# Orthogonal Random Forest for Causal Inference

**Miruna Oprescu** [1]   **Vasilis Syrgkanis** [1]   **Zhiwei Steven Wu** [2]

## Abstract

We propose the *orthogonal random forest*, an algorithm that combines *Neyman-orthogonality* to reduce sensitivity with respect to estimation error of nuisance parameters with generalized random forests (Athey et al., 2017)—a flexible non-parametric method for statistical estimation of conditional moment models using random forests. We provide a consistency rate and establish asymptotic normality for our estimator. We show that under mild assumptions on the consistency rate of the nuisance estimator, we can achieve the same error rate as an oracle with a priori knowledge of these nuisance parameters. We show that when the nuisance functions have a locally sparse parametrization, then a local $\ell_1$-penalized regression achieves the required rate. We apply our method to estimate heterogeneous treatment effects from observational data with discrete treatments or continuous treatments, and we show that, unlike prior work, our method provably allows to control for a high-dimensional set of variables under standard sparsity conditions. We also provide a comprehensive empirical evaluation of our algorithm on both synthetic and real data.

## 1. Introduction

Many problems that arise in causal inference can be formulated in the language of conditional moment models: given a target feature $x$ find a solution $\theta_0(x)$ to a system of conditional moment equations

$$\mathbb{E}\left[\psi(Z; \theta, h_0(x, W)) \mid X = x\right] = 0, \quad (1)$$

given access to $n$ i.i.d. samples from the data generating distribution, where $\psi$ is a known score function and $h_0$ is an

unknown nuisance function that also needs to be estimated from data. Examples include non-parametric regression, heterogeneous treatment effect estimation, instrumental variable regression, local maximum likelihood estimation and estimation of structural econometric models.[1] The study of such conditional moment restriction problems has a long history in econometrics (see e.g. Newey (1993); Ai & Chen (2003); Chen & Pouzo (2009); Chernozhukov et al. (2015)).

In this general estimation problem, the main goal is to estimate the target parameter at a rate that is robust to the estimation error of the nuisance component. This allows the use of flexible models to fit the nuisance functions and enables asymptotically valid inference. Almost all prior work on the topic has focused on two settings: i) they either assume the target function $\theta_0(x)$ takes a parametric form and allow for a potentially high-dimensional parametric nuisance function, e.g. (Chernozhukov et al., 2016; 2017; 2018), ii) or take a non-parametric stance at estimating $\theta_0(x)$ but do not allow for high-dimensional nuisance functions (Wager & Athey, 2015; Athey et al., 2017).

We propose *Orthogonal Random Forest* (ORF), a random forest-based estimation algorithm, which performs non-parametric estimation of the target parameter while permitting more complex nuisance functions with high-dimensional parameterizations. Our estimator is also asymptotically normal and hence allows for the construction of asymptotically valid confidence intervals via plug-in or bootstrap approaches. Our approach combines the notion of *Neyman orthogonality* of the moment equations with a two-stage random forest based algorithm, which generalizes prior work on *Generalized Random Forests* (Athey et al., 2017) and the double machine learning (double ML) approach proposed in (Chernozhukov et al., 2017). To support our general algorithm, we also provide a novel nuisance estimation algorithm—*Forest Lasso*—that effectively recovers high-dimensional nuisance parameters provided they have locally sparse structure. This result combines techniques from Lasso theory (Hastie et al., 2015) with concentration inequalities for $U$-statistics (Hoeffding, 1963).

As a concrete example and as a main application of our approach, we consider the problem of *heterogeneous treat-*

---

[1]Microsoft Research–New England [2]University of Minnesota–Twin Cities. Correspondence to: Miruna Oprescu <moprescu@microsoft.com >, Vasilis Syrgkanis <vasy@microsoft.com>, Zhiwei Steven Wu <zsw@umn.edu>.

---

[1]See e.g. Reiss & Wolak (2007) and examples in Chernozhukov et al. (2016; 2018)

*ment effect* estimation. This problem is at the heart of many decision-making processes, including clinical trial assignment to patients, price adjustments of products, and ad placement by a search engine. In many situations, we would like to take the heterogeneity of the population into account and estimate the *heterogeneous treatment effect* (*HTE*)—the effect of a treatment $T$ (e.g. drug treatment, price discount, and ad position), on the outcome $Y$ of interest (e.g. clinical response, demand, and click-through-rate), as a function of observable characteristics $x$ of the treated subject (e.g. individual patient, product, and ad). HTE estimation is a fundamental problem in causal inference from observational data (Imbens & Rubin, 2015; Wager & Athey, 2015; Athey et al., 2017), and is intimately related to many areas of machine learning, including contextual bandits, off-policy evaluation and optimization (Swaminathan et al., 2016; Wang et al., 2017; Nie & Wager, 2017), and counterfactual prediction (Swaminathan & Joachims, 2015; Hartford et al., 2016).

The key challenge in HTE estimation is that the observations are typically collected by a policy that depends on confounders or control variables $W$, which also directly influence the outcome. Performing a direct regression of the outcome $Y$ on the treatment $T$ and features $x$, without controlling for a multitude of other potential confounders, will produce biased estimation. This leads to a regression problem that in the language of conditional moments takes the form:

$$\mathbb{E}\left[Y - \theta_0(x)T - f_0(x, W) \mid X = x\right] = 0 \qquad (2)$$

where $\theta_0(x)$ is the heterogeneous effect of the treatment $T$ (discrete or continuous) on the outcome $Y$ as a function of the features $x$ and $f_0(x, W)$ is an unknown nuisance function that captures the direct effect of the control variables on the outcome. Moreover, unlike active experimentation settings such as contextual bandits, when dealing with observational data, the actual treatment or logging policy $\mathbb{E}[T|x, W] = g_0(x, W)$ that could potentially be used to de-bias the estimation of $\theta_0(x)$ is also unknown.

There is a surge of recent work at the interplay of machine learning and causal inference that studies efficient estimation and inference of treatment effects. Chernozhukov et al. (2017) propose a two-stage estimation method called *double machine learning* that first orthogonalizes out the effect of high-dimensional confounding factors using sophisticated machine learning algorithms, including Lasso, deep neural nets and random forests, and then estimates the effect of the lower dimensional treatment variables, by running a low-dimensional linear regression between the residualized treatments and residualized outcomes. They show that even if the estimation error of the first stage is not particularly accurate, the second-stage estimate can still be $n^{-1/2}$-asymptotically normal. However, their approach requires a parametric specification of $\theta_0(x)$. In contrast, another line

of work that brings machine learning to causal inference provides fully flexible non-parametric HTE estimation based on random forest techniques (Wager & Athey, 2015; Athey et al., 2017; Powers et al., 2017). However, these methods heavily rely on low-dimensional assumptions.

Our algorithm ORF, when applied to the HTE problem (see Section 6) allows for the non-parametric estimation of $\theta_0(x)$ via forest based approaches while simultaneously allowing for a high-dimensional set of control variables $W$. This estimation problem is of practical importance when a decision maker (DM) wants to optimize a policy that depends only on a small set of variables, e.g. due to data collection or regulatory constraints or due to interpretability of the resulting policy, while at the same time controlling for many potential confounders in the existing data that could lead to biased estimates. Such settings naturally arise in contextual pricing or personalized medicine. In such settings the DM is faced with the problem of estimating a conditional average treatment effect conditional on a small set of variables while controlling for a much larger set. Our estimator provably offers a significant statistical advantage for this task over prior approaches.

In the HTE setting, the ORF algorithm follows the residual-on-residual regression approach analyzed by (Chernozhukov et al., 2016) to formulate a locally Neyman orthogonal moment and then applies our orthogonal forest algorithm to this orthogonal moment. Notably, (Athey et al., 2017) also recommend such a residual on residual regression approach in their empirical evaluation, which they refer to as "local centering", albeit with no theoretical analysis. Our results provide a theoretical foundation of the local centering approach through the lens of Neyman orthogonality. Moreover, our theoretical results give rise to a slightly different overall estimation approach than the one in (Athey et al., 2017): namely we residualize locally around the target estimation point $x$, as opposed to performing an overall residualization step and then calling the Generalized Random Forest algorithm on the residuals. The latter stems from the fact that our results require that the nuisance estimator achieve a good estimation rate *only* around the target point $x$. Hence, residualizing locally seems more appropriate than running a global nuisance estimation, which would typically minimize a non-local mean squared error. Our experimental findings reinforce this intuition (see e.g. comparison between ORF and the GRF-Res benchmark). Another notable work that combines the residualization idea with flexible heterogeneous effect estimation is that of (Nie & Wager, 2017), who formulate the problem as an appropriate residual-based square loss minimization over an arbitrary hypothesis space for the heterogeneous effect function $\theta(x)$. Formally, they show robustness, with respect to nuisance estimation errors, of the mean squared error (MSE) of the resulting estimate in expectation over the distribution $X$ and for the case where
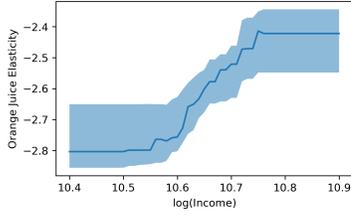
Figure 1: ORF estimates for the effect of orange juice price on demand from a high-dimensional dataset. We depict the estimated heterogeneity in elasticity by income level. The shaded region depicts the 1%-99% confidence interval obtained via bootstrap.

the hypothesis space is a reproducing kernel Hilbert space (RKHS). Our work differs primarily by: i) focusing on sup-norm estimation error at any target point $x$ as opposed to MSE, ii) using forest based estimation as opposed to finding a function in an RKHS, iii) working with the general orthogonal conditional moment problems, and iv) providing asymptotic normality results and hence valid inference.

We provide a comprehensive empirical comparison of ORF with several benchmarks, including three variants of GRF. We show that by setting the parameters according to what our theory suggests, ORF consistently outperforms all of the benchmarks. Moreover, we show that bootstrap based confidence intervals provide good finite sample coverage.

Finally, to motivate the usage of the ORF, we applied our technique to Dominick's dataset, a popular historical dataset of store-level orange juice prices and sales provided by University of Chicago Booth School of Business. The dataset is comprised of a large number of covariates $W$, but economics researchers might only be interested in learning the elasticity of demand as a function of a few variables $x$ such as income or education. We applied our method (see Appendix G for details) to estimate orange juice price elasticity as a function of income, and our results, depicted in Figure 1, unveil the natural phenomenon that lower income consumers are more price-sensitive.

## 2. Estimation via Local Orthogonal Moments

We study non-parametric estimation of models defined via conditional moment restrictions, in the presence of nuisance functions. Suppose we have a set of $2n$ observations $Z_1, \ldots, Z_{2n}$ drawn independently from some underlying distribution $\mathcal{D}$ over the observation domain $\mathcal{Z}$. Each observation $Z_i$ contains a feature vector $X_i \in \mathcal{X} := [0,1]^d$.

Given a target feature $x \in \mathcal{X}$, our goal is to estimate a parameter vector $\theta_0(x) \in \mathbb{R}^p$ that is defined via a local moment condition, i.e. for all $x \in \mathcal{X}$, $\theta_0(x)$ is the unique solution with respect to $\theta$ of:

$$\mathbb{E}\left[\psi(Z; \theta, h_0(x, W)) \mid X = x\right] = 0, \qquad (3)$$

where $\psi \colon \mathcal{Z} \times \mathbb{R}^p \times \mathbb{R}^\ell \to \mathbb{R}^p$ is a score function that maps an observation $Z$, parameter vector $\theta(x) \in \Theta \subset \mathbb{R}^p$, and nuisance vector $h(x,w)$ to a vector-valued score $\psi(z; \theta(x), h(x,w))$ and $h_0 \in H \subseteq \left(\mathbb{R}^d \times \mathbb{R}^L \to \mathbb{R}^\ell\right)$ is an unknown nuisance function that takes as input $X$ and a subvector $W$ of $Z$, and outputs a nuisance vector in $\mathbb{R}^\ell$. For any feature $x \in \mathcal{X}$, parameter $\theta \in \Theta$, and nuisance function $h \in H$, we define the *moment* function as:

$$m(x; \theta, h) = \mathbb{E}\left[\psi(Z; \theta, h(X, W)) \mid X = x\right] \qquad (4)$$

We assume that the dimensions $p, \ell, d$ are constants, while the dimension $L$ of $W$ can be growing with $n$.

We will analyze the following two-stage estimation process.

1. *First stage.* Compute a nuisance estimate $\hat{h}$ for $h_0$ using data $\{Z_{n+1}, \ldots, Z_{2n}\}$ with some guarantee on the conditional root mean squared error:[2]

$$\mathcal{E}(\hat{h}) = \sqrt{\mathbb{E}\left[\|\hat{h}(x, W) - h_0(x, W)\|^2 \mid X = x\right]}$$

2. *Second stage.* Compute a set of similarity weights $\{a_i\}$ over the data $\{Z_1, \ldots, Z_n\}$ that measure the similarity between their feature vectors $X_i$ and the target $x$. Compute the estimate $\hat{\theta}(x)$ using the nuisance estimate $\hat{h}$ via the plug-in weighted moment condition:

$$\hat{\theta}(x) \text{ solves: } \sum_{i=1}^n a_i \psi(Z_i; \theta, \hat{h}(X_i, W_i)) = 0 \quad (5)$$

In practice, our framework permits the use of any method to estimate the nuisance function in the first stage. However, since our description is a bit too abstract let us give a special case, which we will also need to assume for our normality result. Consider the case when the nuisance function $h$ takes the form $h(x, w) = g(w; \nu(x))$, for some known function $g$ but unknown function $\nu \colon \mathcal{X} \to \mathbb{R}^{d_\nu}$ (with $d_\nu$ potentially growing with $n$), i.e. locally around each $x$ the function $h$ is a parametric function of $w$. Moreover, the parameter $\nu_0(x)$ of the true nuisance function $h_0$ is identified as the minimizer of a local loss:

$$\nu_0(x) = \operatorname{argmin}_{\nu \in \mathcal{V}} \mathbb{E}\left[\ell(Z; \nu) \mid X = x\right] \qquad (6)$$

Then we can estimate $\nu_0(x)$ via a locally weighted and penalized empirical loss minimization algorithm. In particular in Section 5 we will consider the case of local $\ell_1$-penalized estimation that we will refer to as *forest lasso* and which provides formal guarantees in the case where $\nu_0(x)$ is sparse.

The key technical condition that allows us to reliably perform the two-stage estimation is the following *local orthogonality* condition, which can be viewed as a localized version of the *Neyman orthogonality* condtion (Neyman,

---

[2]Throughout the paper we denote with $\|\cdot\|$ the euclidean norm and with $\|\cdot\|_p$ the $p$-norm.

1979; Chernozhukov et al., 2017) around the neighborhood of the target feature $x$. Intuitively, the condition says that the score function $\psi$ is insensitive to local perturbations in the nuisance parameters around their true values.

**Definition 2.1** (Local Orthogonality). Fix any estimator $\hat{h}$ for the nuisance function. Then the Gateaux derivative with respect to $h$, denoted $D_\psi[\hat{h} - h_0 \mid x]$, is defined as:

$$\mathbb{E}\left[\nabla_h \psi(Z, \theta_0(x), h_0(x, W))(\hat{h}(x, W) - h_0(x, W)) \mid x\right]$$

where $\nabla_h$ denotes the gradient of $\psi$ with respect to the final $\ell$ arguments. We say that the moment conditions are *locally orthogonal* if for all $x$: $D_\psi[\hat{h} - h_0 \mid x] = 0$.

# 3. Orthogonal Random Forest

We describe our main algorithm *orthogonal random forest* (ORF) for calculating the similarity weights in the second stage of the two stage estimation. In the next section we will see that we will be using this algorithm for the estimation of the nuisance functions, so as to perform a local nuisance estimation. At a high level, ORF can be viewed as an orthogonalized version of GRF that is more robust to the nuisance estimation error. Similar to GRF, the algorithm runs a tree learner over $B$ random *subsamples $S_b$ (without replacement) of size $s < n$*, to build $B$ trees such that each tree indexed by $b$ provides a tree-based weight $a_{ib}$ for each observation $Z_i$ in the input sample. Then the ORF weight $a_i$ for each sample $i$ is the average over the tree-weights $a_{ib}$.

The tree learner starts with a root node that contains the entire $\mathcal{X}$ and recursively grows the tree to split $\mathcal{X}$ into a set of leaves until the number of observations in each leaf is not too small. The set of neighboods defined by the leaves naturally gives a simlarity measure between each observation and the target $x$. Following the same approach of (Tibshirani et al., 2018; Wager & Athey, 2015), we maintain the following tree properties in the process of building a tree.

**Specification 1** (Forest Regularity). *The tree satisfies*

- **Honesty**: *we randomly partition the input sample $S$ into two subsets $S^1$, $S^2$, then uses $S^1$ to place splits in the tree, and uses $S^2$ for estimation.*
- $\rho$**-balanced**: *each split leaves at least a fraction $\rho$ of the observations in $S^2$ on each side of the split for some parameter of $\rho \leq 0.2$.*
- **Minimum leaf size** $r$: *there are between $r$ and $2r - 1$ observations from $S^2$ in each leaf of the tree.*
- $\pi$**-random-split**: *at every step, marginalizing over the internal randomness of the learner, the probability that the next split occurs along the $j$-th feature is at least $\pi/d$ for some $0 < \pi \leq 1$, for all $j = 1, ..., d$.*[3]

---

[3]For example, this can be achieved by uniformly randomizing the splitting variable with probability $\pi$ or via a Poisson sampling

The key modification to GRF's tree learner is our incorporation of orthogonal nuisance estimation in the splitting criterion. While the splitting criterion does not factor into our theoretical analysis (similar to (Tibshirani et al., 2018)), we find it to be an effective practical heuristic.

**Splitting criterion with orthogonalization.** At each internal node $P$ we perform a two-stage estimation over $(P \cap S^1)$, i.e. the set of examples in $S^1$ that reach node $P$: 1) compute a nuisance estimate $\hat{h}_P$ using only data $P \cap S^1$ (e.g. by estimating a parameter $\hat{\nu}_P$ that minimizes $\sum_{i \in (P \cap S^1)} \ell(Z_i; \nu) + \lambda \|\nu\|_1$ and setting $\hat{h}_P(\cdot) = g(\cdot; \hat{\nu}_P)$), and then 2) form estimate $\hat{\theta}_P$ using $\hat{h}_P$:[4]

$$\hat{\theta}_P \in \operatorname{argmin}_{\theta \in \Theta} \left\| \sum_{i \in (P \cap S^1)} \psi(Z_i; \theta, \hat{h}_P(W_i)) \right\|$$

We now generate a large random set of candidate axis-aligned splits (satisfying Specification 1 and we want to find the split into two children $C_1$ and $C_2$ such that if we perform the same two-stage estimation separately at each child, the new estimates $\hat{\theta}_{C_1}$ and $\hat{\theta}_{C_2}$ take on very different values, so that the heterogeneity of the two children nodes is maximized. Performing the two-stage estimation of $\hat{\theta}_{C_1}$ and $\hat{\theta}_{C_2}$ for all candidate splits is too computationally expensive. Instead, we will approximate these estimates by taking a Newton step from the parent node estimate $\hat{\theta}_P$: for any child node $C$ given by a candidate split, our proxy estimate is:

$$\tilde{\theta}_C = \hat{\theta}_P - \frac{1}{|C \cap S^1|} \sum_{i \in C_j \cap S^1} A_P^{-1} \psi(Z_i; \hat{\theta}_P, \hat{h}_P(X_i, W_i))$$

where $A_P = \frac{1}{|P \cap S^1|} \sum_{i \in P \cap S_b^1} \nabla_\theta \psi(Z_i; \hat{\theta}_P, \hat{h}_P(X_i, W_i))$. We select the candidate split that maximizes the following proxy heterogeneity score: for each coordinate $t \in [p]$ let

$$\tilde{\Delta}_t(C_1, C_2) = \sum_{j=1}^2 \frac{1}{|C_j \cap S^1|} \left( \sum_{i \in C_j \cap S^1} \rho_{t,i} \right)^2 \quad (7)$$

where $\rho_{t,i} = A_P^{-1} \psi_t(Z_i; \hat{\theta}_P, \hat{h}_P(X_i, W_i))$. We then create a single heterogeneity score per split as a convex combination that puts weight $\eta$ on the mean and $(1 - \eta)$ on the maximum score across coordinates. $\eta$ is chosen uniformly at random in $[0, 1]$ at each iteration of splitting. Hence, some splits focus on heterogeneity on average, while others focus on creating heterogeneity on individual coordinates.

**ORF weights and estimator.** For each tree indexed $b \in [B]$ based on subsample $S_b$, let $L_b(x) \subseteq \mathcal{X}$ be the leaf that contains the target feature $x$. We assign *tree weight* and *ORF weight* to each observation $i$:

$$a_{ib} = \frac{\mathbf{1}[(X_i \in L_b(x)) \wedge (Z_i \in S_b^2)]}{|L_b(x) \cap S_b^2|}, \quad a_i = \frac{1}{B} \sum_{b=1}^B a_{ib}$$

---

scheme where a random subset of the variables of size $m$ is chosen to consider for candidate splits, with $m \sim \text{Poisson}(\lambda)$.

[4]In our implementation we actually use a cross-fitting approach, where we use half of $P \cap S^1$ to compute a nuisance function to apply to the other half and vice versa.

Wager & Athey (2015) show that under the structural specification of the trees, the tree weights are non-zero only around a small neighborhood of $x$; a property that we will leverage in our analysis.

