.77	0.93	0.81	0.92	mean_pehe
stratified_standardization	ElasticNet	1.00	1.00	1.00	0.60	1.00	ate_rmse
stratified_standardization	ElasticNet	-0.90	-0.82	-0.83	0.00	0.53	mean_pehe
stratified_standardization	KernelRidge	0.81	0.69	0.66	0.44	0.96	ate_rmse
stratified_standardization	KernelRidge	-0.86	-0.81	-0.90	0.00	0.49	mean_pehe
stratified_standardization	Lasso	1.00	1.00	1.00	0.51	1.00	ate_rmse
stratified_standardization	Lasso	-0.86	-0.81	-0.87	0.00	0.64	mean_pehe
stratified_standardization	LinearSVM	0.93	0.85	0.90	0.64	1.00	ate_rmse
stratified_standardization	LinearSVM	0.91	0.82	0.90	0.64	1.00	mean_pehe
stratified_standardization	Ridge	0.86	0.84	0.96	0.38	1.00	ate_rmse
stratified_standardization	Ridge	-0.72	-0.64	-0.75	0.00	0.64	mean_pehe
stratified_standardization	SVM_rbf	-0.78	-0.65	-0.80	0.09	0.44	ate_rmse
stratified_standardization	SVM_rbf	0.79	0.68	0.93	0.62	0.91	mean_pehe
stratified_standardization	SVM_sigmoid	0.95	0.87	1.00	0.93	0.93	ate_rmse
stratified_standardization	SVM_sigmoid	0.99	0.96	1.00	0.98	0.98	mean_pehe
stratified_standardization	Standardized_LinearSVM	0.72	0.57	0.90	0.53	0.89	ate_rmse
stratified_standardization	Standardized_LinearSVM	0.93	0.85	0.98	0.64	1.00	mean_pehe
stratified_standardization	Standardized_SVM_rbf	-0.56	-0.44	-0.34	0.16	0.56	ate_rmse
stratified_standardization	Standardized_SVM_rbf	0.86	0.75	0.61	0.62	0.91	mean_pehe
stratified_standardization	Standardized_SVM_sigmoid	0.92	0.84	1.00	0.87	0.93	ate_rmse
stratified_standardization	Standardized_SVM_sigmoid	0.99	0.98	1.00	0.93	1.00	mean_pehe
stratified_standardization	kNN	0.36	0.32	0.24	0.13	1.00	ate_rmse
stratified_standardization	kNN	0.97	0.92	0.98	0.84	1.00	mean_pehe