**Theorem 3.1** (Kernel shrinkage (Wager & Athey, 2015)). *Suppose the minimum leaf size parameter $r = O(1)$, the tree is $\rho$-balanced and $\pi$-random-split and the distribution of $X$ admits a density in $[0,1]^d$ that is bounded away from zero and infinity. Then the tree weights satisfy*

$$\mathbb{E}\left[\sup\{\|x - x_i\| : a_{ib} > 0\}\right] = O(s^{-\frac{1}{2\alpha d}}),$$

*where*

$$\alpha = \frac{\log(\rho^{-1})}{\pi \log((1-\rho)^{-1})}$$

*and $s$ is size of the subsamples.*

## 4. Convergence and Asymptotic Analysis

The *ORF estimate* $\hat{\theta}$ is computed by solving the weighted moment condition in Equation (5), using the *ORF weights* as described in the previous section. We now provide theoretical guarantees for $\hat{\theta}$ under the following assumption on the moment, score fuction and the data generating process.

**Assumption 4.1.** *The moment condition and the score function satisfy the following:*

1. ***Local Orthogonality.*** *The moment condition satisfies local orthogonality.*
2. ***Identifiability.*** *The moments $m(x; \theta, h_0) = 0$ has a unique solution $\theta_0(x)$.*
3. ***Smooth Signal.*** *The moments $m(x; \theta, h)$ are $O(1)$-Lipschitz in $x$ for any $\theta \in \Theta, h \in H$.*
4. ***Curvature.*** *The Jacobian $\nabla_\theta m(x; \theta_0(x), h_0)$ has minimum eigenvalue bounded away from zero.*
5. ***Smoothness of scores.*** *For every $j \in [p]$ and for all $\theta$ and $h$, the eigenvalues of the expected Hessian $\mathbb{E}\left[\nabla^2_{(\theta,h)} \psi_j(Z; \theta, h(W)) \mid x, W\right]$ are bounded above by a constant $O(1)$. For any $Z$, the score $\psi(Z; \theta, \xi)$ is $O(1)$-Lipschitz in $\theta$ for any $\xi$ and $O(1)$-Lipschitz in $\xi$ for any $\theta$. The gradient of the score with respect to $\theta$ is $O(1)$-Lipschitz in $\xi$.*
6. ***Boundedness.*** *The parameter set $\Theta$ has constant diameter. There exists a bound $\psi_{\max}$ such that for any observation $Z$, the first-stage nuisance estimate $\hat{h}$ satisfies $\|\psi(Z; \theta, \hat{h})\|_\infty \leq \psi_{\max}$ for any $\theta \in \Theta$.*
7. ***Full Support $X$.*** *The distribution of $X$ admits a density that is bounded away from zero and infinity.*

All the results presented in the remainder of the paper will assume these conditions and we omit stating so in each of the theorems. Any extra conditions required for each theorem will be explicitly provided. Note that except for the local orthogonality condition, all of the assumptions are imposing standard boundedness and regularity conditions of the moments.

**Theorem 4.2** ($L^q$-Error Bound). *Suppose that:* $\mathbb{E}\left[\mathcal{E}(\hat{h})^{2q}\right]^{1/2q} \leq \chi_{n,2q}$. *Then:*

$$\mathbb{E}\left[\|\hat{\theta} - \theta_0\|^q\right]^{1/q} = O\left(\frac{1}{s^{\frac{1}{2\alpha d}}} + \sqrt{\frac{s \log(\frac{n}{s})}{n}} + \chi^2_{n,2q}\right)$$

**Theorem 4.3** (High Probability Error Bound). *Suppose that the score is the gradient of a convex loss and let $\sigma > 0$ denote the minimum eigenvalue of the jacobian $M$. Moreover, suppose that the nuisance estimate satisfies that w.p. $1 - \delta$: $\mathcal{E}(\hat{h}) \leq \chi_{n,\delta}$. Then w.p. $1 - 2\delta$:*

$$\|\hat{\theta} - \theta_0\| = \frac{O\left(s^{-\frac{1}{2\alpha d}} + \sqrt{\frac{s \log(\frac{n}{s\delta})}{n}} + \chi^2_{n,\delta}\right)}{\sigma - O(\chi_{n,\delta})} \tag{8}$$

For asymptotic normality we will restrict our framework to the case of parametric nuisance functions, i.e. $h(X, W) = g(W; \nu(X))$ for some known function $g$ and to a particular type of nuisance estimators that recover the true parameter $\nu_0(x)$. Albeit we note that the parameter $\nu(X)$ can be an arbitrary non-parametric function of $X$ and can also be high-dimensional. We will further assume that the moments also have a smooth co-variance structure in $X$, i.e. if we let

$$V = \psi(Z; \theta_0(x), g(W; \nu_0(x)))$$

then $\text{Var}(V \mid X = x')$ is Lipschitz in $x'$ for any $x' \in [0,1]^d$.

**Theorem 4.4** (Asymptotic Normality). *Suppose that $h_0(X, W)$ takes a locally parametric form $g(W; \nu_0(X))$, for some known function $g(\cdot; \nu)$ that is $O(1)$-Lipschitz in $\nu$ w.r.t. the $\ell_r$ norm for some $r \geq 1$ and the nuisance estimate is of the form $\hat{h}(X, W) = g(W; \hat{\nu}(x))$ and satisfies:*

$$\mathbb{E}\left[\|\hat{\nu}(x) - \nu_0(x)\|^4_r\right]^{1/4} \leq \chi_{n,4} = o\left((s/n)^{1/4}\right)$$

*Suppose that $s$ is chosen such that:*

$$s^{-1/(2\alpha d)} = o((s/n)^{1/2-\varepsilon}),$$

*for any $\varepsilon > 0$, and $s = o(n)$. Moreover, $\text{Var}(V \mid X = x')$ is Lipschitz in $x'$ for any $x' \in [0,1]^d$. Then for any coefficient $\beta \in \mathbb{R}^p$, with $\|\beta\| \leq 1$, assuming $\text{Var}(\beta^\intercal M^{-1} V | X = x') > 0$ for any $x' \in [0,1]^d$, there exists a sequence $\sigma_n = \Theta(\sqrt{\text{polylog}(n/s)s/n})$, such that:*

$$\sigma_n^{-1} \left\langle \beta, \hat{\theta} - \theta_0 \right\rangle \to_d \mathcal{N}(0, 1) \tag{9}$$

Given the result in Theorem 4.4, we can follow the same approach of *Bootstrap of Little Bags* by (Athey et al., 2017; Sexton & Laake, 2009) to build valid confidence intervals.

## 5. Nuisance Estimation: Forest Lasso

Next, we study the nuisance estimation problem in the first stage and provide a general nuisance estimation method that leverages locally sparse parameterization of the nuisance function, permitting low error rates even for high-dimensional problems. Consider the case when the nuisance function $h$ takes the form $h(x, w) = g(w; \nu(x))$ for some known functional form $g$, for some known function $g$ but unknown function $\nu : \mathcal{X} \to \mathbb{R}^{d_\nu}$, with $d_\nu$ potentially growing with $n$. Moreover, the parameter $\nu_0(x)$ of the true nuisance function $h_0$ is identified as the minimizer of some local loss, as defined in Equation (6).

We consider the following estimation process: given a set of observations $D_1$, we run the same tree learner in Section 3 over $B$ random subsamples (without replacement) to compute ORF weights $a_i$ for each observation $i$ over $D_1$. Then we apply a local $\ell_1$ penalized $M$-estimation:

$$\hat{\nu}(x) = \arg\min_{\nu \in \mathcal{V}} \sum_{i=1}^{n} a_i \, \ell(Z_i; \nu) + \lambda \|\nu\|_1 \qquad (10)$$

To provide formal guarantees for this method we will need to make the following assumptions.

**Assumption 5.1** (Assumptions for nuisance estimation). *The target parameter and data distribution satisfy:*

- *For any $x \in \mathcal{X}$, $\nu(x)$ is $k$-sparse with support $S(x)$.*
- *$\nu(x)$ is a $O(1)$-Lipschitz in $x$ and the function $\nabla_\nu L(x; \nu) = \mathbb{E}[\nabla_\nu \ell(Z; \nu) \mid X = x]$ is $O(1)$-Lipschitz in $x$ for any $\nu$, with respect to the $\ell_2$ norm.*
- *The data distribution satisfies the conditional restricted eigenvalue condition: for all $\nu \in \mathcal{V}$ and for all $z \in \mathcal{Z}$, for some matrix $\mathcal{H}(z)$ that depends only on the data: $\nabla_{\nu\nu}\ell(z; \nu) \succeq \mathcal{H}(z) \succeq 0$, and for all $x$ and for all $\nu \in C(S(x); 3) \equiv \{\nu \in \mathbb{R}^d : \|\nu_{S(x)^c}\|_1 \le 3\|\nu_{S(x)}\|_1\}$:*

$$\nu^T \mathbb{E}[\mathcal{H}(Z) \mid X = x] \, \nu \ge \gamma \|\nu\|_2^2 \qquad (11)$$

Under Assumption 5.1 we show that the local penalized estimator achieves the following parameter recovery guarantee.

**Theorem 5.2.** *With probability $1 - \delta$:*

$$\|\hat{\nu}(x) - \nu_0(x)\|_1 \le \frac{2\lambda k}{\gamma - 32k\sqrt{s \ln(d_\nu/\delta)/n}}$$

*as long as $\lambda \ge \Theta\left(s^{-1/(2\alpha d)} + \sqrt{\frac{s \ln(d_\nu/\delta)}{n}}\right)$.*

**Example 5.3** (Forest Lasso). *For locally sparse linear regression, $Z_i = (x_i, y_i, W_i)$ and $\ell(Z_i; \nu) = (y_i - \langle \nu, W_i \rangle)^2$. This means, $\nabla_{\nu\nu}\ell(Z_i; \nu) = W_i W_i^T = \mathcal{H}(Z_i)$. Hence, the conditional restricted eigenvalue condition is simply a conditional covariance condition: $\mathbb{E}[WW^\intercal \mid x] \succeq \gamma I$.*

**Example 5.4** (Forest Logistic Lasso). *For locally sparse logistic regression, $Z_i = (x_i, y_i, W_i)$, $y_i \in \{0, 1\}$ and*

$$\ell(Z_i; \nu) = y_i \ln\left(\mathcal{L}(\langle \nu, W_i \rangle)\right) + (1 - y_i) \ln\left(1 - \mathcal{L}(\langle \nu, W_i \rangle)\right),$$

*where $\mathcal{L}(t) = 1/(1 + e^{-t})$ is the logistic function. In this case,*

$$\nabla_{\nu\nu}\ell(Z_i; \nu) = \mathcal{L}(\langle \nu, W_i \rangle)(1 - \mathcal{L}(\langle \nu, W_i \rangle))W_i W_i^\intercal$$
$$\succeq \rho W_i W_i^\intercal = \mathcal{H}(Z_i)$$

*assuming the index $\langle \nu, w \rangle$ is bounded in some finite range. Hence, our conditional restricted eigenvalue condition is the same conditional covariance condition:*

$$\rho \, \mathbb{E}[WW^T \mid x] \succeq \rho\gamma I$$

## 6. Heterogeneous Treatment Effects

Now we apply ORF to the problem of estimating *heterogeneous treatment effects*. We will consider the following *extension* of the *partially linear regression (PLR)* model due to Robinson (1988). [5] We have $2n$ i.i.d. observations $D = \{Z_i = (T_i, Y_i, W_i, X_i)\}_{i=1}^{2n}$ such that for each $i$, $T_i$ represents the treatment applied that can be either real-valued (in $\mathbb{R}^p$) or discrete (taking values in $\{0, e_1, \ldots, e_p\}$, where each $e_j$ denotes the standard basis in $\mathbb{R}^p$), $Y_i \in \mathbb{R}$ represents the outcome, $W_i \in [-1, 1]^{d_\nu}$ represents potential confounding variables (controls), and $X_i \in \mathcal{X} = [0, 1]^d$ is the feature vector that captures the heterogeneity. The set of parameters are related via the following equations:

$$Y = \langle \mu_0(X, W), T \rangle + f_0(X, W) + \varepsilon, \qquad (12)$$
$$T = g_0(X, W) + \eta, \qquad (13)$$

where $\eta, \varepsilon$ are bounded unobserved noises such that $\mathbb{E}[\varepsilon \mid W, X, T] = 0$ and $\mathbb{E}[\eta \mid X, W, \varepsilon] = 0$. In the main equation (12), $\mu_0 : \mathbb{R}^d \times \mathbb{R}^{d_\nu} \to [-1, 1]^p$ represents the treatment effect function. Our goal is to estimate *conditional average treatment effect* (CATE) $\theta_0(x)$ conditioned on target feature $x$:

$$\theta_0(x) = \mathbb{E}[\mu_0(X, W) \mid X = x]. \qquad (14)$$

The confounding equation (13) determines the relationship between treatments variable $T$ and the feature $X$ and confounder $W$. To create an orthogonal moment for identifying $\theta_0(x)$, we follow the classical *residualization* approach similar to (Chernozhukov et al., 2017). First, observe that

$$Y - \mathbb{E}[Y \mid X, W] = \Big\langle \mu_0(X, W), T - \mathbb{E}[T \mid X, W] \Big\rangle + \varepsilon$$

Let us define the function $q_0(X, W) = \mathbb{E}[Y \mid X, W]$, and consider the residuals

$$\tilde{Y} = Y - q_0(X, W)$$
$$\tilde{T} = T - g_0(X, W) = \eta$$

---

Then we can simplify the equation as $\tilde{Y} = \mu_0(X, W) \cdot \tilde{T} + \varepsilon$. As long as $\eta$ is independent of $\mu_0(X, W)$ conditioned on $X$ (e.g. $\eta$ is independent of $W$ or $\mu_0(X, W)$ does not depend on $W$), we also have

$$\mathbb{E}\left[\mu_0(X, W) \mid X, \eta\right] = \mathbb{E}\left[\mu_0(X, W) \mid X\right] = \theta(X).$$

Since $\mathbb{E}\left[\varepsilon \mid X, \eta\right] = \mathbb{E}\left[\mathbb{E}\left[\varepsilon \mid X, W, T\right] \mid X, \eta\right] = 0$, then

$$\mathbb{E}\left[\tilde{Y} \mid X, \tilde{T}\right] = \mathbb{E}\left[\mu_0(X, W) \mid X\right] \cdot \tilde{T} = \theta(X) \cdot \tilde{T}.$$

This relationship suggests that we can obtain an estimate of $\theta(x)$ by regressing $\tilde{Y}$ on $\tilde{T}$ locally around $X = x$. We can thus define the *orthogonalized* score function: for any observation $Z = (T, Y, W, x)$, any parameter $\theta \in \mathbb{R}^p$, any estimates $q$ and $g$ for functions $q_0$ and $g_0$, the score $\psi(Z; \theta, h(X, W))$ is:

$$\{Y - q(X, W) - \theta\left(T - g(X, W)\right)\}\left(T - g(X, W)\right),$$

where $h(X, W) = (q(X, W), g(X, W))$. In the appendix, we show that this moment condition satisfies local orthogonality, and it identifies $\theta_0(x)$ as long as as the noise $\eta$ is independent of $\mu_0(X, W)$ conditioned on $X$ and the expected matrix $\mathbb{E}\left[\eta\eta^\mathsf{T} \mid X = x\right]$ is invertible. Even though the approach applies generically, to obtain formal guarantees on the nuisance estimates via our Forest Lasso, we will restrict their functional form.

**Real-valued treatments.** Suppose $f_0$ and each coordinate $j$ of $g_0$ and $\mu_0$ are given by high-dimensional linear functions: $f_0(X, W) = \langle W, \beta_0(X)\rangle$, $\mu_0^j(X, W) = \langle W, u_0^j(X)\rangle$, $g_0^j(X, W) = \langle W, \gamma_0^j(X)\rangle$, where $\beta_0(X), \gamma_0^j(X), u_0^j(X)$ are $k$-sparse vectors in $\mathbb{R}^{d_\nu}$. Consequently, $q_0(X, W)$ can be written as a $k^2$-sparse linear function over degree-2 polynomial features $\phi_2(W)$ of $W$. Then as long as $\gamma_0, \beta_0$ and $\mu_0$ are Lipschitz in $X$ and the confounders $W$ satisfy

$$\mathbb{E}\left[\phi_2(W)\phi_2(W)^\mathsf{T} \mid X\right] \succeq \Omega(1)I,$$

then we can use Forest Lasso to estimate both $g_0(x, w)$ and $q_0(x, w)$. Hence, we can apply the ORF algorithm to get estimation error rates and asymptotic normality results for $\hat{\theta}$. (see Appendix B for formal statement).

**Discrete treatments.** We now describe how our theory can be applied to discrete treatments. Suppose $f_0$ and each coordinate $j$ of $g_0$ are of the form: $f_0(X, W) = \langle W, \beta_0(X)\rangle$ and $g_0^j(X, W) = \mathcal{L}(\langle W, \gamma_0^j(X)\rangle)$, where $\mathcal{L}(t) = 1/(1 + e^{-t})$ is the logistic function. Note in this case $\eta$ is not independent of $W$ since

$$\mathrm{Var}(\eta_j) = g_0^j(X, W)(1 - g_0^j(X, W)).$$

To maintain the conditional independence between $\mu_0(X, W)$ and $\eta$ conditioned on $X$, we focus on the setting where $\mu_0$ is only a function of $X$, i.e. $\mu(X, W) = \theta(X)$ for all $W, X$. In this setting we can estimate $g_0$ by running a forest logistic lasso for each treatment $j$. Then we can estimate $q_0(x, W)$ as follows: For each $t \in \{e_1, \ldots, e_p\}$ estimate the expected counter-factual outcome function: $m_0^t(x, W) = \mu_0^t(x, W) + f_0(x, W)$, by running a forest lasso between $Y$ and $X, W$ only among the subset of samples that received treatment $t$. Similarly, estimate $f_0(x, W)$ by running a forest lasso between $Y$ and $X, W$ only among the subset of samples that received treatment $t = 0$. Then observe that $q_0(x, W)$ can be written as a function of $f_0$, $g_0^t$ and $m_0^t$. Thus we can combine these estimates to get an estimate of $q_0$. Hence, we can obtain a guarantee similar to that of Corollary B.1 (see appendix).

**Doubly robust moment for discrete treatments.** In the setting where $\mu$ also depends on $W$ and treatments are discrete, we can formulate an alternative orthogonal moment that identifies the CATE even when $\eta$ is correlated with $\mu(X, W)$. This moment is based on first constructing unbiased estimates of the counterfactual outcome $m_0^t(X, W) = \mu_0^t(X, W) + f_0(X, W)$ for every observation $X, W$ and for any potential treatment $t$, i.e. even for $t \neq T$. The latter is done by invoking the doubly robust formula (Robins & Rotnitzky, 1995; Cassel et al., 1976; Kang et al., 2007):

$$Y^{(t)} = m_0^t(X, W) + \frac{(Y - m_0^t(X, W))\mathbf{1}\{T = t\}}{g_0^t(X, W)}$$

with the convention that

$$g_0^0(X, W) = 1 - \sum_{t \neq 0} g_0^t(X, W)$$
$$m_0^0(X, W) = f_0(X, W)$$

Then we can identify the parameter $\theta^t(x)$ using the moment: $\mathbb{E}[Y^{(t)} - Y^{(0)} \mid X = x] = \theta_t(x)$. One can easily show that this moment satisfies the Neyman orthogonality condition with respect to the nuisance functions $m$ and $g$ (see appendix). In fact this property is essentially implied by the fact that the estimates $Y^{(t)}$ satisfy the double robustness property, since double robustness is a stronger condition than orthogonality. We will again consider $\mu_0^j(X, W) = \langle W, u_0^j(X)\rangle$. Then using similar reasoning as in the previous paragraph, we see that with a combination of forest logistic lasso for $g_0^t$ and forest lasso for $m_0^t$, we can estimate these nuisance functions at a sufficiently fast rate for our ORF estimator (based on this doubly robust moment) to be asymptotically normal, assuming they have locally sparse linear or logistic parameterizations.
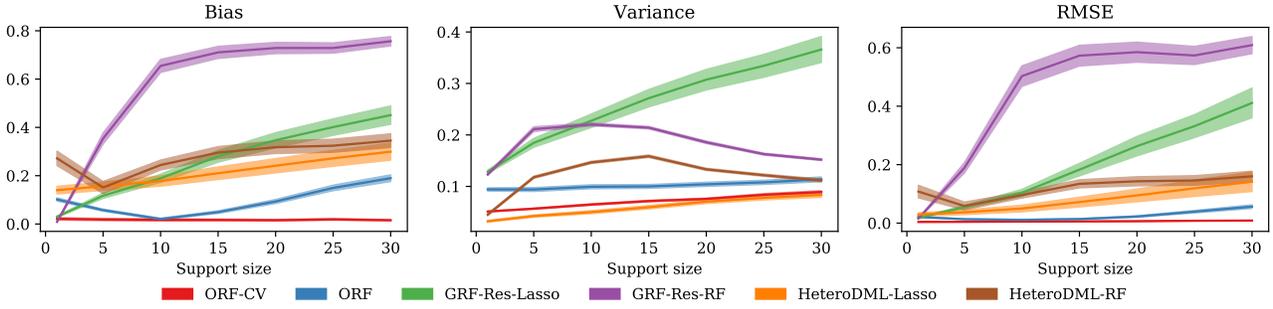
Figure 2: Bias, variance and RMSE as a function of support size for $n = 5000$, $p = 500$, $d = 1$ and a piecewise linear treatment response function. The solid lines represent the mean of the metrics over Monte Carlo experiments and test points, and the filled regions depict the standard deviation, scaled down by 3 for clarity.

## 7. Monte Carlo Experiments

We compare the empirical performance of ORF with other methods in the literature (and their variants).[6] The data generating process we consider is described by the following equations:

$$Y_i = \theta_0(x_i)\, T_i + \langle W_i, \gamma_0 \rangle + \varepsilon_i$$
$$T_i = \langle W_i, \beta_0 \rangle + \eta_i$$

Moreover, $x_i$ is drawn from the uniform distribution $U[0,1]$, $W_i$ is drawn from $\mathcal{N}(0, I_p)$, and the noise terms $\varepsilon_i \sim U[-1,1], \eta_i \sim U[-1,1]$. The $k$-sparse vectors $\beta_0, \gamma_0 \in \mathbb{R}^p$ have coefficients drawn independently from $U[0,1]$. The dimension $p = 500$ and we vary the support size $k$ over the range of $\{1, 5, 10, 15, 20, 25, 30\}$. We examine a treatment function $\theta(x)$ that is continuous and piecewise linear (detailed in Figure 3). In Appendix H we analyze other forms for $\theta(x)$.

For each fixed treatment function, we repeat 100 experiments, each of which consists of generating 5000 observations from the DGP, drawing the vectors $\beta_0$ and $\gamma_0$, and estimating $\hat{\theta}(x)$ at 100 test points $x$ over a grid in $[0, 1]$. We then calculate the bias, variance and root mean squared error (RMSE) of each estimate $\hat{\theta}(x)$. Here we report summary statistics of the median and $5 - 95$ percentiles of these three quantities across test points, so as to evaluate the average performance of each method. We compare two variants of ORF with two variants of GRF (Athey et al., 2017) (see Appendix H for a third variant) and two extensions of double ML methods for heterogeneous treatment effect estimation (Chernozhukov et al., 2017).

*ORF variants.* (1) ORF: We implement ORF as described in Section 3, setting parameters under the guidance of our theoretical result: subsample size $s \approx (n/\log(p))^{1/(2\tau+1)}$, Lasso regularization $\lambda_\gamma, \lambda_q \approx \sqrt{\log(p)s/n}/20$ (for both

tree learner and kernel estimation), number of trees $B = 100 \geq n/s$, a max tree depth of 20, and a minimum leaf size of $r = 5$. (2) ORF with LassoCV (ORF-CV): we replaced the Lasso algorithm in ORF's kernel estimation, with a cross-validated Lasso for the selection of the regularization parameter $\lambda_\gamma$ and $\lambda_q$. ORF-CV provides a more systematic optimization over the parameters.

*GRF variants.* (1) GRF-Res-Lasso: We perform a naive combination of double ML and GRF by first residualizing the treatments and outcomes on both the features $x$ and controls $W$, then running GRF R package by (Tibshirani et al., 2018) on the residualized treatments $\hat{T}$, residualized outcomes $\hat{Y}$, and features $x$. A cross-validated Lasso is used for residualization. (2) GRF-Res-RF: We combine DoubleML and GRF as above, but we use cross-validated Random Forests for calculating residuals $\hat{T}$ and $\hat{Y}$.

*Double ML with Polynomial Heterogeneity (DML-Poly).* An extension of the classic Double ML procedure for heterogeneous treatment effects introduced in (Chernozhukov et al., 2017). This method accounts for heterogeneity by creating an expanded linear base of composite treatments (cross products between treatments and features). (1) Heterogeneous Double ML using LassoCV for first-stage estimation (HeteroDML-Lasso): In this version, we use Lasso with cross-validation for calculating residuals on $x \cup W$ in the first stage. (2) Heterogeneous Double ML using random forest for first-stage estimation (HeteroDML-RF): A more flexible version that uses random forests to perform residualization on treatments and outcomes. The latter performs better when treatments and outcomes have a non-linear relationship with the joint features of $(x, W)$.

We generated data according to the Monte Carlo process above and set the parameters to $n = 5000$ samples, $p = 500$ controls, $d = 1$ features and support size $k \in \{1, 5, 10, 15, 20, 25, 30\}$ and three types of treatment effect functions. In this section, we present the results for a piecewise linear treatment effect function.

---

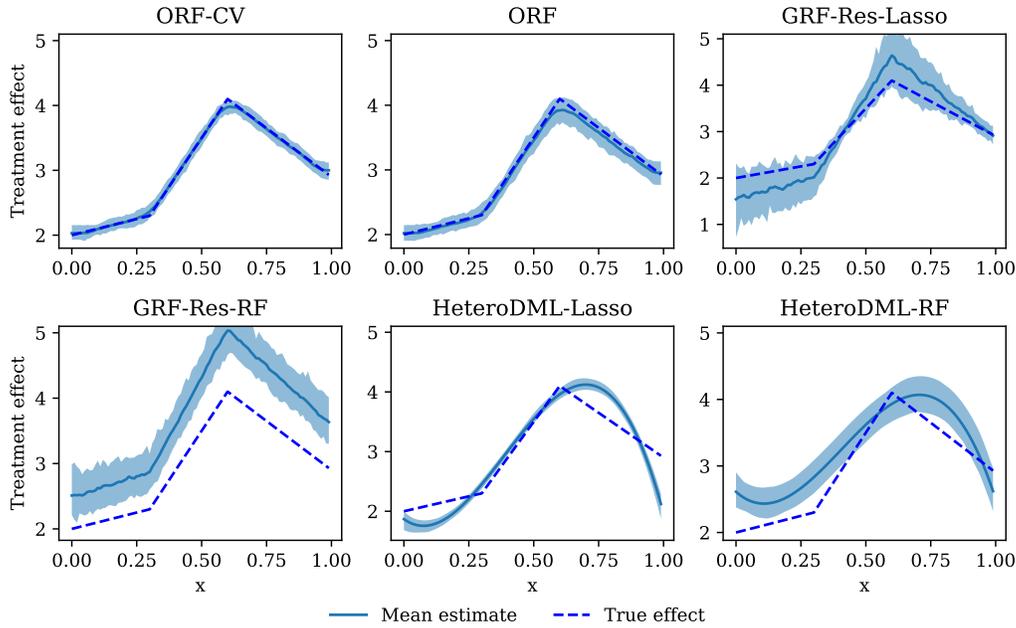[6]The source code for running these experiments is available in the git repo Microsoft/EconML.

Figure 3: Treatment effect estimations for 100 Monte Carlo experiments with parameters $n = 5000$, $p = 500$, $d = 1$, $k = 15$, and $\theta(x) = (x+2)\mathbb{I}_{x \leq 0.3} + (6x+0.5)\mathbb{I}_{x > 0.3 \text{ and } x \leq 0.6} + (-3x+5.9)\mathbb{I}_{x > 0.6}$. The shaded regions depict the mean and the 5%-95% interval of the 100 experiments.

In Figure 3, we inspect the goodness of fit for the chosen estimation methods across 100 Monte Carlo experiments. We note the limitations of two versions of the GRF-Res estimators, GRF-Res-Lasso and GRF-Res-RF, in capturing the treatment effect function well. The GRF-Res-RF estimations have a consistent bias as the Random Forest residualization cannot capture the dependency on the controls $W$ given their high-dimensionality. The HeteroDML methods are not flexible enough to capture the complexity of the treatment effect function. The best performers are the ORF-CV, ORF, and GRF-Res-Lasso, with the latter estimator having a larger bias and variance.

ferent support sizes. The ORF-CV performs very well, with consistent bias and RMSE across support sizes and treatment functions. The bias, variance and RMSE of the ORF grow with support size, but this growth is at a lower rate compared to the alternative estimators. The ORF-CV and ORF algorithms perform better than the GRF-Res methods on all metrics for this example. We observe this pattern for the other choices of support size, sample size and treatment effect function (see Appendix H). In Figure 4, we provide a snapshot of the bootstrap confidence interval coverage for this example.

## Acknowledgements

Figure 4: Sample 1%-99% confidence intervals for 1000 bootstrap iterations with parameters $n = 5000$, $p = 500$, $d = 1$, $k = 15$, and $\theta(x) = (x+2)\mathbb{I}_{x \leq 0.3} + (6x+0.5)\mathbb{I}_{x > 0.3 \text{ and } x \leq 0.6} + (-3x+5.9)\mathbb{I}_{x > 0.6}$. Approximately 90% of the sampled test points are contained in the interval.

We analyze these estimators as we increase the support size of $W$. Figures 2 illustrate the evaluation metrics across dif-

# References

Ai, C. and Chen, X. Efficient estimation of models with conditional moment restrictions containing unknown functions. *Econometrica*, 71 (6):1795–1843, 2003. doi: 10.1111/1468-0262. 00470. URL https://onlinelibrary.wiley. com/doi/abs/10.1111/1468-0262.00470.

Athey, S., Tibshirani, J., and Wager, S. Generalized Random Forests. *ArXiv e-prints*, October 2017.

Cassel, C. M., Särndal, C. E., and Wretman, J. H. Some results on generalized difference estimation and generalized regression estimation for finite populations. *Biometrika*, 63(3):615–620, 1976.

Chen, X. and Pouzo, D. Efficient estimation of semiparametric conditional moment models with possibly nonsmooth residuals. *Journal of Econometrics*, 152(1):46 – 60, 2009. ISSN 0304-4076. doi: https://doi.org/10.1016/j.jeconom.2009.02.002. URL http://www.sciencedirect.com/science/ article/pii/S0304407609000529. Recent Adavances in Nonparametric and Semiparametric Econometrics: A Volume Honouring Peter M. Robinson.

Chernozhukov, V., Newey, W. K., and Santos, A. Constrained Conditional Moment Restriction Models. *arXiv e-prints*, art. arXiv:1509.06311, September 2015.

Chernozhukov, V., Escanciano, J. C., Ichimura, H., Newey, W. K., and Robins, J. M. Locally Robust Semiparametric Estimation. *arXiv e-prints*, art. arXiv:1608.00033, July 2016.

Chernozhukov, V., Chetverikov, D., Demirer, M., Duflo, E., Hansen, C., and Newey, W. Double/debiased/neyman machine learning of treatment effects. *American Economic Review*, 107(5):261–65, May 2017. doi: 10.1257/ aer.p20171038. URL http://www.aeaweb.org/ articles?id=10.1257/aer.p20171038.

Chernozhukov, V., Goldman, M., Semenova, V., and Taddy, M. Orthogonal Machine Learning for Demand Estimation: High Dimensional Causal Inference in Dynamic Panels. *ArXiv e-prints*, December 2017.

Chernozhukov, V., Nekipelov, D., Semenova, V., and Syrgkanis, V. Plug-in Regularized Estimation of High-Dimensional Parameters in Nonlinear Semiparametric Models. *arXiv e-prints*, art. arXiv:1806.04823, June 2018.

Hartford, J., Lewis, G., Leyton-Brown, K., and Taddy, M. Counterfactual Prediction with Deep Instrumental Variables Networks. *ArXiv e-prints*, December 2016.

Hastie, T., Tibshirani, R., and Wainwright, M. *Statistical Learning with Sparsity: The Lasso and Generalizations*. Chapman & Hall/CRC, 2015. ISBN 1498712169, 9781498712163.

Hoeffding, W. Probability inequalities for sums of bounded random variables. *Journal of the American Statistical Association*, 58(301):13–30, 1963. ISSN 01621459. URL http://www.jstor.org/stable/2282952.

Imbens, G. W. and Rubin, D. B. *Causal Inference for Statistics, Social, and Biomedical Sciences: An Introduction*. Cambridge University Press, 2015. doi: 10.1017/CBO9781139025751.

Kang, J. D., Schafer, J. L., et al. Demystifying double robustness: A comparison of alternative strategies for estimating a population mean from incomplete data. *Statistical science*, 22(4):523–539, 2007.

Mentch, L. and Hooker, G. Quantifying uncertainty in random forests via confidence intervals and hypothesis tests. *The Journal of Machine Learning Research*, 17(1): 841–881, 2016.

Newey, W. K. 16 efficient estimation of models with conditional moment restrictions. In *Econometrics*, volume 11 of *Handbook of Statistics*, pp. 419 – 454. Elsevier, 1993. doi: https://doi.org/10.1016/S0169-7161(05)80051-3. URL http://www.sciencedirect.com/ science/article/pii/S0169716105800513.

Neyman, J. C($\alpha$) tests and their use. *Sankhy: The Indian Journal of Statistics, Series A (1961-2002)*, 41(1/2):1–21, 1979. ISSN 0581572X. URL http://www.jstor. org/stable/25050174.

Nie, X. and Wager, S. Learning Objectives for Treatment Effect Estimation. *ArXiv e-prints*, December 2017.

Peel, T., Anthoine, S., and Ralaivola, L. Empirical bernstein inequalities for u-statistics. In Lafferty, J. D., Williams, C. K. I., Shawe-Taylor, J., Zemel, R. S., and Culotta, A. (eds.), *Advances in Neural Information Processing Systems 23*, pp. 1903–1911. Curran Associates, Inc., 2010.

Powers, S., Qian, J., Jung, K., Schuler, A., Shah, N. H., Hastie, T., and Tibshirani, R. Some methods for heterogeneous treatment effect estimation in high-dimensions. *ArXiv e-prints*, July 2017.

Reiss, P. C. and Wolak, F. A. Chapter 64 structural econometric modeling: Rationales and examples from industrial organization. volume 6 of *Handbook of Econometrics*, pp. 4277 – 4415. Elsevier, 2007. doi: https://doi.org/10.1016/S1573-4412(07)06064-3. URL http://www.sciencedirect.com/science/ article/pii/S1573441207060643.

Robins, J. M. and Rotnitzky, A. Semiparametric efficiency in multivariate regression models with missing data. *Journal of the American Statistical Association*, 90(429):122–129, 1995.

Robinson, P. M. Root-n-consistent semiparametric regression. *Econometrica*, 56(4):931–954, 1988. ISSN 00129682, 14680262. URL http://www.jstor.org/stable/1912705.

Sexton, J. and Laake, P. Standard errors for bagged and random forest estimators. *Computational Statistics and Data Analysis*, 53(3):801 – 811, 2009. ISSN 0167-9473. doi: https://doi.org/10.1016/j.csda.2008.08.007. URL http://www.sciencedirect.com/science/article/pii/S0167947308003988. Computational Statistics within Clinical Research.

Swaminathan, A. and Joachims, T. Counterfactual risk minimization: Learning from logged bandit feedback. *CoRR*, abs/1502.02362, 2015. URL http://arxiv.org/abs/1502.02362.

Swaminathan, A., Krishnamurthy, A., Agarwal, A., Dudík, M., Langford, J., Jose, D., and Zitouni, I. Off-policy evaluation for slate recommendation. *CoRR*, abs/1605.04812, 2016. URL http://arxiv.org/abs/1605.04812.

Tibshirani, J., Athey, S., Wager, S., Friedberg, R., Miner, L., and Wright, M. *grf: Generalized Random Forests (Beta)*, 2018. URL https://CRAN.R-project.org/package=grf. R package version 0.10.2.

Wager, S. and Athey, S. Estimation and Inference of Heterogeneous Treatment Effects using Random Forests. *ArXiv e-prints*, October 2015.

Wang, Y.-X., Agarwal, A., and Dudík, M. Optimal and adaptive off-policy evaluation in contextual bandits. 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. 3589–3597, International Convention Centre, Sydney, Australia, 06–11 Aug 2017. PMLR. URL http://proceedings.mlr.press/v70/wang17a.html.
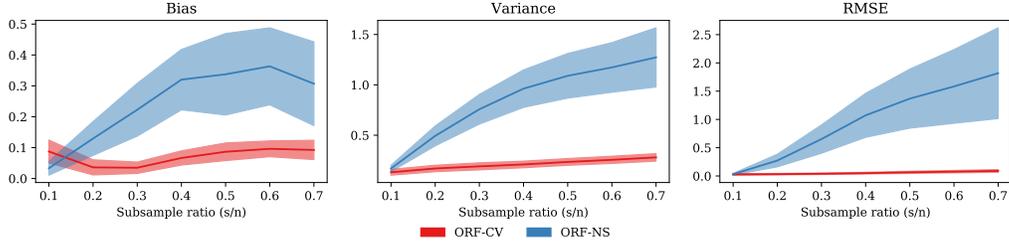
Figure 5: Bias, variance and RMSE versus subsample ratio used for training individual trees. The solid lines represent the means and the filled regions depict the standard deviation for the different metrics across test points, averaged over 100 Monte Carlo experiments.

## A. Two-forest ORF v.s. One-forest ORF

A natural question about ORF is whether it is necessary to have two separate forests for the two-stage estimation. We investigate this question by implementing a variant of ORF without sample splitting (ORF-NS)—it builds only one random forest over the entire dataset, and perform the two-stage estimation using the same set of importance weights. We empirically compare ORF-CV with ORF-NS. In Figure 5, we note that the bias, variance and RMSE of the ORF-NS increase drastically with the subsample ratio $(s/n)$, whereas the same metrics are almost constant for the ORF-CV. This phenomenon is consistent with the theory, since larger subsamples induce a higher probability of collision between independently drawn samples, and the "spill-over" can incur large bias and error.

## B. Formal Guarantee for Real-Valued Treatments

**Corollary B.1** (Accuracy for real-valued treatments). *Suppose that $\beta_0(X)$ and each coorindate $u_0^j(X), \gamma_0^j(X)$ are Lipschitz in $X$ and have $\ell_1$ norms bounded by $O(1)$ for any $X$. Assume that distribution of $X$ admits a density that is bounded away from zero and infinity. For any feature $X$, the conditional covariance matrices satisfy $\mathbb{E}\left[\eta\eta^\intercal \mid X\right] \succeq \Omega(1)I_p$, $\mathbb{E}\left[WW^\intercal \mid X\right] \succeq \Omega(1)I_{d_\nu}$ and $\mathbb{E}\left[\varphi_2(W)\varphi_2(W)^\intercal \mid X\right] \succeq \Omega(1)I_{d_\nu^2+d_\nu}$, where $\varphi_2(W)$ denotes the degree-2 polynomial feature vector of $W$. Then with probability $1-\delta$, ORF returns an estimator $\hat\theta$ such that*

$$\|\hat\theta - \theta_0\| \leq O\left(n^{\frac{-1}{2+2\alpha d}}\sqrt{\log(nd_\nu/\delta)}\right)$$

*as long as the sparsity $k \leq O\left(n^{\frac{1}{8+8\alpha d}}\right)$ and the subsampling rate of ORF $s = \Theta(n^{\alpha d/(1+\alpha d)})$. Moreover, for any $b \in \mathbb{R}^p$ with $\|b\| \leq 1$, there exists a sequence $\sigma_n = \Theta(\sqrt{\text{polylog}(n)}n^{-1/(1+\alpha d)})$ such that*

$$\sigma_n^{-1}\left\langle b, \hat\theta - \theta\right\rangle \to_d \mathcal{N}(0,1),$$

*as long as the sparsity $k = o\left(n^{\frac{1}{8+8\alpha d}}\right)$ and the subsampling rate of ORF $s = \Theta(n^{\varepsilon+\alpha d/(1+\alpha d)})$ for any $\varepsilon > 0$.*

## C. Uniform Convergence of Lipschitz $U$-Processes

**Lemma C.1** (Stochastic Equicontinuity for $U$-statistics via Bracketing). *Consider a parameter space $\Theta$ that is a bounded subset of $\mathbb{R}^p$, with $\text{diam}(\Theta) = \sup_{\theta,\theta'\in\Theta} \|\theta - \theta'\|_2 \leq R$. Consider the $U$-statistic over $n$ samples of order $s$:*

$$\mathbb{G}_{s,n}f(\cdot;\theta) = \binom{n}{s}^{-1} \sum_{1 \leq i_1 \leq \ldots \leq i_s \leq n} f(z_{i_1}, \ldots, z_{i_s}; \theta) \tag{15}$$

*where $f(\cdot;\theta) : \mathcal{Z}^s \to \mathbb{R}$ is a known symmetric function in its first $s$ sarguments and $L$-Lipschitz in $\theta$. Suppose that $\sup_{\theta\in\Theta} \sqrt{\mathbb{E}\left[f(Z_1,\ldots,Z_s;\theta)^2\right]} \leq \eta$ and $\sup_{\theta\in\Theta,Z_1,\ldots,Z_s\in\mathcal{Z}^s} f(Z_1,\ldots,Z_s;\theta) \leq G$. Then w.p. $1-\delta$:*

$$\sup_{\theta\in\Theta} |\mathbb{G}_{s,n}f(\cdot;\theta) - \mathbb{E}[f(Z_{1:s};\theta)]| = O\left(\eta\sqrt{\frac{s\left(\log(n/s)+\log(1/\delta)\right)}{n}} + (G+L\,R)\frac{s(\log(n/s)+\log(1/\delta))}{n}\right) \tag{16}$$

*Proof of Lemma C.1.* Note that for any fixed $\theta \in \Theta$, $\Psi_s(\theta, Z_{1:n})$ is a U-statistic of order $s$. Therefore by the Bernstein inequality for $U$-statistics (see e.g. Theorem 2 of (Peel et al., 2010)), for any fixed $\theta \in \Theta$, w.p. $1 - \delta$

$$|\mathbb{G}_{s,n} f(\cdot; \theta) - \mathbb{E}[f(Z_{1:s}; \theta)]| \leq \eta \sqrt{\frac{2 \log(1/\delta)}{n/s}} + G \frac{2 \log(1/\delta)}{3(n/s)}$$

Since $\text{diam}(\Theta) \leq R$, we can find a finite space $\Theta_\varepsilon$ of size $R/\varepsilon$, such that for any $\theta \in \Theta$, there exists $\theta_\varepsilon \in \Theta_\varepsilon$ with $\|\theta - \theta_\varepsilon\| \leq \varepsilon$. Moreover, since $f$ is $L$-Lipschitz with respect to $\theta$:

$$|\mathbb{G}_{s,n} f(\cdot; \theta) - \mathbb{E}[f(Z_{1:s}; \theta)]| \leq |\mathbb{G}_{s,n} f(\cdot; \theta_\varepsilon) - \mathbb{E}[f(Z_{1:s} \theta_\varepsilon)]| + 2L\|\theta - \theta_\varepsilon\|$$

Thus we have that:

$$\sup_{\theta \in \Theta} |\mathbb{G}_{s,n} f(\cdot; \theta) - \mathbb{E}[f(Z_{1:s}; \theta)]| \leq \sup_{\theta \in \Theta_\varepsilon} |\mathbb{G}_{s,n} f(\cdot; \theta_\varepsilon) - \mathbb{E}[f(Z_{1:s} \theta_\varepsilon)]| + 2L\varepsilon$$

Taking a union bound over $\theta \in \Theta_\varepsilon$, we have that w.p. $1 - \delta$:

$$\sup_{\theta \in \Theta_\varepsilon} |\mathbb{G}_{s,n} f(\cdot; \theta_\varepsilon) - \mathbb{E}[f(Z_{1:s} \theta_\varepsilon)]| \leq \eta \sqrt{\frac{2 \log(R/(\varepsilon \, \delta))}{n/s}} + G \frac{2 \log(R/(\varepsilon \, \delta))}{3(n/s)}$$

Choosing $\varepsilon = \frac{sR}{n}$ and applying the last two inequalities, yields the desired result. $\qquad\square$

## D. Estimation Error and Asymptotic Normality

Since throughout the section we will fix the target vector $x$, we will drop it from the notation when possible, e.g. we will let $\theta_0 = \theta_0(x)$ and $\hat{\theta} = \hat{\theta}(x)$. We begin by introducing some quantities that will be useful throughout our theoretical analysis. First we denote with $\omega$ the random variable that corresponds to the internal randomness of the tree-splitting algorithm. Moreover, when the tree splitting algorithm is run with target $x$, an input dataset of $\{Z_i\}_{i=1}^s$ and internal randomness $\omega$, we denote with $\alpha_i(\{Z_i\}_{i=1}^s, \omega)$ the weight that it assigns to the sample with index $i$. Finally, for each sub-sample $b = 1 \ldots B$ we denote with $S_b$ the index of the samples chosen and $\omega_b$ the internal randomness that was drawn.