Table 17: Lalonde CPS Propensity Score Model Correlations

meta-estimator	prop_score_model	spearman	kendall	pearson	prob_better_better	prob_better_or_equal	class_score
ipw	DecisionTree	-0.25	-0.20	-0.10	0.36	0.44	mean_test_f1
ipw	DecisionTree	0.71	0.51	0.69	0.72	0.78	mean_test_average_precision
ipw	DecisionTree	-0.24	-0.20	-0.08	0.36	0.44	mean_test_balanced_accuracy
ipw	LogisticRegression_l1_liblinear	0.69	0.53	0.93	0.67	0.82	mean_test_f1
ipw	LogisticRegression_l1_liblinear	0.88	0.80	0.98	0.64	1.00	mean_test_average_precision
ipw	LogisticRegression_l1_liblinear	0.67	0.50	0.89	0.62	0.82	mean_test_balanced_accuracy
ipw	LogisticRegression_l1_saga				0.00	1.00	mean_test_f1
ipw	LogisticRegression_l1_saga	-0.78	-0.73	-0.74	0.00	0.64	mean_test_average_precision
ipw	LogisticRegression_l1_saga				0.00	1.00	mean_test_balanced_accuracy
ipw	LogisticRegression_l2	0.86	0.74	0.89	0.78	0.91	mean_test_f1
ipw	LogisticRegression_l2	0.80	0.68	0.97	0.73	0.89	mean_test_average_precision
ipw	LogisticRegression_l2	0.82	0.70	0.87	0.71	0.91	mean_test_balanced_accuracy
ipw	LogisticRegression_l2_liblinear	0.69	0.58	0.75	0.73	0.82	mean_test_f1
ipw	LogisticRegression_l2_liblinear	0.75	0.68	0.93	0.73	0.89	mean_test_average_precision
ipw	LogisticRegression_l2_liblinear	0.74	0.65	0.72	0.64	0.91	mean_test_balanced_accuracy
ipw	SVM_sigmoid				0.00	1.00	mean_test_f1
ipw	SVM_sigmoid	-0.04	0.00	0.01	0.44	0.56	mean_test_average_precision
ipw	SVM_sigmoid				0.00	1.00	mean_test_balanced_accuracy
ipw	Standardized_SVM_rbf	-0.22	-0.17	0.11	0.36	0.49	mean_test_f1
ipw	Standardized_SVM_rbf	0.80	0.61	0.65	0.78	0.82	mean_test_average_precision
ipw	Standardized_SVM_rbf	-0.38	-0.42	0.08	0.22	0.40	mean_test_balanced_accuracy
ipw	Standardized_SVM_sigmoid	-0.49	-0.28	-0.39	0.27	0.49	mean_test_f1
ipw	Standardized_SVM_sigmoid	0.10	0.09	0.37	0.31	0.76	mean_test_average_precision
ipw	Standardized_SVM_sigmoid	-0.21	-0.18	-0.21	0.20	0.67	mean_test_balanced_accuracy
ipw	kNN	0.89	0.84	0.81	0.70	1.00	mean_test_f1
ipw	kNN	0.82	0.74	0.80	0.80	0.90	mean_test_average_precision
ipw	kNN	0.89	0.84	0.91	0.70	1.00	mean_test_balanced_accuracy
ipw_trimeps.01	DecisionTree	-0.83	-0.67	-0.62	0.14	0.22	mean_test_f1