We then consider the weighted empirical score, weighted by the sub-sampled ORF weights:

$$\Psi(\theta, h) = \sum_{i=1}^n a_i \, \psi(Z_i; \theta, h(W_i)) = \frac{1}{B} \sum_{b=1}^B \sum_{i \in S_b} \alpha_i(\{Z_i\}_{i \in S_b}, \omega_b) \, \psi(Z_i; \theta, h(W_i)) \tag{17}$$

We will also be considering the complete multi-dimensional $U$-statistic, where we average over all sub-samples of size $s$:

$$\Psi_0(\theta, h) = \binom{n}{s}^{-1} \sum_{1 \leq i_1 \leq \ldots \leq i_s \leq n} \mathbb{E}_\omega \left[ \sum_{t=1}^s \alpha_{i_t}(\{Z_{i_t}\}_{t=1}^s, \omega) \, \psi(Z_{i_t}; \theta, h(W_{i_t})) \right]. \tag{18}$$

and we denote with:

$$f(Z_{i_1}, \ldots, Z_{i_s}; \theta, h) = \mathbb{E}_\omega \left[ \sum_{t=1}^s \alpha_{i_t}(\{Z_{i_t}\}_{t=1}^s, \omega) \, \psi(Z_{i_t}; \theta, h(W_{i_t})) \right] \tag{19}$$

First, we will bound the estimation error as a sum of $m(x; \theta, \hat{h})$ and second order terms. The proof follows from the Taylor expansion of the moment function and the mean-value theorem.

**Lemma D.1.** *Under Assumption 4.1, for any nuisance estimate $\hat{h}$ and for the ORF estimate $\hat{\theta}$ estimated with plug-in nuisance estimate $\hat{h}$:*

$$\hat{\theta} - \theta_0 = M^{-1}\left(m(x; \hat{\theta}, \hat{h}) - \Psi(\hat{\theta}, \hat{h})\right) + \xi$$

*where $\xi$ satisfies $\|\xi\| = \mathcal{O}\left(\mathbb{E}\left[\|\hat{h}(W) - h_0(W)\|^2 \mid x\right] + \|\hat{\theta} - \theta_0\|^2\right)$.*

**Proof outline of main theorems.** We will now give a rough outline of the proof of our main results. In doing so we will also present some core technical Lemmas that we will use in the formal proofs of these theorems in the subsequent corresponding subsections.

Lemma D.1 gives rise to the following core quantity:

$$\Lambda(\theta, h) = m(x; \theta, h) - \Psi(\theta, h) \tag{20}$$

Suppose that our first stage estimation rate guarantees a local root-mean-squared-error (RMSE) of $\chi_n$, i.e.:

$$\mathcal{E}(h) = \sqrt{\mathbb{E}\left[\|\hat{h}(W) - h_0(W)\|^2 \mid x\right]} \leq \chi_n \tag{21}$$

Then we have that:

$$\hat{\theta} - \theta_0 = M^{-1}\Lambda(\hat{\theta}, \hat{h}) + O(\chi_n^2 + \|\hat{\theta} - \theta_0\|^2)$$

Thus to understand the estimation error of $\hat{\theta}$ and its asymptotic distribution, we need to analyze the concentration of $\Lambda(\hat{\theta}, \hat{h})$ around zero and its asymptotic distribution. Subsequently, invoking consistency of $\hat{\theta}$ and conditions on a sufficiently fast nuisance estimation rate $\chi_n$, we would be able to show that the remainder terms are asymptotically negligible.

Before delving into our two main results on mean absolute error (MAE) and asymptotic normality we explore a bit more the term $\Lambda(\theta, h)$ and decompose it into three main quantities, that we will control each one separately via different arguments.

**Lemma D.2** (Error Decomposition). *For any $\theta, h$, let $\mu_0(\theta, h) = \mathbb{E}[\Psi_0(\theta, h)]$. Then:*

$$\Lambda(\theta, h) = \underbrace{m(x; \theta, h) - \mu_0(\theta, h)}_{\Gamma(\theta,h)\,=\,\text{kernel error}} + \underbrace{\mu_0(\theta, h) - \Psi_0(\theta, h)}_{\Delta(\theta,h)\,=\,\text{sampling error}} + \underbrace{\Psi_0(\theta, h) - \Psi(\theta, h)}_{E(\theta,h)\,=\,\text{subsampling error}}. \tag{22}$$

When arguing about the MAE of our estimator, the decomposition presented in Lemma D.2 is sufficient to give us the final result by arguing about concentration of each of the terms. However, for asymptotic normality we need to further refine the decomposition into terms that when scaled appropriately converge to zero in probability and terms that converge to a normal random variable. In particular, we need to further refine the sampling error term $\Delta(\theta, h)$ as follows:

$$\Delta(\theta, h) = \underbrace{\Delta(\theta_0, \tilde{h}_0)}_{\text{asymptotically normal term}} + \underbrace{\Delta(\theta, h) - \Delta(\theta_0, \tilde{h}_0)}_{F(\theta,h)\,=\,\text{stochastic equicontinuity term}} \tag{23}$$

for some appropriately defined fixed function $\tilde{h}_0$. Consider for instance the case where $\theta$ is a scalar. If we manage to show that there exists a scaling $\sigma_n$, such that $\sigma_n^{-1}\Delta(\theta_0, \tilde{h}_0) \to_d \mathcal{N}(0, 1)$, and all other terms $\Gamma, E, F$ and $\chi_n^2$ converge to zero in probability when scaled by $\sigma_n^{-1}$, then we will be able to conclude by Slutzky's theorem that: $\sigma_n^{-1}M\left(\hat{\theta} - \theta_0\right) \to_d \mathcal{N}(0, 1)$ and establish the desired asymptotic normality result.

Since controlling the convergence rate to zero of the terms $\Gamma, \Delta, E$ would be useful in both results, we provide here three technical lemmas that control these rates.

**Lemma D.3** (Kernel Error). *If the ORF weights when trained on a random sample $\{Z_i\}_{i=1}^s$, satisfy that:*

$$\mathbb{E}\left[\sup\{\|X_i - x\| : a_i(\{Z_i\}_{i=1}^s, \omega) > 0\}\right] \leq \varepsilon(s) \tag{24}$$

*where expectation is over the randomness of the samples and the internal randomness $\omega$ of the ORF algorithm. Then*

$$\sup_{\theta,h}\|\Gamma(\theta, h)\| = \sqrt{p}\,L\,\varepsilon(s) \tag{25}$$

**Lemma D.4** (Sampling Error). *Under Assumption 4.1, conditional on any nuisance estimate $\hat{h}$ from the first stage, with probability $1 - \delta$:*

$$\sup_{\theta}\|\Delta(\theta, \hat{h})\| = O\left(\sqrt{\frac{s\,(\log(n/s) + \log(1/\delta))}{n}}\right) \tag{26}$$

*Proof.* Since $\Delta(\theta, \hat{h})$ is a $U$-statistic as it can be written as: $\mathbb{G}_{s,n} f(\cdot; \theta, \hat{h}) - \mathbb{E}\left[f(Z_{1:s}; \theta, \hat{h})\right]$. Moreover, under Assumption 4.1, the function $f(\cdot; \theta, \hat{h})$ satisfies the conditions of Lemma C.1 with $\eta = G = \psi_{\max} = O(1)$. Moreover, $f(\cdot; \theta, \hat{h})$ is $L$-Lipschitz for $L = O(1)$, since it is a convex combination of $O(1)$-Lipschitz functions. Finally, $\text{diam}(\Theta) = O(1)$. Thus applying Lemma C.1, we get, the lemma. $\qquad\square$

**Lemma D.5** (Subsampling Error). *If the ORF weights are built on $B$ randomly drawn sub-samples with replacement, then*

$$\sup_{\theta, h} \|E(\theta, \hat{h})\| = O\left(\frac{\log(B) + \log(1/\delta)}{\sqrt{B}}\right) \tag{27}$$

### D.1. Consistency of ORF Estimate

**Theorem D.6** (Consistency). *Assume that the nuisance estimate satisfies:*

$$\mathbb{E}\left[\mathcal{E}(\hat{h})\right] = o(1) \tag{28}$$

*and that $B \geq n/s$, $s = o(n)$ and $s \to \infty$ as $n \to \infty$. Then the ORF estimate $\hat{\theta}$ satisfies:*

$$\|\hat{\theta} - \theta_0(x)\| = o_p(1)$$

*Moreover, for any constant integer $q \geq 1$:*

$$\left(\mathbb{E}[\|\hat{\theta} - \theta_0\|^{2q}]\right)^{1/q} = o\left(\left(\mathbb{E}\left[\|\hat{\theta} - \theta_0\|^q\right]\right)^{1/q}\right) \tag{29}$$

*Proof.* By the definition of $\hat{\theta}$, we have that: $\Psi(\hat{\theta}, \hat{h}) = 0$. Thus we have that:

$$\left\|m(x; \hat{\theta}, \hat{h})\right\| = \left\|m(x; \hat{\theta}, \hat{h}) - \Psi(\hat{\theta}, \hat{h})\right\| = \left\|\Lambda(\hat{\theta}, \hat{h})\right\|$$

By Lemmas 3.1, D.2, D.3, D.4 and D.5, we have that with probability $1 - 2\delta$:

$$\left\|\Lambda(\hat{\theta}, \hat{h})\right\| = O\left(s^{-1/(2\alpha d)} + \sqrt{\frac{s\left(\log(n/s) + \log(1/\delta)\right)}{n}} + \sqrt{\frac{\log(B) + \log(1/\delta)}{B}}\right)$$

Integrating this tail bound we get that:

$$\mathbb{E}\left[\left\|\Lambda(\hat{\theta}, \hat{h})\right\|\right] = O\left(s^{-1/(2\alpha d)} + \sqrt{\frac{s\log(n/s)}{n}} + \sqrt{\frac{\log(B)}{B}}\right)$$

Thus if $B \geq n/s$, $s = o(n)$ and $s \to \infty$ then all terms converge to zero as $n \to \infty$.

Since $\psi(x; \theta, h(w))$ is $L$-Lipschitz in $h(w)$ for some constant $L$:

$$\|m(x; \hat{\theta}, h_0) - m(x; \hat{\theta}, \hat{h})\| = L\mathbb{E}\left[\left\|\hat{\theta}(W) - \hat{h}(W)\right\| \mid x\right] \leq L\sqrt{\mathbb{E}\left[\left\|\hat{\theta}(W) - \hat{h}(W)\right\|^2 \mid x\right]} = L\mathcal{E}(\hat{h})$$

Moreover, by our consistency guarantee on $\hat{h}$:

$$\mathbb{E}\left[\|m(x; \hat{\theta}, h_0) - m(x; \hat{\theta}, \hat{h})\|\right] \leq L\mathbb{E}\left[\mathcal{E}(\hat{h})\right] = o(1)$$

Thus we conclude that:

$$\mathbb{E}[\|m(x; \hat{\theta}, h_0)\|] = o(1)$$

which implies that $\|m(x; \hat{\theta}, h_0)\| = o_p(1)$.

By our first assumption, for any $\varepsilon$, there exists a $\delta$, such that: $\Pr[\|\hat{\theta} - \theta_0(x)\| \geq \varepsilon] \leq \Pr[\|m(x; \hat{\theta}, h_0)\| \geq \delta]$. Since $\|m(x; \hat{\theta}, h_0)\| = o_p(1)$, the probability on the right-hand-side converges to $0$ and hence also the left hand side. Hence, $\|\hat{\theta} - \theta_0(x)\| = o_p(1)$.

We now prove the second part of the theorem which is a consequence of consistency. By consistency of $\hat{\theta}$, we have that for any $\varepsilon$ and $\delta$, there exists $n^*(\varepsilon, \delta)$ such that for all $n \geq n(\varepsilon, \delta)$:

$$\Pr\left[\|\hat{\theta} - \theta_0\| \geq \varepsilon\right] \leq \delta$$

Thus for any $n \geq n^*(\varepsilon, \delta)$:

$$\mathbb{E}[\|\hat{\theta} - \theta_0\|^{2q}] \leq \varepsilon^q \mathbb{E}\left[\|\hat{\theta} - \theta_0\|^q\right] + \delta \mathbb{E}\left[|\hat{\theta} - \theta_0\|^{2q}\right]$$

Choosing $\varepsilon = (4C)^{-q}$ and $\delta = (4C)^{-1}(\mathrm{diam}(\Theta))^{-q}$ yields that:

$$\mathbb{E}[\|\hat{\theta} - \theta_0\|^{2q}] \leq \frac{1}{2C}\mathbb{E}\left[\|\hat{\theta} - \theta_0\|^q\right]$$

Thus for any constant $C$ and for $n \geq n^*((4C)^{-q}, (4C)^{-1}(\mathrm{diam}(\Theta))^{-q}) = O(1)$, we get that:

$$\left(\mathbb{E}[\|\hat{\theta} - \theta_0\|^{2q}]\right)^{1/q} = \frac{1}{(2C)^{1/q}}\left(\mathbb{E}\left[\|\hat{\theta} - \theta_0\|^q\right]\right)^{1/q}$$

which concludes the claim that $\left(\mathbb{E}[\|\hat{\theta} - \theta_0\|^{2q}]\right)^{1/q} = o\left(\left(\mathbb{E}\left[\|\hat{\theta} - \theta_0\|^q\right]\right)^{1/q}\right)$. $\qquad \square$

### D.2. Proof of Theorem 4.2: Mean $L^q$ Estimation Error

*Proof.* Applying Lemma D.1 and the triangle inequality for the $L^q$ norm, we have that:

$$\left(\mathbb{E}\left[\|\hat{\theta} - \theta_0\|^q\right]\right)^{1/q} = O\left(\left(\mathbb{E}\left[\|\Lambda(\hat{\theta}, \hat{h})\|^q\right]\right)^{1/q} + \left(\mathcal{E}(\hat{h})^{2q}\right)^{1/q} + \left(\mathbb{E}\left[\|\hat{\theta} - \theta_0\|^{2q}\right]\right)^{1/q}\right)$$

By assumption $\left(\mathcal{E}(\hat{h})^{2q}\right)^{1/q} \leq \chi_{n,2q}^2$. By the consistency Theorem D.6:

$$\left(\mathbb{E}\left[\|\hat{\theta} - \theta_0\|^{2q}\right]\right)^{1/q} = o\left(\left(\mathbb{E}\left[\|\hat{\theta} - \theta_0\|^q\right]\right)^{1/q}\right).$$

and therefore this term can be ignored for $n$ larger than some constant. Moreover, by Lemma D.2:

$$\left(\mathbb{E}\left[\|\Lambda(\hat{\theta}, \hat{h})\|^q\right]\right)^{1/q} = O\left(\left(\mathbb{E}\left[\|\Gamma(\hat{\theta}, \hat{h})\|^q\right]\right)^{1/q} + \left(\mathbb{E}\left[\|\Delta(\hat{\theta}, \hat{h})\|^q\right]\right)^{1/q} + \left(\mathbb{E}\left[\|E(\hat{\theta}, \hat{h})\|^q\right]\right)^{1/q}\right)$$

By Lemma D.3 and Lemma 3.1 we have:

$$\left(\mathbb{E}\left[\|\Gamma(\hat{\theta}, \hat{h})\|^q\right]\right)^{1/q} = O\left(\varepsilon(s)\right) = O\left(s^{-1/(2\alpha d)}\right)$$

for a constant $\alpha = \frac{\log(\rho^{-1})}{\pi \log((1-\rho)^{-1})}$. Moreover, by integrating the exponential tail bound provided by the high probability statements in Lemmas D.4 and D.5, we have that for any constant $q$:

$$\left(\mathbb{E}\left[\|\Delta(\hat{\theta}, \hat{h})\|^q\right]\right)^{1/q} = O\left(\sqrt{\frac{s\log(n/s)}{n}}\right)$$

$$\left(\mathbb{E}\left[\|E(\hat{\theta}, \hat{h})\|^q\right]\right)^{1/q} = O\left(\sqrt{\frac{\log(B)}{B}}\right)$$

For $B > n/s$, the second term is negligible compared to the first and can be ignored. Combining all the above inequalities:

$$\left(\mathbb{E}\left[\|\hat{\theta} - \theta_0\|^q\right]\right)^{1/q} = O\left(s^{-1/(2\alpha d)} + \sqrt{\frac{s\log(n/s)}{n}}\right)$$

$\qquad \square$

**D.3. Proof of Theorem 4.3: Finite Sample High Probability Error Bound for Gradients of Convex Losses**

*Proof.* We condition on the event that $\mathcal{E}(\hat{h}) \leq \chi_{n,\delta}$, which occurs with probability $1 - \delta$. Since the Jacobian of $m(x; \theta, h_0)$ has eigenvalues lower bounded by $\sigma$ and each entry of the Jacobian of $\psi$ is $L$-Lispchitz with respect to the nuisance for some constant $L$, we have that for every vector $\nu \in \mathbb{R}^p$ (with $p$ the dimension of $\theta_0$):

$$\frac{\nu^T \nabla_\theta m(x; \theta, \hat{h}) \nu}{\|\nu\|^2} \geq \frac{\nu^T \nabla_\theta m(x; \theta, h_0) \nu}{\|\nu\|^2} + \frac{\nu^T \nabla_\theta \left( m(x; \theta, \hat{h}) - m(x; \theta, h_0) \right) \nu}{\|\nu\|^2}$$

$$\geq \sigma - L \cdot \mathbb{E}\left[ \|\hat{h}(W) - h_0(W)\| \mid x \right] \frac{\|\nu\|_1^2}{\|\nu\|^2}$$

$$\geq \sigma - L \chi_{n,\delta} \, p = \sigma - O(\chi_{n,\delta})$$

Where in the last inequality we also used Holder's inequality to upper bound the $L^1$ norm by the $L^2$ norm of the first stage error. Thus the expected loss function $L(\theta) = \mathbb{E}\left[ \ell(Z; \theta, \hat{h}(W) \mid x \right]$ is $\hat{\sigma} = \sigma - O(\chi_{n,\delta})$ strongly convex, since $\nabla_\theta m(x; \theta, \hat{h})$ is the Hessian of $L(\theta)$. We then have:

$$L(\hat{\theta}) - L(\theta_0) \geq \nabla_\theta L(\theta_0)'(\hat{\theta} - \theta_0) + \frac{\hat{\sigma}}{2} \|\hat{\theta} - \theta_0\|^2 = m(x; \theta_0, \hat{h})'(\hat{\theta} - \theta_0) + \frac{\hat{\sigma}}{2} \|\hat{\theta} - \theta_0\|^2$$

Moreover, by convexity of $L(\theta)$, we have:

$$L(\theta_0) - L(\hat{\theta}) \geq \nabla_\theta L(\hat{\theta})'(\theta_0 - \hat{\theta}) = m(x; \hat{\theta}, \hat{h})'(\theta_0 - \hat{\theta})$$

Combining the above we get:

$$\frac{\hat{\sigma}}{2} \|\hat{\theta} - \theta_0\|^2 \leq (m(x; \hat{\theta}, \hat{h}) - m(x; \theta_0, \hat{h}))'(\hat{\theta} - \theta_0) \leq \|m(x; \hat{\theta}, \hat{h}) - m(x; \theta_0, \hat{h})\| \, \|\hat{\theta} - \theta_0\|$$

Dividing over by $\|\hat{\theta} - \theta_0\|$, we get:

$$\|\hat{\theta} - \theta_0\| \leq \frac{2}{\hat{\sigma}} \left( \|m(x; \hat{\theta}, \hat{h})\| + \|m(x; \theta_0, \hat{h})\| \right)$$

The term $\|m(x; \hat{\theta}, \hat{h})\|$ is upper bounded by $\|\Lambda(\hat{\theta}, \hat{h})\|$ (since $\Psi(\hat{\theta}, \hat{h}) = 0$). Hence, by Lemmas 3.1, D.2, D.3, D.4 and D.5 and our assumptions on the choice of $s, B$, we have that with probability $1 - 2\delta$:

$$\left\| \Lambda(\hat{\theta}, \hat{h}) \right\| = O\left( s^{-1/(2\alpha d)} + \sqrt{\frac{s\left(\log(n/s) + \log(1/\delta)\right)}{n}} \right)$$

Subsequently, using a second order Taylor expansion around $h_0$ and orthogonality argument almost identical to the proof of Lemma D.1, we can show that the second term $\|m(x; \theta_0, \hat{h})\|$ is it upper bounded by $O(\chi_{n,\delta}^2)$. More formally, since $m(x; \theta_0, h_0) = 0$ and the moment is locally orthogonal, invoking a second order Taylor expansion:

$$m_j(x; \theta_0, \hat{h}) = m_j(x; \theta_0, h_0) + D_{\psi_j}[\hat{h} - h_0 \mid x] + \underbrace{\frac{1}{2} \mathbb{E}\left[ (\hat{h}(W) - h_0(W))^\intercal \nabla_h^2 \psi_j(Z; \theta_0, \tilde{h}^{(j)}(W))(\hat{h}(W) - h_0(W)) \mid x \right]}_{\rho_j}$$

$$= \rho_j$$

for some function $\tilde{h}_j$ implied by the mean value theorem. Since the moment is smooth, we have: $\|\rho\| = O\left( \mathbb{E}\left[ \left\| \hat{h}(W) - h_0(W) \right\|^2 \mid x \right] \right) = O\left( \chi_{n,\delta}^2 \right)$. Thus $\left\| m_j(x; \theta_0, \hat{h}) \right\| = O\left( \chi_{n,\delta}^2 \right)$. Combining all the latter inequalities yields the result. □

**D.4. Proof of Theorem 4.4: Asymptotic Normality**

*Proof.* We want to show asymptotic normality of any fixed projection $\left\langle \beta, \hat{\theta} \right\rangle$ with $\|\beta\| \leq 1$. First consider the random variable $V = \left\langle \beta, M^{-1}\Delta(\theta_0, \tilde{h}_0) \right\rangle$, where $\tilde{h}_0(X, W) = g(W; \nu_0(x))$, i.e. the nuisance function $\tilde{h}_0$ ignores the input $X$ and uses the parameter $\nu_0(x)$ for the target point $x$. Asymptotic normality of $V$ follows by identical arguments as in (Wager & Athey, 2015) or (Mentch & Hooker, 2016), since this term is equivalent to the estimate of a random forest in a regression setting, where we want to estimate $\mathbb{E}[Y \mid X = x]$ and where the observation of sample $i$ is:

$$Y_i = \left\langle \beta, M^{-1} \left( m(X_i; \theta_0, \tilde{h}_0) - \psi(Z_i; \theta_0, \tilde{h}_0(X_i, W_i)) \right) \right\rangle \tag{30}$$

By Theorem 1 of (Wager & Athey, 2015) and the fact that our forest satisfies Specification 1 and under our set of assumptions, we have, that there exists a sequence $\sigma_n$, such that:

$$\sigma_n^{-1} V \to \mathcal{N}(0, 1) \tag{31}$$

for $\sigma_n = \Theta\left( \sqrt{\operatorname{polylog}(n/s)^{-1} s/n} \right)$. More formally, we check that each requirement of Theorem 1 of (Wager & Athey, 2015) is satisfied:

(i) We assume that the distribution of $X$ admits a density that is bounded away from zero and infinity,

(ii) $\mathbb{E}[Y|X = x^*] = 0$ and hence is continuous in $x^*$ for any $x^*$,

(iii) The variance of the $Y$ conditional on $X = x^*$ for some $x^*$ is:

$$\operatorname{Var}(Y|X = x^*) = \mathbb{E}\left[ \left\langle \beta, M^{-1}\psi(Z; \theta_0, \tilde{h}_0(X, W)) \right\rangle^2 \mid X = x^* \right] - \mathbb{E}\left[ \left\langle \beta, M^{-1}\psi(Z; \theta_0, \tilde{h}_0(X, W)) \mid X = x^* \right\rangle \right]^2$$

The second term is $O(1)$-Lipschitz in $x^*$ by Lipschitzness of $m(x^*; \theta_0, \tilde{h}_0) = \mathbb{E}\left[ \psi(Z; \theta_0, \tilde{h}_0(X, W)) \mid X = x^* \right]$. For simplicity of notation consider the random variable $V = \psi(Z; \theta_0, \tilde{h}_0(X, W))$. Then the first part is equal to some linear combination of the covariance terms:

$$Q(x^*) \triangleq \mathbb{E}\left[ \beta^{\intercal} M^{-1} V V^T (M^{-1})^{\intercal} \beta \mid X = x^* \right] = \beta^{\intercal} M^{-1} \mathbb{E}\left[ V V^T \mid X = x^* \right] (M^{-1})^{\intercal} \beta =$$

Thus by Lipschitzness of the covariance matrix of $\psi$, we have that: $\|\mathbb{E}\left[ V V^{\intercal} \mid X = x^* \right] - \mathbb{E}\left[ V V^{\intercal} \mid X = \tilde{x} \right]\|_F \leq L\|x^* - \tilde{x}\|$ and therefore by the Cauchy-Schwarz inequality and the lower bound $\sigma > 0$ on the eigenvalues of $M$:

$$|Q(x^*) - Q(\tilde{x})| \leq L\|x^* - \tilde{x}\| \|\beta^{\intercal} M^{-1}\|^2 \leq \frac{L}{\sigma^2}\|x^* - \tilde{x}\|$$

Thus $\operatorname{Var}(Y|X = x^*)$ is $O(1)$-Lipschitz continuous in $x^*$ and hence also $\mathbb{E}[Y^2|X = x^*]$ is $O(1)$-Lipschitz continuous.

(iv) The fact that $\mathbb{E}[|Y - \mathbb{E}[Y|X = x]|^{2+\delta}|X = x] \leq H$ for some constant $\delta, H$ follows by our assumption on the boundedness of $\psi$ and the lower bound on the eigenvalues of $M$,

(v) The fact that $\operatorname{Var}[Y|X = x'] > 0$ follows from the fact that $\operatorname{Var}\left( \beta^{\intercal} M^{-1}\psi(Z; \theta_0, \tilde{h}_0(X, W)) \mid X = x' \right) > 0$,

(vi) The fact that tree is honest, $\alpha$-balanced with $\alpha \leq 0.2$ and symmetric follows by Specification 1,

(vii) From our assumption on $s$ that $s^{-1/(2\alpha d)} = o((s/n)^{1/2-\varepsilon})$, it follows that $s = \Theta(n^\beta)$ for some $\beta \in \left( 1 - \frac{1}{1+\alpha d}, 1 \right]$.

Since, by Lemmas D.1, D.2 and Equation (23):

$$\left\| \left\langle \beta, \hat{\theta} - \theta_0 \right\rangle - V \right\| = O\left( \|\Gamma(\hat{\theta}, \hat{h})\| + \|\Delta(\hat{\theta}, \hat{h}) - \Delta(\theta_0, \tilde{h}_0)\| + \|\mathcal{E}(\hat{h})\|^2 + \|\hat{\theta} - \theta_0\|^2 \right)$$

it suffices to show that:

$$\sigma_n^{-1} \mathbb{E}\left[ \|\Gamma(\hat{\theta}, \hat{h})\| + \|\Delta(\hat{\theta}, \hat{h}) - \Delta(\theta_0, \tilde{h}_0)\| + \|\mathcal{E}(\hat{h})\|^2 + \|\hat{\theta} - \theta_0\|^2 \right] \to 0$$

as then by Slutzky's theorem we have that $\sigma_n^{-1} \left\langle \beta, \hat{\theta} - \theta_0 \right\rangle \to_d \mathcal{N}(0,1)$. The first term is of order $O(s^{-1/(2\alpha d)})$, hence by our assumption on the choice of $s$, it is $o(\sigma_n)$. The third term is $O(\chi_{n,2}^2) = O(\chi_{n,4}^2)$, which by assumption is also $o(\sigma_n)$. The final term, by applying our $L^q$ estimation error result for $q = 2$ and the assumption on our choice of $s$, we get that it is of order $O\left(\frac{s\log(n/s)}{n}\right) = o(\sigma_n)$.

Thus it remains to bound the second term. For that we will invoke the stochastic equicontinuity Lemma C.1. Observe that each coordinate $j$ of the term corresponds to the deviation from its mean of a $U$ statistic with respect to the class of functions:

$$\gamma_j(\cdot; \theta, \hat{h}) = f_j(\cdot; \theta, \hat{h}) - f_j(\cdot; \theta_0, h_0) \tag{32}$$

Observe that by Lipschitzness of $\psi$ with respect to $\theta$ and the output of $h$ and the locally parametric form of $h$, we have that:

$$
\begin{aligned}
|\gamma_j(Z_{1:s}; \theta, h)| &= \left| \mathbb{E}_\omega \left[ \sum_{t=1}^s \alpha_t \left( \{Z_t\}_{t=1}^s, \omega \right) \left( \psi_j(Z_t; \theta, \hat{h}(W_t)) - \psi_j(Z_t; \theta, \tilde{h}_0(W_t)) \right) \right] \right| \\
&\leq \mathbb{E}_\omega \left[ \sum_{t=1}^s \alpha_t \left( \{Z_t\}_{t=1}^s, \omega \right) \left| \psi_j(Z_t; \theta, \hat{h}(W_t)) - \psi_j(Z_t; \theta_0, \tilde{h}_0(W_t)) \right| \right] \\
&\leq L\, \mathbb{E}_\omega \left[ \sum_{t=1}^s \alpha_t \left( \{Z_t\}_{t=1}^s, \omega \right) \left( \|\theta - \theta_0\| + \|g(W_t; \nu) - g(W_t; \nu_0(x))\| \right) \right] \\
&\leq L\, \mathbb{E}_\omega \left[ \sum_{t=1}^s \alpha_t \left( \{Z_t\}_{t=1}^s, \omega \right) \left( \|\theta - \theta_0\| + L \|\nu - \nu_0(x)\| \right) \right] \\
&= L \left( \|\theta - \theta_0\| + L \|\nu - \nu_0(x)\| \right)
\end{aligned}
$$

Thus by Jensen's inequality and the triangle inequality:

$$\sqrt{\mathbb{E}\left[ |\gamma_j(Z_{1:s}; \theta, h)|^2 \right]} \leq L\|\theta - \theta_0\| + L^2 \|\nu - \nu_0(x)\|$$

Thus:

$$\sup_{\theta : \|\theta - \theta_0\| \leq \eta, \|\nu - \nu_0(x)\| \leq \gamma} \sqrt{\mathbb{E}\left[ |\gamma_j(Z_{1:s}; \theta, g(\cdot; \nu))|^2 \right]} = O(\eta + \gamma)$$

By our $L^q$ error result and Markov's inequality, we have that with probability $1 - \delta$: $\|\hat{\theta} - \theta_0\| \leq \eta = O(\sigma_n/\delta)$. Similarly, by our assumption on the nuisance error $\left( \mathbb{E}\left[ \|\hat{\nu} - \nu_0(x)\|^4 \right] \right)^{1/4} \leq \chi_{n,4}$ and Markov's inequality we have that with probability $1 - \delta$: $\|\hat{\nu} - \nu_0(x)\| \leq O(\chi_{n,4}/\delta)$. Thus applying Lemma C.1, we have that conditional on the event that $\|\hat{\nu} - \nu_0(x)\| \leq O(\chi_{n,4}/\delta)$, w.p. $1 - \delta$:

$$
\begin{aligned}
\sup_{\theta : \|\theta - \theta_0\| \leq \sigma_n/\delta} \sqrt{\mathbb{E}\left[ |\gamma_j(Z_{1:s}; \theta, \hat{h})|^2 \right]} &= O\left( (\sigma_n/\delta + \chi_{n,4}/\delta) \sqrt{\frac{s(\log(n/s) + \log(1/\delta))}{n}} + \frac{s(\log(n/s) + \log(1/\delta))}{n} \right) \\
&= O\left( \sigma_n^2 \operatorname{polylog}(n/s)/\delta + \chi_{n,4}\sigma_n \operatorname{polylog}(n/s)/\delta + \frac{s(\log(n/s) + \log(1/\delta))}{n} \right) \\
&= O(\sigma_n^{3/2} \operatorname{polylog}(n/s)/\delta)
\end{aligned}
$$

where we used the fact that $\chi_{n,4}^2 = o(\sigma_n)$, $\sqrt{\log(1/\delta)} \leq 1/\delta$ and that $\sigma_n = \Theta\left( \sqrt{\operatorname{polylog}(n/s)^{-1} s/n} \right)$. By a union bound we have that w.p. $1 - 3\delta$:

$$\|\Delta(\hat{\theta}, \hat{h}) - \Delta(\theta_0, \tilde{h}_0)\| = O(\sigma_n^{3/2} \operatorname{polylog}(n/s)/\delta)$$

Integrating this tail bound and using the boundedness of the score we get:

$$\mathbb{E}\left[ \|\Delta(\hat{\theta}, \hat{h}) - \Delta(\theta_0, \tilde{h}_0)\| \right] = O(\sigma_n^{3/2} \operatorname{polylog}(n/s) \log(1/\sigma_n)) = o(\sigma_n) \tag{33}$$

This completes the proof of the theorem. $\qquad\square$

## D.5. Omitted Proofs of Technical Lemmas

*Proof of Lemma D.1.* Fix a conditioning vector $x$. By performing a second order Taylor expansion of each coordinate $j \in [p]$ of the expected score function $m_j$ around the true parameters $\theta_0 = \theta_0(x)$ and $h_0$ and applying the multi-dimensional mean-value theorem, we can write that for any $\theta \in \Theta$:

$$
\begin{aligned}
m_j(x; \theta, \hat{h}) = {} & m_j(x; \theta_0, h_0) + \nabla_\theta m_j(x; \theta_0, h_0)'(\theta - \theta_0) + D_{\psi_j}[\hat{h} - h_0 \mid x] \\
& + \underbrace{\frac{1}{2} \mathbb{E}\left[ (\theta - \theta_0, \hat{h}(W) - h_0(W))^\intercal \nabla_{\theta, h}^2 \psi_j(Z; \tilde{\theta}^{(j)}, \tilde{h}^{(j)}(W))(\theta - \theta_0, \hat{h}(W) - h_0(W)) \mid x \right]}_{\rho_j}
\end{aligned}
$$

where each $\tilde{\theta}^{(j)}$ is some convex combination of $\theta$ and $\theta_0$ and each $\tilde{h}^{(j)}(W)$ is some convex combination of $\hat{h}(W)$ and $h_0(W)$. Note that $m(x; \theta_0, h_0) = 0$ by definition and $D_{\psi_j}[\hat{h} - h_0 \mid x] = 0$ by local orthogonality. Let $\rho$ denote the vector of second order terms. We can thus write the above set of equations in matrix form as:

$$
M(\theta - \theta_0) = m(x; \theta, \hat{h}) - \rho
$$

where we remind that $M = \nabla_\theta m(x; \theta_0, h_0)$ is the Jacobian of the moment vector. Since by our assumptions $M$ is invertible and has eigenvalues bounded away from zero by a constant, we can write:

$$
(\theta - \theta_0) = M^{-1} m(x; \theta, \hat{h}) - M^{-1} \rho
$$

Letting $\xi = -M^{-1} \rho$, we have that by the boundedness of the eigenvalues of $M^{-1}$:

$$
\|\xi\| = O(\|\rho\|)
$$

By our bounded eigenvalue Hessian assumption on $\mathbb{E}\left[ \nabla_{\theta, h}^2 \psi_j(Z; \tilde{\theta}^{(j)}, \tilde{h}^{(j)}(W)) \mid x, W \right]$, we know that:

$$
\|\rho\|_\infty = O\left( \mathbb{E}\left[ \|\hat{h}(W) - h_0(W)\|^2 \mid x \right] + \|\theta - \theta_0\|^2 \right)
$$

Combining the above two equations and using the fact that $\|\rho\| \le \sqrt{p}\|\rho\|_\infty$, yields that for any $\theta \in \Theta$:

$$
\theta - \theta_0 = M^{-1}\left( m(x; \theta, \hat{h}) - \Psi(\theta, \hat{h}) \right) + \xi
$$

Evaluating the latter at $\theta = \hat{\theta}$ and also observing that by the definition of $\hat{\theta}$, $\Psi(\hat{\theta}, \hat{h}) = 0$ yields the result. $\qquad \square$

*Proof of Lemma D.3.* First we argue that by invoking the honesty of the ORF weights we can re-write $\mu_0(\theta, h)$ as:

$$
\mu_0(\theta, h) = \mathbb{E}\left[ \binom{n}{s}^{-1} \sum_{1 \le i_1 \le \dots \le i_s \le n} \mathbb{E}_\omega \left[ \sum_{t=1}^s \alpha_{i_t}\left( \{Z_{i_t}\}_{t=1}^s, \omega \right) m(X_{i_t}; \theta, h) \right] \right] \tag{34}
$$

To prove this claim, it suffices to show that for any subset of $s$ indices:

$$
\mathbb{E}\left[ \alpha_{i_t}\left( \{Z_{i_t}\}_{t=1}^s, \omega \right) \psi(Z_{i_t}; \theta, h) \right] = \mathbb{E}\left[ \alpha_{i_t}\left( \{Z_{i_t}\}_{t=1}^s, \omega \right) m(X_{i_t}; \theta, h) \right] \tag{35}
$$

By honesty of the ORF weights, we know that either $i_t \in S^1$, in which case $\alpha_{i_t}\left( \{Z_{i_t}\}_{t=1}^s, \omega \right) = 0$, or otherwise $i_t \in S^2$ and then $\alpha_{i_t}\left( \{Z_{i_t}\}_{t=1}^s, \omega \right)$ is independent of $Z_{i_t}$, conditional on $X_{i_t}, Z_{-i_t}, \omega$. Thus in any case $\alpha_{i_t}\left( \{Z_{i_t}\}_{t=1}^s, \omega \right)$ is independent of $Z_{i_t}$, conditional on $X_{i_t}, Z_{-i_t}, \omega$. Moreover since $Z_{i_t}$ is independent of $Z_{-i_t}, \omega$ conditional on $X_{i_t}$:

$$
\mathbb{E}\left[ \psi(Z_{i_t}; \theta, h) \mid X_{i_t}, Z_{-i_t}, \omega \right] = \mathbb{E}\left[ m(X_{i_t}; \theta, h) \mid X_{i_t}, Z_{-i_t}, \omega \right]
$$

By the law of iterated expectation and the independence properties claimed above, we can write:

$$
\begin{aligned}
\mathbb{E}\left[ \alpha_{i_t}\left( \{Z_{i_t}\}_{t=1}^s, \omega \right) \psi(Z_{i_t}; \theta, h) \right] = {} & = \mathbb{E}\left[ \mathbb{E}\left[ \alpha_{i_t}\left( \{Z_{i_t}\}_{t=1}^s, \omega \right) \mid X_{i_t}, Z_{-i_t}, \omega \right] \mathbb{E}\left[ \psi(Z_{i_t}; \theta, h) \mid X_{i_t}, Z_{-i_t}, \omega \right] \right] \\
& = \mathbb{E}\left[ \mathbb{E}\left[ \alpha_{i_t}\left( \{Z_{i_t}\}_{t=1}^s, \omega \right) \mid X_{i_t}, Z_{-i_t}, \omega \right] \mathbb{E}\left[ m(X_{i_t}; \theta, h) \mid X_{i_t}, Z_{-i_t}, \omega \right] \right] \\
& = \mathbb{E}\left[ \alpha_{i_t}\left( \{Z_{i_t}\}_{t=1}^s, \omega \right) m(X_{i_t}; \theta, h) \right]
\end{aligned}
$$