ipw_trimeps.01	DecisionTree	0.40	0.23	0.56	0.58	0.64	mean_test_average_precision
ipw_trimeps.01	DecisionTree	-0.84	-0.67	-0.61	0.14	0.22	mean_test_balanced_accuracy
ipw_trimeps.01	LogisticRegression_l1_liblinear	0.43	0.24	0.89	0.53	0.69	mean_test_f1
ipw_trimeps.01	LogisticRegression_l1_liblinear	0.56	0.36	0.95	0.47	0.82	mean_test_average_precision
ipw_trimeps.01	LogisticRegression_l1_liblinear	0.47	0.30	0.83	0.53	0.73	mean_test_balanced_accuracy
ipw_trimeps.01	LogisticRegression_l1_saga				0.00	1.00	mean_test_f1
ipw_trimeps.01	LogisticRegression_l1_saga	-0.78	-0.73	-0.74	0.00	0.64	mean_test_average_precision
ipw_trimeps.01	LogisticRegression_l1_saga				0.00	1.00	mean_test_balanced_accuracy
ipw_trimeps.01	LogisticRegression_l2	0.90	0.79	0.93	0.80	0.93	mean_test_f1
ipw_trimeps.01	LogisticRegression_l2	0.82	0.73	0.97	0.76	0.91	mean_test_average_precision
ipw_trimeps.01	LogisticRegression_l2	0.87	0.75	0.92	0.73	0.93	mean_test_balanced_accuracy
ipw_trimeps.01	LogisticRegression_l2_liblinear	0.63	0.49	0.84	0.69	0.78	mean_test_f1
ipw_trimeps.01	LogisticRegression_l2_liblinear	0.66	0.53	0.97	0.67	0.82	mean_test_average_precision
ipw_trimeps.01	LogisticRegression_l2_liblinear	0.55	0.39	0.79	0.53	0.80	mean_test_balanced_accuracy
ipw_trimeps.01	SVM_sigmoid				0.00	1.00	mean_test_f1
ipw_trimeps.01	SVM_sigmoid	0.11	0.12	0.19	0.49	0.62	mean_test_average_precision
ipw_trimeps.01	SVM_sigmoid				0.00	1.00	mean_test_balanced_accuracy
ipw_trimeps.01	Standardized_SVM_rbf	-0.57	-0.36	-0.61	0.27	0.40	mean_test_f1
ipw_trimeps.01	Standardized_SVM_rbf	0.75	0.61	0.84	0.78	0.82	mean_test_average_precision
ipw_trimeps.01	Standardized_SVM_rbf	-0.72	-0.61	-0.68	0.13	0.31	mean_test_balanced_accuracy
ipw_trimeps.01	Standardized_SVM_sigmoid	-0.57	-0.33	-0.43	0.24	0.47	mean_test_f1
ipw_trimeps.01	Standardized_SVM_sigmoid	-0.03	-0.03	0.26	0.27	0.71	mean_test_average_precision
ipw_trimeps.01	Standardized_SVM_sigmoid	-0.36	-0.30	-0.22	0.16	0.62	mean_test_balanced_accuracy
ipw_trimeps.01	kNN	-0.24	-0.12	-0.49	0.40	0.49	mean_test_f1
ipw_trimeps.01	kNN	0.71	0.58	0.90	0.78	0.80	mean_test_average_precision
ipw_trimeps.01	kNN	-0.30	-0.16	-0.67	0.38	0.47	mean_test_balanced_accuracy