Finally, by a repeated application of the triangle inequality and the lipschitz property of the conditional moments, we have:

$$\|\Gamma(\theta,h)\| \leq \binom{n}{s}^{-1} \sum_{1 \leq i_1 \leq \dots \leq i_s \leq n} \mathbb{E}\left[\sum_{t=1}^{s} \alpha_{i_t}\left(\{Z_{i_t}\}_{t=1}^{s},\omega\right) \|m(x;\theta,h) - m(X_{i_t};\theta,h)\|\right]$$

$$\leq \sqrt{p}\, L \binom{n}{s}^{-1} \sum_{1 \leq i_1 \leq \dots \leq i_s \leq n} \mathbb{E}\left[\sum_{t=1}^{s} \alpha_{i_t}\left(\{Z_{i_t}\}_{t=1}^{s},\omega\right) \|X_{i_t} - x\|\right]$$

$$\leq \sqrt{p}\, L \binom{n}{s}^{-1} \sum_{1 \leq i_1 \leq \dots \leq i_s \leq n} \mathbb{E}\left[\sup\{\|X_{i_t} - x\| : \alpha_{i_t}\left(\{Z_{i_t}\}_{t=1}^{s},\omega\right) > 0\}\right]$$

$$\leq \sqrt{p}\, L\, \varepsilon(s)$$

$\square$

*Proof of Lemma D.5.* We prove that the concentration holds conditional on the samples $Z_{1:n}$ and $\hat{h}$, the result then follows. Let

$$\tilde{f}(S_b,\omega_b;\theta,h) = \sum_{i \in S_b} \alpha_i\left(S_b,\omega_b\right) \psi(Z_i;\theta,h(W_i)).$$

Observe that conditional on $Z_{1:n}$ and $\hat{h}$, the random variables $\tilde{f}(S_1,\omega_1;\theta,h),\dots,\tilde{f}(S_B,\omega_B;\theta,h)$ are conditionally independent and identically distributed (where the randomness is over the choice of the set $S_b$ and the internal algorithm randomness $\omega_b$). Then observe that we can write $\Psi(\theta,h) = \frac{1}{B}\sum_{b=1}^{B} \tilde{f}(S_b,\omega_b;\theta,h)$. Thus conditional on $Z_{1:n}$ and $\hat{h}$, $\Psi(\theta,\hat{h})$ is an average of $B$ independent and identically distributed random variables. Moreover, since $S_b$ is drawn uniformly at random among all sub-samples of $[n]$ of size $s$ and since the randomness of the algorithm is drawn identically and independently on each sampled tree:

$$\mathbb{E}\left[\tilde{f}(S_b,\omega_b;\theta,\hat{h}) \mid Z_{1:n}\right] = \binom{n}{s}^{-1} \sum_{1 \leq i_1 \leq \dots \leq i_s \leq n} f(\{Z_{i_t}\}_{t=1}^{s};\theta,\hat{h}) = \Psi_0(\theta,h)$$

Finally, observe that under Assumption 4.1, $|\tilde{f}(S_b,\omega_b;\theta,\hat{h})| \leq \psi_{\max} = O(1)$ a.s.. Thus by a Chernoff bound, we have that for any fixed $\theta \in \Theta$, w.p. $1 - \delta$:

$$\|\Psi(\theta,\hat{h}) - \Psi_0(\theta,\hat{h})\| \leq O\left(\sqrt{\frac{\log(1/\delta)}{B}}\right)$$

Since $\Theta$ has constant diameter, we can construct an $\varepsilon$-cover of $\Theta$ of size $O(1/\varepsilon)$. By Lipschitzness of $\psi$ with respect to $\theta$ and following similar arguments as in the proof of Lemma C.1, we can also get a uniform concentration:

$$\|\Psi(\theta,\hat{h}) - \Psi_0(\theta,\hat{h})\| \leq O\left(\sqrt{\frac{\log(B) + \log(1/\delta)}{B}}\right)$$

$\square$

# E. Omitted Proofs from Section 5

*Proof of Theorem 5.2.* By convexity of the loss $\ell$ and the fact that $\hat{\nu}(x)$ is the minimizer of the weighted penalized loss, we have:

$$\lambda\left(\|\nu_0(x)\|_1 - \|\hat{\nu}(x)\|_1\right) \geq \sum_{i=1}^{n} a_i(x)\,\ell(Z_i;\hat{\nu}(x)) - \sum_{i=1}^{n} a_i(x)\,\ell(Z_i;\nu_0(x)) \qquad \text{(optimality of } \hat{\nu}(x))$$

$$\geq \sum_{i=1}^{n} a_i(x)\,\langle \nabla_\nu \ell(z_i;\nu_0(x)), \hat{\nu}(x) - \nu_0(x)\rangle \qquad \text{(convexity of } \ell)$$

$$\geq -\left\|\sum_{i} a_i(x)\nabla_\nu \ell(z_i;\nu_0(x))\right\|_\infty \|\hat{\nu}(x) - \nu_0(x)\|_1 \qquad \text{(Cauchy-Schwarz)}$$

$$\geq -\frac{\lambda}{2}\|\hat{\nu}(x) - \nu_0(x)\|_1 \qquad \text{(assumption on } \lambda)$$

If we let $\rho(x) = \hat{\nu}(x) - \nu_0(x)$, then observe that by the definition of the support $S$ of $\nu_0(x)$ and the triangle inequality, we have:

$$
\begin{aligned}
\|\nu_0(x)\|_1 - \|\hat{\nu}(x)\|_1 &= \|\nu_0(x)_S\|_1 + \|\nu_0(x)_{S^c}\|_1 - \|\hat{\nu}(x)_S\|_1 - \|\hat{\nu}(x)_{S^c}\|_1 && \text{(separability of } \ell_1 \text{ norm)} \\
&= \|\nu_0(x)_S\|_1 - \|\hat{\nu}(x)_S\|_1 - \|\hat{\nu}(x)_{S^c}\|_1 && \text{(definition of support)} \\
&= \|\nu_0(x)_S\|_1 - \|\hat{\nu}(x)_S\|_1 - \|\hat{\nu}(x)_{S^c} - \nu_0(x)_{S^c}\|_1 && \text{(definition of support)} \\
&= \|\nu_0(x)_S - \hat{\nu}(x)_S\|_1 - \|\hat{\nu}(x)_{S^c} - \nu_0(x)_{S^c}\|_1 && \text{(triangle inequality)} \\
&\leq \|\rho(x)_S\|_1 - \|\rho(x)_{S^c}\|_1 && \text{(definition of } \rho(x))
\end{aligned}
$$

Thus re-arranging the terms in the latter series of inequalities, we get that $\rho(x) \in C(S(x); 3)$.

We now show that the weighted empirical loss function satisfies a conditional restricted strong convexity property with constant $\hat{\gamma} = \gamma - k\sqrt{s\ln(d_\nu/\delta)/n}$ with probability $1 - \delta$. This follows from observing that:

$$
H = \nabla_{\nu\nu} \sum_{i=1}^n a_i(x) \ell(z_i; \nu) = \sum_{i=1}^n a_i \nabla_{\nu\nu} \ell(z_i; \nu) \succeq \sum_{i=1}^n a_i \mathcal{H}(z_i) = \frac{1}{B} \sum_b \sum_{i \in b} a_{ib}(x) \mathcal{H}(z_i)
$$

Thus the Hessian is lower bounded by a matrix whose entries correspond to a Monte-Carlo approximation of the $U$-statistic:

$$
U = \frac{1}{\binom{n}{s}} \sum_{S \subseteq [n]:|S|=s} \frac{1}{s!} \sum_{i \in S} \mathbb{E}_\omega \left[ a_i(S, \omega) \mathcal{H}(z_i) \right] \tag{36}
$$

where $\Pi_s$ denotes the set of permutations of $s$ elements, $S_\pi$ denotes the permuted elements of $S$ according to $\pi$ and $S_\pi^1, S_\pi^2$ denotes the first and second half of the ordered elements of $S$ according to $\pi$. Finally, $a_i(S, \omega)$ denotes the tree weight assigned to point $i$ by a tree learner trained on $S$ under random seed $\omega$.

Hence, for sufficiently large $B$, by a $U$-statistic concentration inequality (Hoeffding, 1963) and a union bound, each entry will concentrate around the expected value of the $U$ statistic to within $2\sqrt{s\ln(d_\nu/\delta)/n}$, i.e.: with probability $1 - \delta$:

$$
\left\| \frac{1}{B} \sum_b \sum_{i \in b} a_{ib}(x) \mathcal{H}(z_i) - \mathbb{E}[U] \right\|_\infty \leq 2\sqrt{\frac{s\ln(d_\nu/\delta)}{n}} \tag{37}
$$

Moreover, observe that by the tower law of expectation and by honesty of the ORF trees we can write:

$$
\mathbb{E}[U] = \mathbb{E}\left[ \frac{1}{\binom{n}{s}} \sum_{S \subseteq [n]:|S|=s} \frac{1}{s!} \sum_{i \in S} a_i(S, \omega) \mathbb{E}[\mathcal{H}(z_i) \mid x_i] \right] \tag{38}
$$

Since each $\mathbb{E}[\mathcal{H}(z_i) \mid x_i]$ satisfies the restricted eigenvalue condition with constant $\gamma$, we conclude that $\mathbb{E}[U]$ also satisfies the same condition as it is a convex combination of these conditional matrices. Thus for any vector $\rho \in C(S(x); 3)$, we have w.p. $1 - \delta$:

$$
\begin{aligned}
\rho^T H \rho &\geq \rho^T \left( \sum_{i=1}^n a_i(x) \mathcal{H}(z_i) \right) \rho && \text{(lower bound on Hessian)} \\
&\geq \rho^T \mathbb{E}[U] \rho - 2\sqrt{\frac{s\ln(d_\nu/\delta)}{n}} \|\rho\|_1^2 && (U\text{-statistic matrix concentration)} \\
&\geq \gamma \|\rho\|_2^2 - 2\sqrt{\frac{s\ln(d_\nu/\delta)}{n}} \|\rho\|_1^2 && \text{(restricted strong convexity of population)} \\
&\geq \left( \gamma - 32k\sqrt{\frac{s\ln(d_\nu/\delta)}{n}} \right) \|\rho\|_2^2 && (\rho \in C(S(x); 3) \text{ and sparsity, imply: } \|\rho\|_1 \leq 4\sqrt{k}\|\rho\|_2)
\end{aligned}
$$

Since $\rho(x) \in C(S(x); 3)$ and since the weighted empirical loss satisfies a $\hat{\gamma}$ restricted strong convexity:

$$
\begin{aligned}
\sum_{i=1}^n a_i(x) \ell(z_i; \hat{\nu}(x)) - \sum_{i=1}^n a_i(x) \ell(z_i; \nu_0(x)) &\geq \sum_{i=1}^n a_i(x) \langle \nabla_\nu \ell(z_i; \nu_0(x)), \hat{\nu} - \nu_0(x) \rangle + \hat{\gamma} \|\rho(x)\|_2^2 \\
&\geq -\frac{\lambda}{2} \|\rho(x)\|_1^2 + \hat{\gamma} \|\rho(x)\|_2^2 && \text{(assumption on } \lambda)
\end{aligned}
$$

Combining with the upper bound of $\lambda \left( \|\rho(x)_{S(x)}\|_1 - \|\rho(x)_{S(x)^c}\|_1 \right)$ on the difference of the two weighted empirical losses via the chain of inequalities at the beginning of the proof, we get that:

$$\hat{\gamma}\|\rho(x)\|_2^2 \geq \frac{3\lambda}{2}\|\rho(x)_{S(x)}\|_1 - \frac{\lambda}{2}\|\rho(x)_{S(x)^c}\|_1 \leq \frac{3\lambda}{2}\|\rho(x)_{S(x)}\|_1 \leq \frac{3\lambda\sqrt{k}}{2}\|\rho(x)_{S(x)}\|_2 \leq \frac{3\lambda\sqrt{k}}{2}\|\rho(x)\|_2$$

Dividing both sides by $\|\rho(x)\|_2$ and combining with the fact that $\|\rho(x)\|_1 \leq 4\sqrt{k}\|\rho(x)\|_2$ yields the first part of the theorem.

**Bounding the gradient.**   Let $\tau = 1/(2\alpha d)$. We first upper bound the expected value of each entry of the gradient. By the shrinkage property of the ORF weights:

$$\left| \sum_{i=1}^n \mathbb{E}\left[ a_i(x)\nabla_{\nu_j}\ell(z_i; \nu_0(x)) \mid x_i \right] \right| \leq \left| \mathbb{E}\left[ \nabla_{\nu_j}\ell(z; \nu_0(x)) \mid x \right] \right| + \mathbb{E}\left[ \sum_{i=1}^n a_i(x) \left| \mathbb{E}\left[ \nabla_{\nu_j}\ell(z_i; \nu_0(x)) \mid x_i \right] - \mathbb{E}\left[ \nabla_{\nu_j}\ell(z; \nu_0(x)) \mid x \right] \right| \right]$$

$$\leq \left| \nabla_{\nu_j}L(\nu_0(x); x) \right| + L\,\mathbb{E}\left[ \sum_{i=1}^n a_i(x) \|x_i - x\| \right] \qquad \text{(Lipschitzness of } \nabla_\nu L(\nu; x))$$

$$\leq \left| \nabla_{\nu_j}L(\nu_0(x); x) \right| + L\,s^{-\tau} \qquad\qquad\qquad \text{(Kernel shrinkage)}$$

$$\leq L\,s^{-\tau} \qquad\qquad\qquad \text{(First order optimality condition of } \nu_0(x))$$

Moreover, since the quantity $\sum_i a_i(x)\nabla_{\nu_j}\ell(z_i; \nu_0(x))$ is also a Monte-Carlo approximation to an appropriately defined $U$-statistic (defined analogous to quantity $U$), for sufficiently large $B$, it will concentrate around its expectation to within $\sqrt{s\ln(1/\delta)/n}$, w.p. $1 - \delta$. Since the absolute value of its expectation is at most $Ls^{-\tau}$, we get that the absolute value of each entry w.p. $1 - \delta$ is at most $Ls^{-\tau} + \sqrt{s\ln(1/\delta)/n}$. Thus with a union bound over the $p$ entries of the gradient, we get that uniformly, w.p. $1 - \delta$ all entries have absolute values bounded within $Ls^{-\tau} + \sqrt{s\ln(d_\nu/\delta)/n}$. $\qquad\square$

# F. Omitted Proofs from Heterogeneous Treatment Effects Estimation

We now verify the moment conditions for our CATE estimation satisfies the required conditions in Assumption 4.1.

## F.1. Local Orthogonality

Recall that for any observation $Z = (T, Y, W, X)$, any parameters $\theta \in \mathbb{R}^p$, nuisance estimate $\hat{h}$ parameterized by functions $q, g$, we first consider the following *residualized* score function for PLR is defined as:

$$\psi(Z; \theta, h(X, W)) = \{Y - q(X, W) - \langle \theta, (T - g(X, W)) \rangle\} (T - g(X, W)), \tag{39}$$

with $h(X, W) = (q(X, W), g(X, W))$.

For discrete treatments, we also consider the following *doubly robust* score function, with each coordinate indexed by treatment $t$ defined as:

$$\psi^t(Z; \theta, h(X, W)) = m^t(X, W) + \frac{(Y - m^t(X, W))\,\mathbf{1}[T = t]}{g^t(X, W)} - m^0(X, W) - \frac{(Y - m^0(X, W))\,\mathbf{1}[T = 0]}{g^0(X, W)} - \theta^t \tag{40}$$

where $h(X, W) = (m(X, W), g(X, W))$.

**Lemma F.1** (Local orthogonality for residualized moments). *The moment condition with respect to the score function $\psi$ defined in (39) satisfies conditional orthogonality.*

*Proof.* We establish local orthogonality via an even stronger *conditional orthogonality*:

$$\mathbb{E}\left[ \nabla_h \psi(Z, \theta_0(x), h_0(X, W)) \mid W, x \right] = 0 \tag{41}$$

In the following, we will write $\nabla_h \psi$ to denote the gradient of $\psi$ with respect to the nuisance argument. For any $W, x$, we can write

$$\mathbb{E}\left[ \nabla_h \psi\left( Z; \theta_0(x), (q_0(x, W), g_0(x, W)) \right) \mid W, x \right] = \mathbb{E}\left[ (T - g_0(x, W), -Y + q_0(x, W) + 2\theta_0(x)^\intercal (T - g_0(x, W))) \mid W, x \right]$$

Furthermore, we have $\mathbb{E}\left[T - g_0(x, W) \mid W, x\right] = \mathbb{E}\left[\eta \mid W, x\right] = 0$ and

$$\mathbb{E}\left[-Y + q_0(x, W) + 2\theta_0(x)^\mathsf{T}\left(T - g_0(x, W)\right) \mid W, x\right] = \mathbb{E}\left[q_0(x, W) - Y + 2\theta(x)^\mathsf{T}\eta \mid W, x\right] = 0$$

where the last equality follows from that $\mathbb{E}\left[\eta \mid W, x\right] = 0$ and $\mathbb{E}\left[\langle W, q_0\rangle - Y \mid W, x\right] = 0$. $\qquad\square$

**Lemma F.2** (Local orthogonality for doubly robust moments). *The moment condition with respect to the score function $\psi$ defined in* (44) *satisfies conditional orthogonality.*

*Proof.* For every coordinate (or treatment) $t$, we have

$$\mathbb{E}\left[\nabla_g \psi^t\left(Z; \theta_0(x), (m_0(x, W), g_0(x, W))\right) \mid W, x\right]$$

$$= \mathbb{E}\left[-\frac{(Y - m_0^t(X, W))\mathbf{1}[T = t]}{(g_0^t(x, W))^2} + \frac{(Y - m_0^0(X, W))\mathbf{1}[T = t]}{(g_0^0(x, W))^2} \mid W, x\right]$$

$$= \mathbb{E}\left[-\frac{(Y - m_0^t(X, W))}{(g_0^t(x, W))^2} \mid W, x, T = t\right]\Pr[T = t \mid W, x]$$

$$+ \mathbb{E}\left[\frac{(Y - m_0^0(X, W))}{(g_0^0(x, W))^2} \mid W, x, T = 0\right]\Pr[T = 0 \mid W, x] = 0$$

and

$$\mathbb{E}\left[\nabla_m \psi^t\left(Z; \theta_0(x), (m_0(x, W), g_0(x, W))\right) \mid W, x\right]$$

$$= \mathbb{E}\left[\nabla_m\left(m_0^t(x, W) + \frac{(Y - m_0^t(x, W))\mathbf{1}[T = t]}{g_0^t(x, W)} - m_0^0(x, W) - \frac{(Y - m_0^0(x, W))\mathbf{1}[T = 0]}{g_0^0(x, W)}\right) \mid W, x\right]$$

$$= \mathbb{E}\left[\nabla_m\left(m_0^t(x, W) + \frac{(-m_0^t(x, W))\mathbf{1}[T = t]}{g_0^t(x, W)} - m_0^0(x, W) - \frac{(-m_0^0(x, W))\mathbf{1}[T = 0]}{g_0^0(x, W)}\right) \mid W, x\right]$$

$$= \nabla_m\left(\mathbb{E}\left[m_0^t(x, W) - m_0^t(x, W) - m_0^0(x, W) + m_0^0(x, W) \mid W, x\right]\right) = 0$$

This completes the proof. $\qquad\square$

### F.2. Identifiability

**Lemma F.3** (Identifiability for residualized moments.). *As long as $\mu(X, W)$ is independent of $\eta$ conditioned on $X$ and the matrix $\mathbb{E}\left[\eta\eta^\mathsf{T} \mid X = x\right]$ is invertible for any $x$, the parameter $\theta(x)$ is the unique solution to $m(x; \theta, h) = 0$.*

*Proof.* The moment conditions $m(x; \theta, h) = 0$ can be written as

$$\mathbb{E}\left[\{Y - q_0(X, W) - \theta^\mathsf{T}\left(T - g_0(X, W)\right)\}\left(T - g_0(X, W)\right) \mid X = x\right] = 0$$

The left hand side can re-written as

$$\mathbb{E}\left[\{\langle\eta, \mu_0(X, W)\rangle + \varepsilon - \langle\theta, \eta\rangle)\}\eta \mid X = x\right] = \mathbb{E}\left[\{\langle\eta, \mu_0(X, W)\rangle - \langle\theta, \eta\rangle)\}\eta \mid X = x\right]$$

$$= \mathbb{E}\left[\langle\mu_0(X, W) - \theta, \eta\rangle)\eta \mid X = x\right]$$

$$= \mathbb{E}\left[\eta\eta^\mathsf{T} \mid X = x\right]\mathbb{E}\left[\mu_0(X, W) - \theta \mid X = x\right]$$

Since the conditional expected covariance matrix $\mathbb{E}\left[\eta\eta^\mathsf{T} \mid X = x\right]$ is invertible, the expression above equals to zero only if $\mathbb{E}\left[\mu_0(X, W) - \theta \mid X = x\right] = 0$. This implies that $\theta = \mathbb{E}\left[\mu_0(X, W) \mid X = x\right] = \theta_0(x)$. $\qquad\square$

**Lemma F.4** (Identifiability for doubly robust moments.). *As long as $\mu(X, W)$ is independent of $\eta$ conditioned on $X$ and the matrix $\mathbb{E}\left[\eta\eta^\mathsf{T} \mid X = x\right]$ is invertible for any $x$, the parameter $\theta(x)$ is the unique solution to $m(x; \theta, h) = 0$.*

*Proof.* For each coordinate $t$, the moment condition can be written as

$$\mathbb{E}\left[m_0^t(X,W) + \frac{(Y - m_0^t(X,W))\,\mathbf{1}[T=t]}{g_0^t(X,W)} - m_0^0(X,W) - \frac{(Y - m_0^0(X,W))\,\mathbf{1}[T=0]}{g_0^0(X,W)} - \theta^t \mid X = x\right] = 0 \quad (42)$$

Equivalently,

$$\mathbb{E}\left[m_0^t(X,W) - m_0^0(X,W) - \theta^t \mid X = x\right] = \mathbb{E}_W\left[\mathbb{E}\left[-\frac{(Y - m_0^t(X,W))\,\mathbf{1}[T=t]}{g_0^t(X,W)} + \frac{(Y - m_0^0(X,W))\,\mathbf{1}[T=0]}{g_0^0(X,W)} \mid W, X = x\right]\right]$$

The inner expectation of the right hand side can be written as:

$$\mathbb{E}\left[-\frac{(Y - m_0^t(X,W))\,\mathbf{1}[T=t]}{g_0^t(X,W)} + \frac{(Y - m_0^0(X,W))\,\mathbf{1}[T=0]}{g_0^0(X,W)} \mid W, X = x\right] = 0$$

This means the moment condition is equivalent to

$$\mathbb{E}\left[m_0^t(X,W) - m_0^0(X,W) \mid X = x\right] = \theta^t.$$

This completes the proof. $\qquad\square$

### F.3. Smooth Signal

Now we show that the moments $m(x; \theta, h)$ are $O(1)$-Lipschitz in $x$ for any $\theta$ and $h$ under standard boundedness conditions on the parameters.