Table 18: Lalonde PSID Propensity Score Model Correlations

meta-estimator	prop_score_model	spearman	kendall	pearson	prob_better_better	prob_better_or_equal	class_score
ipw	DecisionTree	0.88	0.70	0.81	0.83	0.86	mean_test_f1
ipw	DecisionTree	-0.90	-0.74	-0.81	0.11	0.17	mean_test_average_precision
ipw	DecisionTree	0.84	0.67	0.80	0.78	0.86	mean_test_balanced_accuracy
ipw	LogisticRegression_l1_liblinear	-0.86	-0.73	-0.72	0.04	0.36	mean_test_f1
ipw	LogisticRegression_l1_liblinear	-0.22	-0.11	-0.47	0.27	0.64	mean_test_average_precision
ipw	LogisticRegression_l1_liblinear	-0.87	-0.75	-0.73	0.00	0.44	mean_test_balanced_accuracy
ipw	LogisticRegression_l1_saga				0.00	1.00	mean_test_f1
ipw	LogisticRegression_l1_saga	-0.86	-0.81	-0.74	0.00	0.49	mean_test_average_precision
ipw	LogisticRegression_l1_saga				0.00	1.00	mean_test_balanced_accuracy
ipw	LogisticRegression_l2				0.00	1.00	mean_test_f1
ipw	LogisticRegression_l2	0.24	0.21	0.34	0.27	0.87	mean_test_average_precision
ipw	LogisticRegression_l2				0.00	1.00	mean_test_balanced_accuracy
ipw	LogisticRegression_l2_liblinear				0.00	1.00	mean_test_f1
ipw	LogisticRegression_l2_liblinear	0.83	0.77	0.93	0.51	1.00	mean_test_average_precision
ipw	LogisticRegression_l2_liblinear				0.00	1.00	mean_test_balanced_accuracy
ipw	SVM_sigmoid	0.37	0.37	0.10	0.51	0.80	mean_test_f1
ipw	SVM_sigmoid	-0.49	-0.42	-0.34	0.09	0.62	mean_test_average_precision
ipw	SVM_sigmoid	0.43	0.37	0.22	0.40	0.87	mean_test_balanced_accuracy
ipw	Standardized_SVM_rbf	0.68	0.53	0.63	0.62	0.84	mean_test_f1
ipw	Standardized_SVM_rbf	0.36	0.26	0.86	0.58	0.67	mean_test_average_precision
ipw	Standardized_SVM_rbf	0.67	0.51	0.62	0.60	0.84	mean_test_balanced_accuracy
ipw	Standardized_SVM_sigmoid	0.06	0.08	0.41	0.42	0.64	mean_test_f1
ipw	Standardized_SVM_sigmoid	-0.52	-0.42	-0.54	0.18	0.47	mean_test_average_precision
ipw	Standardized_SVM_sigmoid	0.16	0.16	0.41	0.42	0.71	mean_test_balanced_accuracy
ipw	kNN	0.90	0.76	0.91	0.86	0.89	mean_test_f1
ipw	kNN	-0.84	-0.74	-0.78	0.11	0.18	mean_test_average_precision
ipw	kNN	0.87	0.74	0.88	0.82	0.89	mean_test_balanced_accuracy
ipw_trimeps.01	DecisionTree	0.81	0.65	0.61	0.81	0.83	mean_test_f1
ipw_trimeps.01	DecisionTree	-0.87	-0.69	-0.68	0.14	0.19	mean_test_average_precision
ipw_trimeps.01	DecisionTree	0.77	0.61	0.61	0.75	0.83	mean_test_balanced_accuracy
ipw_trimeps.01	LogisticRegression_l1_liblinear	-0.89	-0.79	-0.80	0.02	0.33	mean_test_f1
ipw_trimeps.01	LogisticRegression_l1_liblinear	-0.22	-0.11	-0.44	0.27	0.64	mean_test_average_precision
ipw_trimeps.01	LogisticRegression_l1_liblinear	-0.87	-0.75	-0.81	0.00	0.44	mean_test_balanced_accuracy
ipw_trimeps.01	LogisticRegression_l1_saga				0.00	1.00	mean_test_f1
ipw_trimeps.01	LogisticRegression_l1_saga	-0.91	-0.88	-0.74	0.00	0.49	mean_test_average_precision
ipw_trimeps.01	LogisticRegression_l1_saga				0.00	1.00	mean_test_balanced_accuracy
ipw_trimeps.01	LogisticRegression_l2				0.00	1.00	mean_test_f1
ipw_trimeps.01	LogisticRegression_l2	0.24	0.21	0.33	0.27	0.87	mean_test_average_precision
ipw_trimeps.01	LogisticRegression_l2				0.00	1.00	mean_test_balanced_accuracy
ipw_trimeps.01	LogisticRegression_l2_liblinear				0.00	1.00	mean_test_f1
ipw_trimeps.01	LogisticRegression_l2_liblinear	0.86	0.81	0.92	0.51	1.00	mean_test_average_precision
ipw_trimeps.01	LogisticRegression_l2_liblinear				0.00	1.00	mean_test_balanced_accuracy
ipw_trimeps.01	SVM_sigmoid	0.20	0.19	-0.30	0.42	0.73	mean_test_f1
ipw_trimeps.01	SVM_sigmoid	-0.19	-0.16	0.09	0.18	0.71	mean_test_average_precision
ipw_trimeps.01	SVM_sigmoid	0.29	0.25	-0.20	0.36	0.82	mean_test_balanced_accuracy
ipw_trimeps.01	Standardized_SVM_rbf	0.68	0.53	0.59	0.62	0.84	mean_test_f1
ipw_trimeps.01	Standardized_SVM_rbf	0.36	0.26	0.88	0.58	0.67	mean_test_average_precision
ipw_trimeps.01	Standardized_SVM_rbf	0.67	0.51	0.59	0.60	0.84	mean_test_balanced_accuracy
ipw_trimeps.01	Standardized_SVM_sigmoid	-0.06	-0.03	0.40	0.38	0.60	mean_test_f1
ipw_trimeps.01	Standardized_SVM_sigmoid	-0.22	-0.16	-0.54	0.29	0.58	mean_test_average_precision
ipw_trimeps.01	Standardized_SVM_sigmoid	0.00	0.05	0.40	0.38	0.67	mean_test_balanced_accuracy
ipw_trimeps.01	kNN	0.66	0.52	0.74	0.73	0.78	mean_test_f1
ipw_trimeps.01	kNN	-0.65	-0.52	-0.55	0.22	0.27	mean_test_average_precision
ipw_trimeps.01	kNN	0.59	0.46	0.63	0.69	0.76	mean_test_balanced_accuracy

Table 19: Twins Propensity Score Model Correlations

meta-estimator	prop_score_model	spearman	kendall	pearson	prob_better_better	prob_better_or_equal	class_score
ipw	DecisionTree	-0.63	-0.57	-0.55	0.03	0.58	mean_test_f1
ipw	DecisionTree	-0.27	-0.26	-0.28	0.06	0.81	mean_test_average_precision
ipw	DecisionTree	0.56	0.54	0.58	0.25	1.00	mean_test_balanced_accuracy
ipw	LogisticRegression_l1_liblinear	0.57	0.56	0.59	0.18	1.00	mean_test_f1
ipw	LogisticRegression_l1_liblinear	0.87	0.83	0.97	0.51	1.00	mean_test_average_precision
ipw	LogisticRegression_l1_liblinear				0.00	1.00	mean_test_balanced_accuracy
ipw	LogisticRegression_l1_saga				0.00	1.00	mean_test_f1
ipw	LogisticRegression_l1_saga	0.83	0.76	0.94	0.47	1.00	mean_test_average_precision
ipw	LogisticRegression_l1_saga	0.90	0.85	0.73	0.53	1.00	mean_test_balanced_accuracy
ipw	LogisticRegression_l2				0.00	1.00	mean_test_f1
ipw	LogisticRegression_l2	1.00	1.00	0.99	0.38	1.00	mean_test_average_precision
ipw	LogisticRegression_l2	0.25	0.24	0.24	0.09	1.00	mean_test_balanced_accuracy
ipw	LogisticRegression_l2_liblinear				0.00	1.00	mean_test_f1
ipw	LogisticRegression_l2_liblinear	0.31	0.22	0.67	0.36	0.80	mean_test_average_precision
ipw	LogisticRegression_l2_liblinear				0.00	1.00	mean_test_balanced_accuracy
ipw	SVM_sigmoid				0.00	1.00	mean_test_f1
ipw	SVM_sigmoid				0.00	1.00	mean_test_average_precision
ipw	SVM_sigmoid				0.00	1.00	mean_test_balanced_accuracy
ipw	Standardized_SVM_rbf	0.96	0.93	0.90	0.56	1.00	mean_test_f1
ipw	Standardized_SVM_rbf	0.38	0.37	0.30	0.22	0.96	mean_test_average_precision
ipw	Standardized_SVM_rbf	-0.93	-0.86	-0.86	0.00	0.44	mean_test_balanced_accuracy
ipw	Standardized_SVM_sigmoid	0.37	0.35	0.34	0.13	1.00	mean_test_f1
ipw	Standardized_SVM_sigmoid	0.60	0.58	0.87	0.20	1.00	mean_test_average_precision
ipw	Standardized_SVM_sigmoid	-0.31	-0.30	-0.30	0.00	0.89	mean_test_balanced_accuracy
ipw	kNN	0.45	0.41	0.61	0.31	0.93	mean_test_f1
ipw	kNN	0.00	0.00	-0.05	0.13	0.87	mean_test_average_precision
ipw	kNN	-0.76	-0.76	-0.76	0.00	0.69	mean_test_balanced_accuracy
ipw_trimeps.01	DecisionTree	-0.63	-0.57	-0.55	0.03	0.58	mean_test_f1
ipw_trimeps.01	DecisionTree	-0.27	-0.26	-0.28	0.06	0.81	mean_test_average_precision
ipw_trimeps.01	DecisionTree	0.56	0.54	0.58	0.25	1.00	mean_test_balanced_accuracy
ipw_trimeps.01	LogisticRegression_l1_liblinear	0.57	0.56	0.59	0.18	1.00	mean_test_f1
ipw_trimeps.01	LogisticRegression_l1_liblinear	0.87	0.83	0.97	0.51	1.00	mean_test_average_precision
ipw_trimeps.01	LogisticRegression_l1_liblinear				0.00	1.00	mean_test_balanced_accuracy
ipw_trimeps.01	LogisticRegression_l1_saga				0.00	1.00	mean_test_f1
ipw_trimeps.01	LogisticRegression_l1_saga	0.83	0.76	0.94	0.47	1.00	mean_test_average_precision
ipw_trimeps.01	LogisticRegression_l1_saga	0.90	0.85	0.73	0.53	1.00	mean_test_balanced_accuracy
ipw_trimeps.01	LogisticRegression_l2				0.00	1.00	mean_test_f1
ipw_trimeps.01	LogisticRegression_l2	1.00	1.00	0.99	0.38	1.00	mean_test_average_precision
ipw_trimeps.01	LogisticRegression_l2	0.25	0.24	0.24	0.09	1.00	mean_test_balanced_accuracy
ipw_trimeps.01	LogisticRegression_l2_liblinear				0.00	1.00	mean_test_f1
ipw_trimeps.01	LogisticRegression_l2_liblinear	0.31	0.22	0.67	0.36	0.80	mean_test_average_precision
ipw_trimeps.01	LogisticRegression_l2_liblinear				0.00	1.00	mean_test_balanced_accuracy
ipw_trimeps.01	SVM_sigmoid				0.00	1.00	mean_test_f1
ipw_trimeps.01	SVM_sigmoid				0.00	1.00	mean_test_average_precision
ipw_trimeps.01	SVM_sigmoid				0.00	1.00	mean_test_balanced_accuracy
ipw_trimeps.01	Standardized_SVM_rbf	0.96	0.93	0.90	0.56	1.00	mean_test_f1
ipw_trimeps.01	Standardized_SVM_rbf	0.38	0.37	0.30	0.22	0.96	mean_test_average_precision
ipw_trimeps.01	Standardized_SVM_rbf	-0.93	-0.86	-0.86	0.00	0.44	mean_test_balanced_accuracy
ipw_trimeps.01	Standardized_SVM_sigmoid	0.37	0.35	0.34	0.13	1.00	mean_test_f1
ipw_trimeps.01	Standardized_SVM_sigmoid	0.60	0.58	0.87	0.20	1.00	mean_test_average_precision
ipw_trimeps.01	Standardized_SVM_sigmoid	-0.31	-0.30	-0.30	0.00	0.89	mean_test_balanced_accuracy
ipw_trimeps.01	kNN	0.45	0.41	0.61	0.31	0.93	mean_test_f1
ipw_trimeps.01	kNN	0.00	0.00	-0.05	0.13	0.87	mean_test_average_precision
ipw_trimeps.01	kNN	-0.76	-0.76	-0.76	0.00	0.69	mean_test_balanced_accuracy