First, we consider the residualized moment function is defined as

$$\psi(Z; \theta, h(X,W)) = \{Y - q(X,W) - \theta^\mathsf{T}(T - g(X,W))\}(T - g(X,W)),$$

Then for any $\theta$ and $h$ given by functions $g$ and $q$,

$$m(x; \theta, h) = \mathbb{E}\left[\{Y - q(x,W) - \theta^\mathsf{T}(T - g(x,W))\}(T - g(x,W)) \mid X = x\right]$$

**Real-valued treatments**   In the real-valued treatment case, each coordinate $j$ of $g$ is given by a high-dimensional linear function: $g^j(x,W) = \langle W, \gamma^j \rangle$, where $\gamma^j$ is a $k$-sparse vectors in $\mathbb{R}^{d_\nu}$ with $\ell_1$ norm bounded by a constant, and $q(x,W)$ can be written as a $\langle q', \phi_2(W) \rangle$ with $q'$ is a $k^2$-sparse vector in $\mathbb{R}^{d_\nu^2}$ and $\phi_2(W)$ denotes the degree-2 polynomial feature vector of $W$.

$$m_j(x; \theta, h) = \mathbb{E}\left[\{Y - \langle q', \phi_2(W) \rangle - \theta_j(T - \langle \gamma^j, W \rangle)\}(T - g(x,W)) \mid X = x\right]$$

Note that as long as we restrict the space $\Theta$ and $H$ to satisfy $\|\theta\| \leq O(1)$, $\|\gamma\|_1, \|q'\|_1 \leq 1$, we know each coordinate $m_j$ is smooth in $x$.

**Discrete treatments with residualized moments.**   In the discrete treatment case, each coordinate $j$ of $g$ is of the form $g^j(x,W) = \mathcal{L}(\langle W, \gamma^j \rangle)$, where $\mathcal{L}(t) = 1/(1 + e^{-t})$ is the logistic function. The estimate $q$ consists of several components. First, consider function $f$ of the form $f(x,W) = \langle W, \beta \rangle$ as an estimate for the outcome of the null treament. For each $t \in \{e_1, \ldots, e_p\}$, we also have an estimate $m^t(x,W)$ for the expected counter-factual outcome function $\mu^t(x) + f(x,W)$, which takes the form of $\langle b, W \rangle$. Then the estimate $q$ is defined as:

$$q(x,W) = \sum_{t=1}^{p}(m^t(x,W) - f(x,W))g^t(x,W) + f(x,W).$$

With similar reasoning, as long as we restrict $\Theta$ and $H$ to satisfy $|\theta| \leq O(1)$, $\|\gamma^j\|_1 \leq 1$ for all $j$, and $\|\beta\|_1, \|b\|_1 \leq 1$, we know each coordinate $m_j$ is smooth in $x$.

**Discrete treatments with doubly robust moments.** Redcall that for each coordinate $t$, the moment function with input $\theta$ and nuisance parameters $m, g$ is defined as

$$\mathbb{E}\left[ m^t(x, W) + \frac{(Y - m^t(x, W))\,\mathbf{1}[T = t]}{g^t(x, W)} - m^0(x, W) - \frac{(Y - m^0(x, W))\,\mathbf{1}[T = 0]}{g^0(x, W)} - \theta^t \mid X = x \right] \tag{43}$$

where each $m^t(x, W)$ takes the form of $\langle b, W \rangle$ and each $g^t(x, W) = \mathcal{L}(\langle W, \gamma^t \rangle)$, with $\mathcal{L}$ denoting the logistic function. Then as long as we restrict the parameter space and $H$ to satisfy $\|\gamma^t\|_1$ for all $t$, then we know that $|\langle \gamma^t, W \rangle| \le O(1)$ and so $g^t(x, W) \ge \Omega(1)$. Furthermore, if we restrict the vector $b$ to satisfy $\|b\|_1 \le 1$, we know each coordinate $m_j$ is smooth in $x$.

### F.4. Curvature

Now we show that the jacobian $\nabla_\theta m(x; \theta_0(x), h_0)$ has minimum eigenvalues bounded away from 0.

**Residualized moments.** First, we consider the residualized moment function is defined as

$$\psi(Z; \theta, h(X, W)) = \{Y - q(X, W) - \theta^\mathsf{T}(T - g(X, W))\}(T - g(X, W)),$$

Then for any $\theta$ and $h$ given by functions $g$ and $q$,

$$m(x; \theta, h) = \mathbb{E}\left[\{Y - q(x, W) - \theta^\mathsf{T}(T - g(x, W))\}(T - g(x, W)) \mid X = x\right]$$

Let $J$ be the expected Jacobian $\nabla_\theta m(x; \theta_0(x), h_0)$, and we can write

$$J_{jj'} = \mathbb{E}\left[(T_j - g_0^j(x, W))(T_{j'} - g_0^{j'}(x, W)) \mid X = x\right]$$

Then for any $v \in \mathbb{R}^p$ with unit $\ell_2$ norm, we have

$$
\begin{aligned}
vJv^\mathsf{T} &= \mathbb{E}\left[ \sum_j (T_j - g_0^j(x, W))^2 v_j^2 + 2\sum_{j,j'} (T_j - g_0^j(x, W))(T_{j'} - g_0^{j'}(x, W)) v_j v_{j'} \mid X = x \right] \\
&= \mathbb{E}\left[ \left(\sum_j (T_j - g_0^j(x, W)) v_j \right)^2 \mid X = x \right] \\
&= \mathbb{E}\left[ v^\mathsf{T}(\eta \eta^\mathsf{T}) v \mid X = x \right]
\end{aligned}
$$

Then as long as the conditional expected covariance matrix $\mathbb{E}[\eta \eta^\mathsf{T} \mid X = x]$ has minimum eigenvalue bounded away from zero, we will also have $\min_v vJv^\mathsf{T}$ bounded away from zero.

**Discrete treatments with doubly robust moments.** Redcall that for each coordinate $t$, the moment function with input $\theta$ and nuisance parameters $m, g$ is defined as

$$\mathbb{E}\left[ m^t(x, W) + \frac{(Y - m^t(x, W))\,\mathbf{1}[T = t]}{g^t(x, W)} - m^0(x, W) - \frac{(Y - m^0(x, W))\,\mathbf{1}[T = 0]}{g^0(x, W)} - \theta^t \mid X = x \right] \tag{44}$$

Then $\nabla_\theta m(x; \theta_0(x), h_0) = -I$, which implies the minimum eigenvalue is 1.

### F.5. Smoothness of scores

**Residualized moments.** First, we consider the residualized moment function with each coordinate defined as

$$\psi_j(Z; \theta, h(X, W)) = \{Y - q(X, W) - \theta^\mathsf{T}(T - g(X, W))\}(T_j - g_j(X, W)),$$

Observe that for both real-valued and discrete treatments, the scales of $\theta$, $q(X, W)$, and $g(X, W)$ are bounded by $O(1)$. Thus, the smoothness condition immediately follows.

**Doubly robust moments.** For every treatment $t$,

$$\psi_t(Z; \theta, h(X, W)) = m^t(X, W) + \frac{\left(Y - m^t(X, W)\right) \mathbf{1}[T = t]}{g^t(X, W)} - m^0(X, W) - \frac{\left(Y - m^0(X, W)\right) \mathbf{1}[T = 0]}{g^0(X, W)} - \theta^t$$

Recall that each $m^t(x, W)$ takes the form of $\langle b, W \rangle$ and each $g^t(x, W) = \mathcal{L}(\langle W, \gamma^t \rangle)$, with $\mathcal{L}$ denoting the logistic function. Then as long as we restrict the parameter space and $H$ to satisfy $\|\gamma^t\|_1$ for all $t$, then we know that $|\langle \gamma^t, W \rangle| \leq O(1)$ and so $g^t(X, W) \geq \Omega(1)$. Furthermore, if we restrict the vector $b$ to satisfy $\|b\|_1 \leq 1$, we know each $m_j(X, W) \leq O(1)$. Therefore, the smoothness condition also holds.

### F.6. Accuracy for discrete treatments

For both score functions, we require that each discrete treatment (including the null treatment) is assigned with constant probability.

**Corollary F.5** (Accuracy for residualized scores). *Suppose that $\beta_0(X)$ and each coorindate $\beta_0(X), \gamma_0^j(X)$ and $\theta(X)$ are Lipschitz in $X$ and have $\ell_1$ norms bounded by $O(1)$ for any $X$. Assume that distribution of $X$ admits a density that is bounded away from zero and infinity. For any feature $X$, the conditional covariance matrices satisfy $\mathbb{E}\left[\eta\eta^\intercal \mid X\right] \succeq \Omega(1)$, $\mathbb{E}\left[WW^\intercal \mid X\right] \succeq \Omega(1)I_{d_\nu}$. Then with probability $1 - \delta$, ORF returns an estimator $\hat{\theta}$ such that*

$$\|\hat{\theta} - \theta_0\| \leq O\left(n^{\frac{-1}{2+2\alpha d}} \sqrt{\log(nd_\nu/\delta)}\right)$$

*as long as the sparsity $k \leq O\left(n^{\frac{1}{4+4\alpha d}}\right)$ and the subsampling rate of ORF $s = \Theta(n^{\alpha d/(1+\alpha d)})$. Moreover, for any $b \in \mathbb{R}^p$ with $\|b\| \leq 1$, there exists a sequence $\sigma_n = \Theta(\sqrt{\text{polylog}(n)}n^{-1/(1+\alpha d)})$ such that*

$$\sigma_n^{-1}\left\langle b, \hat{\theta} - \theta \right\rangle \to_d \mathcal{N}(0, 1),$$

*as long as the sparsity $k = o\left(n^{\frac{1}{4+4\alpha d}}\right)$ and the subsampling rate of ORF $s = \Theta(n^{\varepsilon + \alpha d/(1+\alpha d)})$ for any $\varepsilon > 0$.*

**Corollary F.6** (Accuracy for doubly robust scores). *Suppose that $\beta_0(X)$ and each coorindate $u_0^j(X), \gamma_0^j(X)$ are Lipschitz in $X$ and have $\ell_1$ norms bounded by $O(1)$ for any $X$. Assume that distribution of $X$ admits a density that is bounded away from zero and infinity. For any feature $X$, the conditional covariance matrices satisfy $\mathbb{E}\left[\eta\eta^\intercal \mid X\right] \succeq \Omega(1)$, $\mathbb{E}\left[WW^\intercal \mid X\right] \succeq \Omega(1)I_{d_\nu}$. Then with probability $1 - \delta$, ORF returns an estimator $\hat{\theta}$ such that*

$$\|\hat{\theta} - \theta_0\| \leq O\left(n^{\frac{-1}{2+2\alpha d}} \sqrt{\log(nd_\nu/\delta)}\right)$$

*as long as the sparsity $k \leq O\left(n^{\frac{1}{4+4\alpha d}}\right)$ and the subsampling rate of ORF $s = \Theta(n^{\alpha d/(1+\alpha d)})$. Moreover, for any $b \in \mathbb{R}^p$ with $\|b\| \leq 1$, there exists a sequence $\sigma_n = \Theta(\sqrt{\text{polylog}(n)}n^{-1/(1+\alpha d)})$ such that*

$$\sigma_n^{-1}\left\langle b, \hat{\theta} - \theta \right\rangle \to_d \mathcal{N}(0, 1),$$

*as long as the sparsity $k = o\left(n^{\frac{1}{4+4\alpha d}}\right)$ and the subsampling rate of ORF $s = \Theta(n^{\varepsilon + \alpha d/(1+\alpha d)})$ for any $\varepsilon > 0$.*

# G. Orange Juice Experiment

Dominick's orange juice dataset (provided by the University of Chicago Booth School of Business) contains 28,947 entries of store-level, weekly prices and sales of different brands of orange juice. The dataset also contains 15 continuous and categorical variables that encode store-level customer information such as the mean age, income, education level, etc, as well as brand information. The goal is to learn the elasticity of orange juice as a function of income (or education, etc) in the presence of high-dimensional controls.

In the experiment depicted in Figure 1, we trained the ORF using 500 trees, a minimum leaf size of 50, subsample ratio of 0.02, with Lasso models for both residualization and kernel estimation. We evaluated the resulting algorithm on 50 $\log(Income)$ points between 10.4 and 10.9. We then followed-up with 100 experiments on bootstrap samples of the original dataset to build bootstrap confidence intervals. The emerging trend in the elasticity as a function of income follows our intuition: higher income levels correspond to a more inelastic demand.

# H. All Experimental Results

We present all experimental results for the parameter choices described in Section 7. We vary the support size $k \in \{1, 5, 10, 15, 20, 25, 30\}$, the dimension $d \in \{1, 2\}$ of the feature vector $x$ and the treatment response function $\theta \in$ {piecewise linear, piecewise constant and piecewise polynomial}. We measure the bias, variance and root mean square error (RMSE) as evaluation metrics for the different estimators we considered in Section 7. In addition, we add another version of the GRF (GRF-xW) where we run the GRF R package directly on the observations, using features and controls $(x, W)$ jointly as the covariates. For the parameter space we consider, the ORF-CV and the ORF algorithms outperform the other estimators on all regimes.

## H.1. Experimental results for one-dimensional, piecewise linear $\theta_0$

Consider a piecewise linear function: $\theta_0(x) = (x + 2)\mathbb{I}_{x \leq 0.3} + (6x + 0.5)\mathbb{I}_{x > 0.3 \text{ and } x \leq 0.6} + (-3x + 5.9)\mathbb{I}_{x > 0.6}$.



Figure 6: Bias, variance and RMSE as a function of support size for $n = 5000$, $p = 500$, $d = 1$ and $\theta_0$. The solid lines represent the mean of the metrics across test points, averaged over the 100 experiments, and the filled regions depict the standard deviation, scaled down by a factor of 3 for clarity.
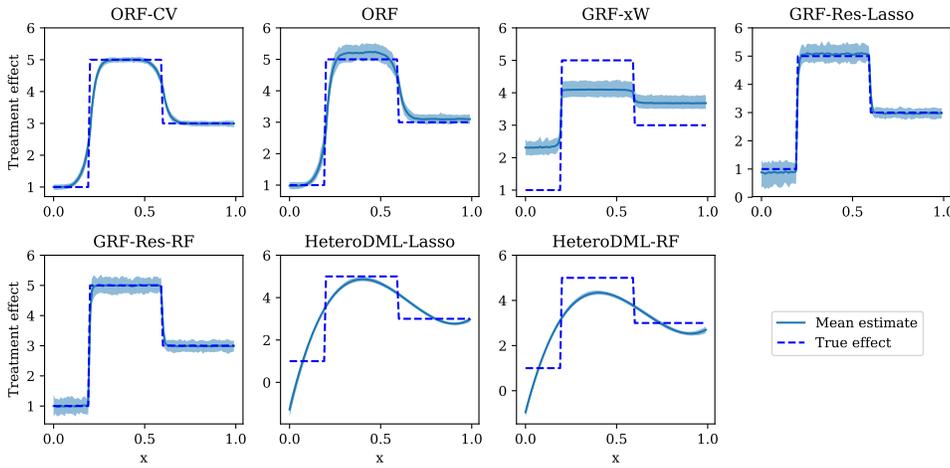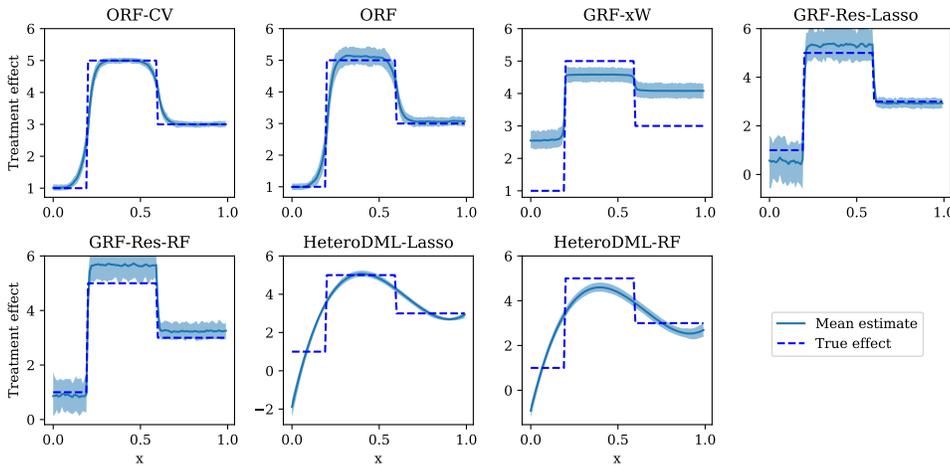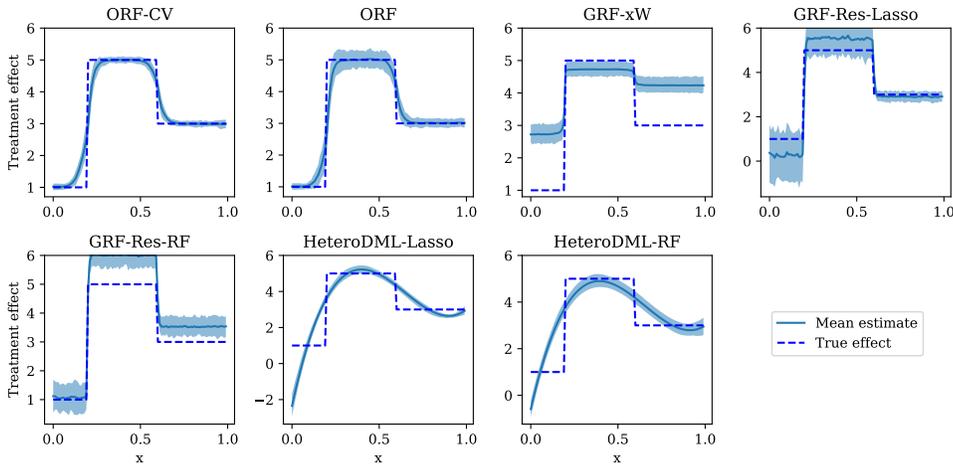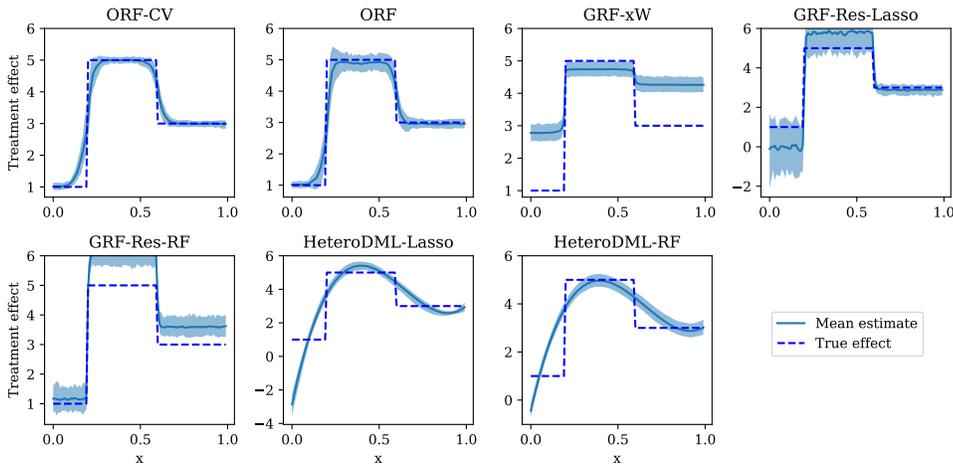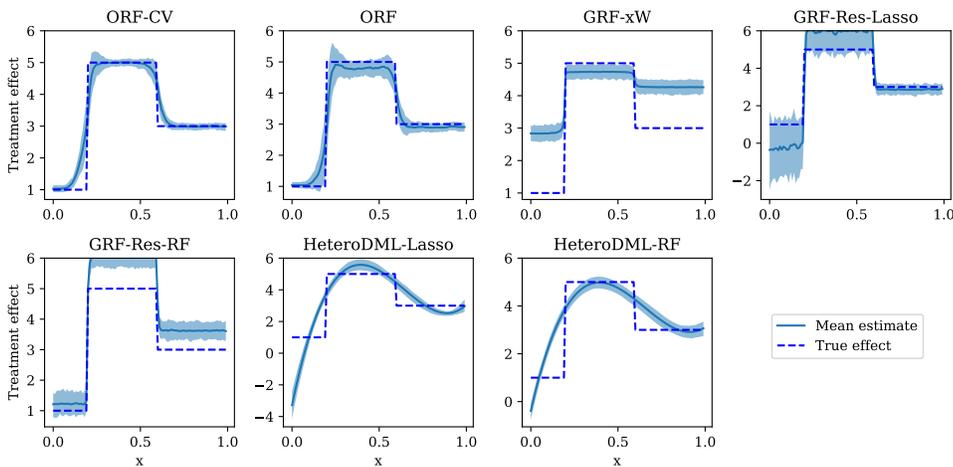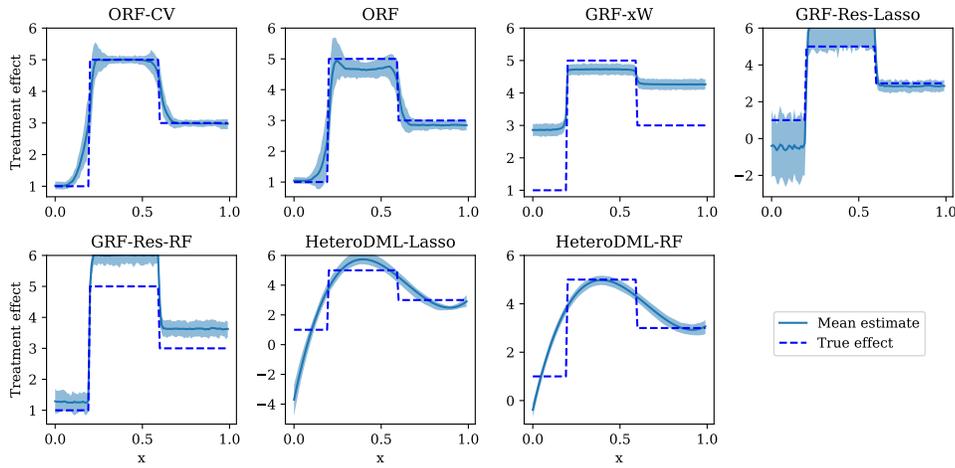


Figure 7: Treatment effect estimations for 100 Monte Carlo experiments with parameters $n = 5000$, $p = 500$, $d = 1$, $\mathbf{k = 1}$, and a piecewise linear treatment response. The shaded regions depict the mean and the 5%-95% interval of the 100 experiments.