589 **D.3 Model Selection**Table 20: Correlations for predictive RMSE metric with ATE RMSE and mean PEHE among **standardization** estimators whose hyperparameters were selected by **cross-validation on the predictive RMSE**.

dataset	spearman	kendall	pearson	prob_better_better	prob_better_or_equal	causal_score
lalonde_cps	0.02	0.01	-0.20	0.48	0.53	ate_rmse
lalonde_cps	0.03	0.02	-0.22	0.49	0.54	mean_pehe
lalonde_psid	-0.25	-0.16	-0.06	0.39	0.46	ate_rmse
lalonde_psid	0.04	0.03	-0.03	0.48	0.55	mean_pehe
twins	0.94	0.88	0.88	0.76	0.98	ate_rmse
twins	0.41	0.33	0.89	0.51	0.77	mean_pehe

Table 21: Correlations for three predictive classification metrics with ATE RMSE among **IPW** estimators whose hyperparameters were selected by **cross-validation on the balanced F score**.

dataset	spearman	kendall	pearson	prob_better	prob_better_or_eq	class_score
lalonge_cps	0.08	0.07	-0.05	0.50	0.57	mean_test_f1
lalonge_cps	0.39	0.27	0.03	0.54	0.71	mean_test_average_precision
lalonge_cps	-0.09	-0.07	-0.34	0.43	0.50	mean_test_balanced_accuracy
lalonge_psid	-0.19	-0.15	-0.49	0.37	0.50	mean_test_f1
lalonge_psid	-0.36	-0.21	-0.24	0.35	0.45	mean_test_average_precision
lalonge_psid	-0.16	-0.11	-0.57	0.36	0.54	mean_test_balanced_accuracy
twins	0.62	0.54	0.65	0.26	0.99	mean_test_f1
twins	0.42	0.32	0.04	0.48	0.78	mean_test_average_precision
twins	-0.12	-0.09	-0.56	0.20	0.75	mean_test_balanced_accuracy

Table 22: Correlations for three predictive classification metrics with ATE RMSE among **IPW** estimators whose hyperparameters were selected by **cross-validation on the average precision score**.

dataset	spearman	kendall	pearson	prob_better	prob_better_or_eq	class_score
lalonge_cps	-0.09	-0.07	-0.27	0.42	0.51	mean_test_f1
lalonge_cps	0.27	0.18	0.01	0.52	0.65	mean_test_average_precision
lalonge_cps	-0.11	-0.03	-0.43	0.44	0.53	mean_test_balanced_accuracy
lalonge_psid	-0.03	0.01	-0.12	0.34	0.67	mean_test_f1
lalonge_psid	0.40	0.33	0.60	0.59	0.72	mean_test_average_precision
lalonge_psid	-0.03	0.01	-0.14	0.33	0.68	mean_test_balanced_accuracy
twins	0.63	0.56	0.67	0.26	0.99	mean_test_f1
twins	0.84	0.72	0.28	0.58	0.94	mean_test_average_precision
twins	-0.20	-0.16	-0.60	0.15	0.75	mean_test_balanced_accuracy

Table 23: Correlations for three predictive classification metrics with ATE RMSE among **IPW** estimators whose hyperparameters were selected by **cross-validation on the balanced accuracy score**.

dataset	spearman	kendall	pearson	prob_better	prob_better_or_eq	class_score
lalonge_cps	-0.04	-0.02	-0.08	0.46	0.52	mean_test_f1
lalonge_cps	0.40	0.27	0.07	0.57	0.69	mean_test_average_precision
lalonge_cps	-0.13	-0.10	-0.33	0.41	0.50	mean_test_balanced_accuracy
lalonge_psid	-0.17	-0.12	-0.49	0.38	0.51	mean_test_f1
lalonge_psid	-0.38	-0.22	-0.24	0.35	0.44	mean_test_average_precision
lalonge_psid	-0.16	-0.11	-0.57	0.36	0.54	mean_test_balanced_accuracy
twins	0.78	0.67	0.43	0.57	0.95	mean_test_f1
twins	0.61	0.47	0.13	0.52	0.85	mean_test_average_precision
twins	-0.30	-0.24	-0.34	0.19	0.65	mean_test_balanced_accuracy

## Appendix E Causal Estimators

### E.1 Meta-Estimators

#### E.1.1 Standardization

We rewrite Equation 1 using the mean conditional outcome  $\mu(t, w)$  from Equation 3:

$$\tau = \mathbb{E}_W [\mathbb{E}[Y|T = 1, W] - \mathbb{E}[Y|T = 0, W]] \quad (4)$$

$$= \mathbb{E}_W [\mu(1, W) - \mu(0, W)] \quad (5)$$

594 By plugging in an arbitrary model  $\hat{\mu}$ , we get the following:

$$\hat{\tau} = \frac{1}{m} \sum_w (\hat{\mu}(1, w) - \hat{\mu}(0, w)) , \quad (6)$$

595 where  $m$  is the number of observations. This is a *meta-estimator* or *nonparametric estimator* of  
 596  $\tau$ , since we have not specified any parametric assumptions for  $\hat{\mu}$ .<sup>4</sup> We call the class of estimators  
 597 that result from using an arbitrary model for  $\hat{\mu}$  *standardization estimators* (Hernán & Robins, 2020).  
 598 These estimators also go by other names such as G-computation estimators (see, e.g., Snowden et al.,  
 599 2011) and S-learner (Künzel et al., 2019). The models  $\hat{\mu}(t, w)$  also provide estimates of CATEs by  
 600 simply plugging in to Equation 3:

$$\hat{\tau}(w) = \hat{\mu}(1, w) - \hat{\mu}(0, w) . \quad (7)$$

### 601 E.1.2 Stratified Standardization

602 Rather than modeling  $\mu$  with a single model  $\hat{\mu}$ , we could model  $\mu$  with a model for each value of  
 603 treatment. Specifically, we can model  $\mu(1, w)$  with a model  $\hat{\mu}_1(w)$ , and we can model  $\mu(0, w)$  with  
 604 another model  $\hat{\mu}_0(w)$ . Then, we can estimate ATEs and CATEs as follows:

$$\hat{\tau} = \frac{1}{m} \sum_w (\hat{\mu}_1(w) - \hat{\mu}_0(w)) \quad (8)$$

$$\hat{\tau}(w) = \hat{\mu}_1(w) - \hat{\mu}_0(w) . \quad (9)$$

605 We call the class of estimators that model  $\mu$  via  $\hat{\mu}_1(w)$  and  $\hat{\mu}_0(w)$  *stratified standardization estimators*.  
 606 This is also sometimes called the T-learner (Künzel et al., 2019).

### 607 E.1.3 Inverse probability weighting

608 We can also adjust for confounding by reweighting the population such that each observation is  
 609 weighted by  $\frac{1}{P(t_i|w_i)}$ . This reweighted population is referred to as the *pseudo-population* (Hernán &  
 610 Robins, 2020). This general technique (Horvitz & Thompson, 1952) is known as *inverse probability*  
 611 *weighting* (IPW). It can be shown that

$$\tau = \mathbb{E} \left[ \frac{I(T=1)Y}{P(T|W)} - \frac{I(T=0)Y}{P(T|W)} \right] , \quad (10)$$

612 where  $I$  denotes an indicator random variable, which takes the value 1 if its argument is true and takes  
 613 the value 0 otherwise. Now, a natural estimator is a plug-in estimator for this estimand. Note that  
 614  $P(T=1|W=w)$  is the propensity score  $e(w)$ . We can then estimate  $\tau$  using an arbitrary model  
 615  $\hat{e}(w)$  via the following plug-in estimator:

$$\hat{\tau} = \frac{1}{m} \sum_w \left( \frac{I(T=1)Y}{\hat{e}(w)} - \frac{I(T=0)Y}{1 - \hat{e}(w)} \right) \quad (11)$$

616 The weights  $\frac{1}{\hat{e}(w)}$  for the treated units and  $\frac{1}{1 - \hat{e}(w)}$  for the untreated units can sometimes be very large,  
 617 leading to high variance. Therefore, it is common to trim them. We consider estimators that trim  
 618 them when  $\hat{e}(w) < 0.01$  or when  $\hat{e}(w) > 0.99$ .

619 Another option is to “stabilize” the weight by multiplying them by fractions that are independent of  
 620  $W$  and sum to 1. A natural option is  $P(T)$ . This yields the *stabilized IPW estimator*:

$$\hat{\tau} = \frac{1}{m} \sum_w \left( \frac{I(T=1)P_T(1)Y}{\hat{e}(w)} - \frac{I(T=0)P_T(0)Y}{1 - \hat{e}(w)} \right) \quad (12)$$

## 621 E.2 Outcome and Propensity Score Models

622 We model  $\hat{\mu}(t, w)$ ,  $\hat{\mu}_1(w)$ ,  $\hat{\mu}_0(w)$ , and  $\hat{e}(w)$  using a large variety of models available in *scikit-learn*  
 623 (Pedregosa et al., 2011).

<sup>4</sup>Künzel et al. (2019) refer to nonparametric estimators as “meta-learners.”

624 We consider the following models for  $\hat{\mu}$ ,  $\hat{\mu}_1$ , and  $\hat{\mu}_0$ : ordinary least-squares (OLS) linear regression,  
625 lasso regression, ridge regression, elastic net, support-vector machines (SVMs) with radial basis  
626 function (RBF) kernels, SVMs with sigmoid kernels, linear SVMs, SVMs with standardized inputs  
627 (because SVMs are sensitive to input scale), kernel ridge regression, and decision trees.

628 We consider the following models for  $\hat{e}$ : vanilla logistic regression, logistic regression with L2  
629 regularization, logistic regression with L1 regularization, logistic regression with variants on the  
630 optimizer such as liblinear or SAGA,  $k$ -nearest neighbors (kNN), decision trees, Gaussian Naive  
631 Bayes, and quadratic discriminant analysis (QDA).