Figure 8: Treatment effect estimations for 100 Monte Carlo experiments with parameters $n = 5000$, $p = 500$, $d = 1$, $\mathbf{k = 5}$, and a piecewise linear treatment response. The shaded regions depict the mean and the 5%-95% interval of the 100 experiments.
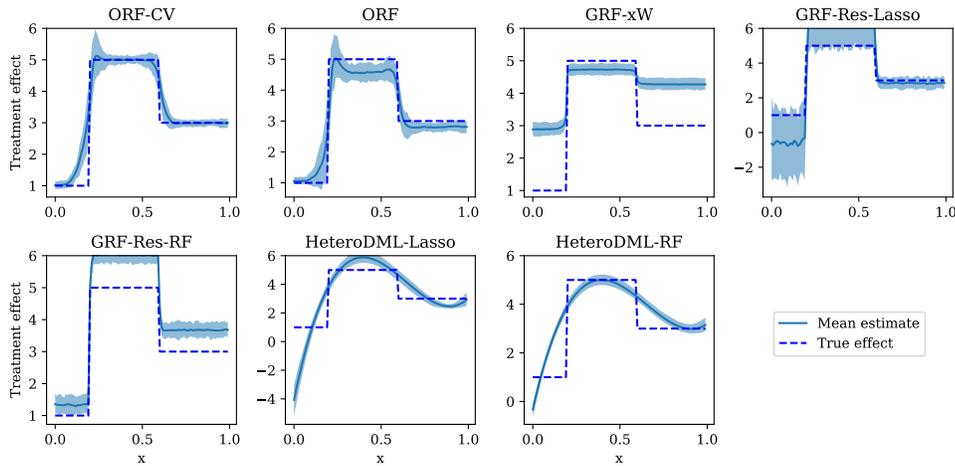


Figure 9: Treatment effect estimations for 100 Monte Carlo experiments with parameters $n = 5000$, $p = 500$, $d = 1$, $\mathbf{k = 10}$, and a piecewise linear treatment response. The shaded regions depict the mean and the 5%-95% interval of the 100 experiments.



Figure 10: Treatment effect estimations for 100 Monte Carlo experiments with parameters $n = 5000$, $p = 500$, $d = 1$, $\mathbf{k = 15}$, and a piecewise linear treatment response. The shaded regions depict the mean and the 5%-95% interval of the 100 experiments.

Figure 11: Treatment effect estimations for 100 Monte Carlo experiments with parameters $n = 5000$, $p = 500$, $d = 1$, $\mathbf{k = 20}$, and a piecewise linear treatment response. The shaded regions depict the mean and the 5%-95% interval of the 100 experiments.



Figure 12: Treatment effect estimations for 100 Monte Carlo experiments with parameters $n = 5000$, $p = 500$, $d = 1$, $\mathbf{k = 25}$, and a piecewise linear treatment response. The shaded regions depict the mean and the 5%-95% interval of the 100 experiments.



Figure 13: Treatment effect estimations for 100 Monte Carlo experiments with parameters $n = 5000$, $p = 500$, $d = 1$, $\mathbf{k = 30}$, and a piecewise linear treatment response. The shaded regions depict the mean and the 5%-95% interval of the 100 experiments.

## H.2. Experimental results for one-dimensional, piecewise constant $\theta_0$

We introduce the results for a piecewise constant function given by:

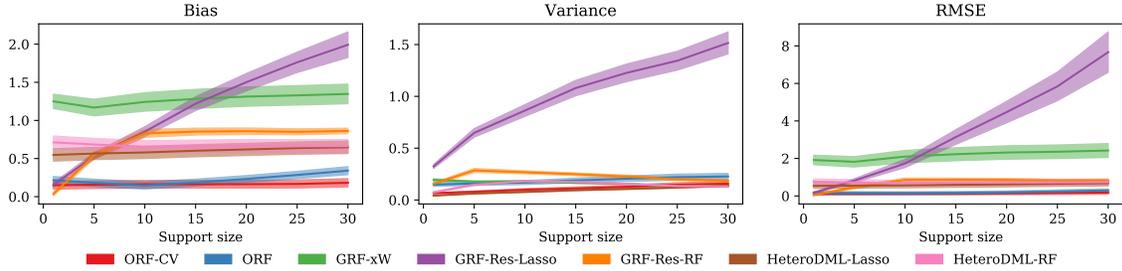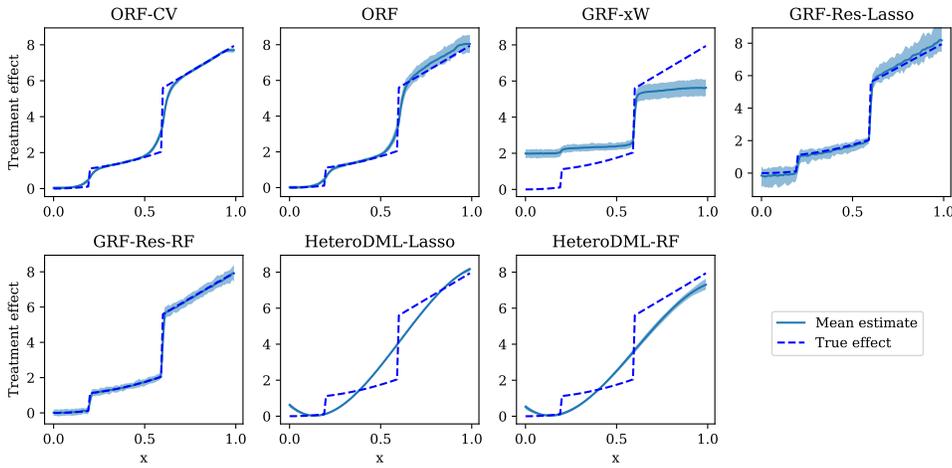$$\theta_0(x) = \mathbb{I}_{x \leq 0.2} + 5\mathbb{I}_{x > 0.2 \text{ and } x \leq 0.6} + 3\mathbb{I}_{x > 0.6}$$



Figure 14: Bias, variance and RMSE as a function of support size for $n = 5000$, $p = 500$, $d = 1$ and a piecewise constant treatment function. The solid lines represent the mean of the metrics across test points, averaged over the Monte Carlo experiments, and the filled regions depict the standard deviation, scaled down by a factor of 3 for clarity.
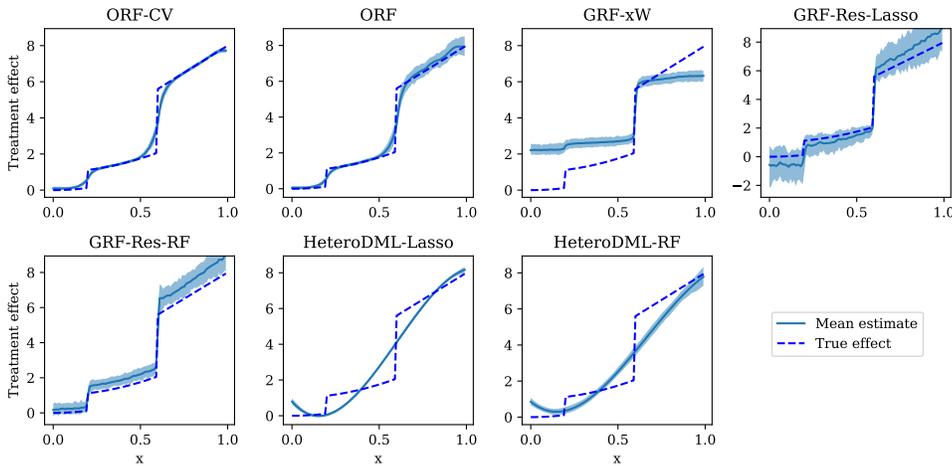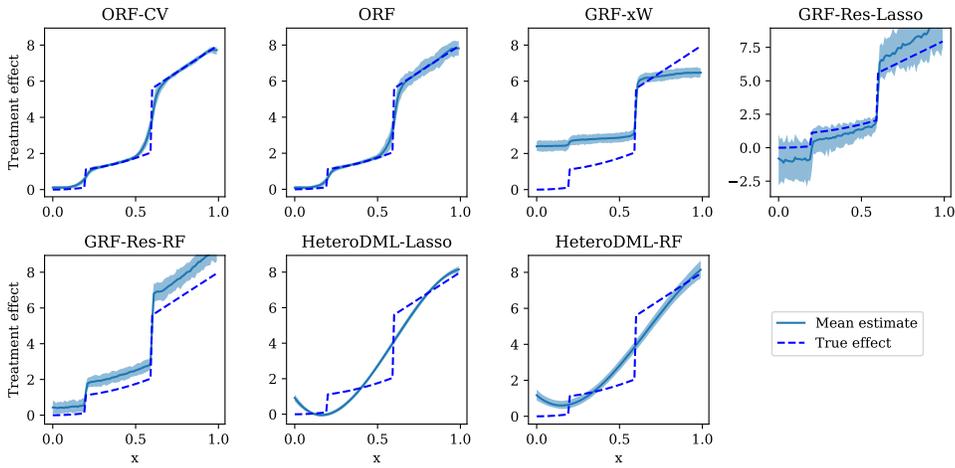


Figure 15: Treatment effect estimations for 100 Monte Carlo experiments with parameters $n = 5000$, $p = 500$, $d = 1$, $\mathbf{k = 1}$, and a piecewise constant treatment response. The shaded regions depict the mean and the 5%-95% interval of the 100 experiments.



Figure 16: Treatment effect estimations for 100 Monte Carlo experiments with parameters $n = 5000$, $p = 500$, $d = 1$, $\mathbf{k = 5}$, and a piecewise constant treatment response. The shaded regions depict the mean and the 5%-95% interval of the 100 experiments.
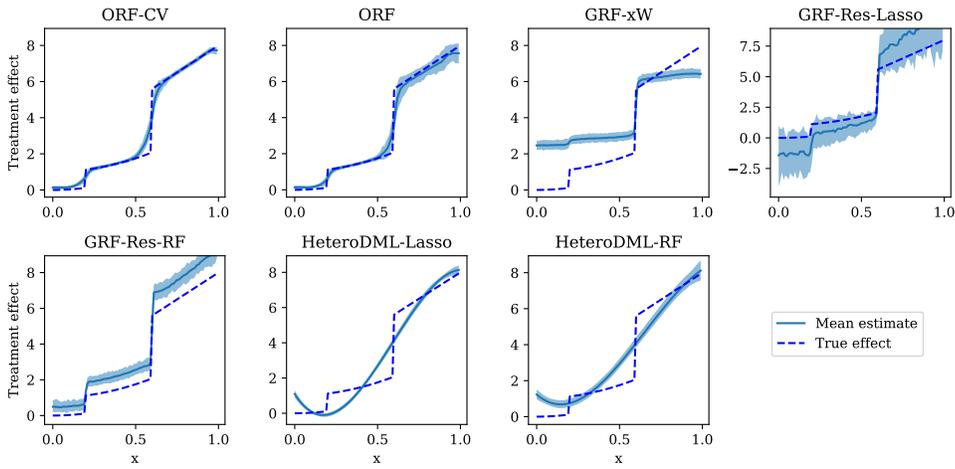
Figure 17: Treatment effect estimations for 100 Monte Carlo experiments with parameters $n = 5000$, $p = 500$, $d = 1$, $\mathbf{k = 10}$, and a piecewise constant treatment response. The shaded regions depict the mean and the 5%-95% interval of the 100 experiments.
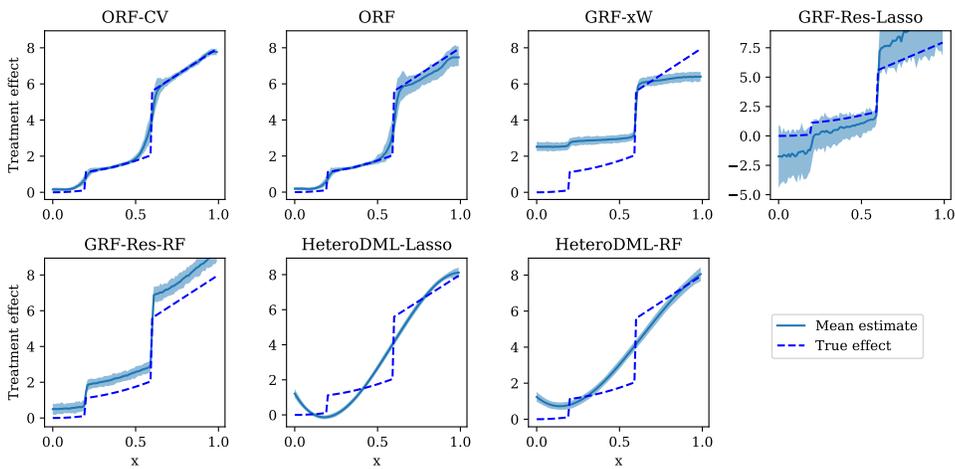
Figure 18: Treatment effect estimations for 100 Monte Carlo experiments with parameters $n = 5000$, $p = 500$, $d = 1$, $\mathbf{k = 15}$, and a piecewise constant treatment response. The shaded regions depict the mean and the 5%-95% interval of the 100 experiments.

Figure 19: Treatment effect estimations for 100 Monte Carlo experiments with parameters $n = 5000$, $p = 500$, $d = 1$, $\mathbf{k = 20}$, and a piecewise constant treatment response. The shaded regions depict the mean and the 5%-95% interval of the 100 experiments.
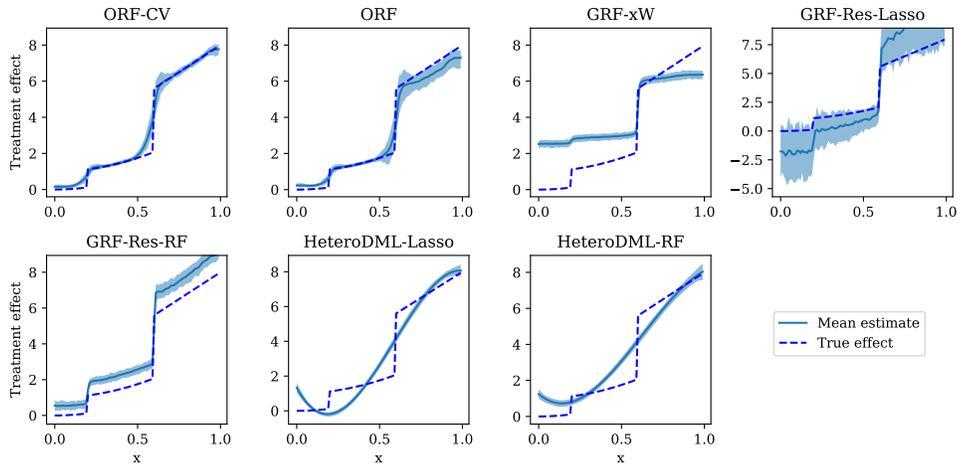
Figure 20: Treatment effect estimations for 100 Monte Carlo experiments with parameters $n = 5000$, $p = 500$, $d = 1$, $\mathbf{k = 25}$, and a piecewise constant treatment response. The shaded regions depict the mean and the 5%-95% interval of the 100 experiments.
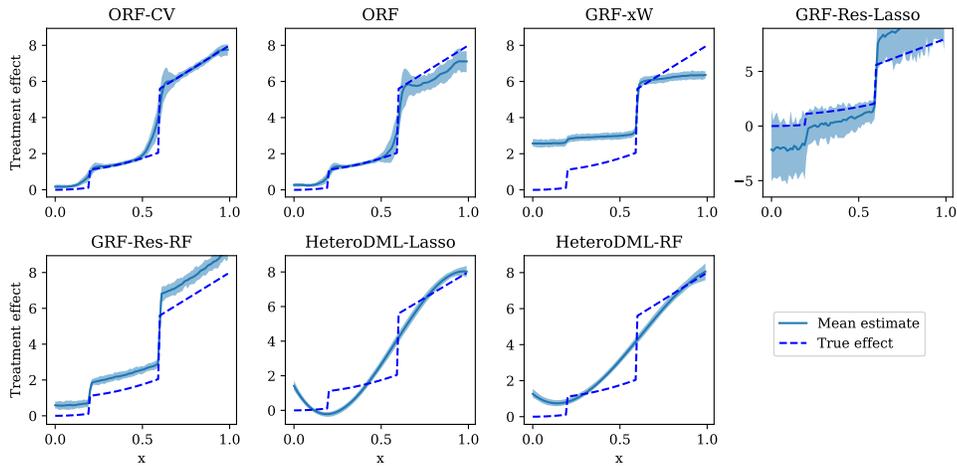


Figure 21: Treatment effect estimations for 100 Monte Carlo experiments with parameters $n = 5000$, $p = 500$, $d = 1$, $\mathbf{k = 30}$, and a piecewise constant treatment response. The shaded regions depict the mean and the 5%-95% interval of the 100 experiments.

## H.3. Experimental results for one-dimensional, piecewise polynomial $\theta_0$

We present the results for a piecewise polynomial function given by:

$$\theta_0(x) = 3x^2 \mathbb{I}_{x \leq 0.2} + (3x^2 + 1)\mathbb{I}_{x>0.2 \text{ and } x \leq 0.6} + (6x + 2)\mathbb{I}_{x>0.6}$$



Figure 22: Bias, variance and RMSE as a function of support size for $n = 5000$, $p = 500$, $d = 1$ and a piecewise polynomial treatment function. The solid lines represent the mean of the metrics across test points, averaged over the Monte Carlo experiments, and the filled regions depict the standard deviation, scaled down by a factor of 3 for clarity.
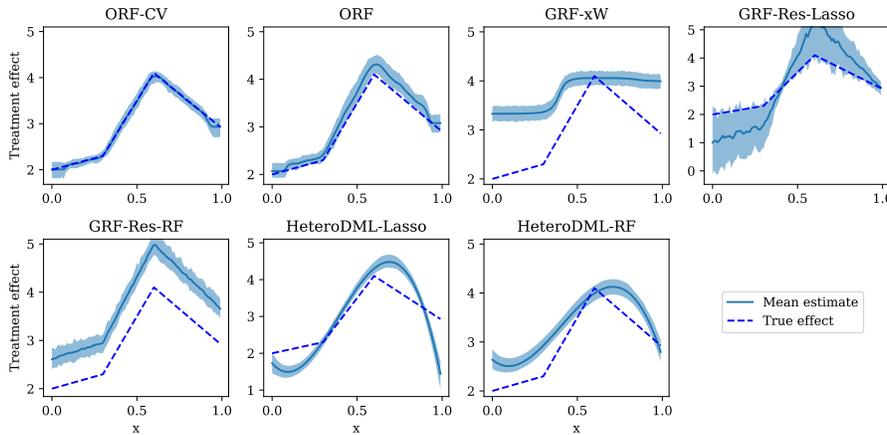


Figure 23: Treatment effect estimations for 100 Monte Carlo experiments with parameters $n = 5000$, $p = 500$, $d = 1$, $\mathbf{k = 1}$, and a piecewise polynomial treatment response. The shaded regions depict the mean and the 5%-95% interval of the 100 experiments.



Figure 24: Treatment effect estimations for 100 Monte Carlo experiments with parameters $n = 5000$, $p = 500$, $d = 1$, $\mathbf{k = 5}$, and a piecewise polynomial treatment response. The shaded regions depict the mean and the 5%-95% interval of the 100 experiments.

Figure 25: Treatment effect estimations for 100 Monte Carlo experiments with parameters $n = 5000$, $p = 500$, $d = 1$, $\mathbf{k = 10}$, and a piecewise polynomial treatment response. The shaded regions depict the mean and the 5%-95% interval of the 100 experiments.
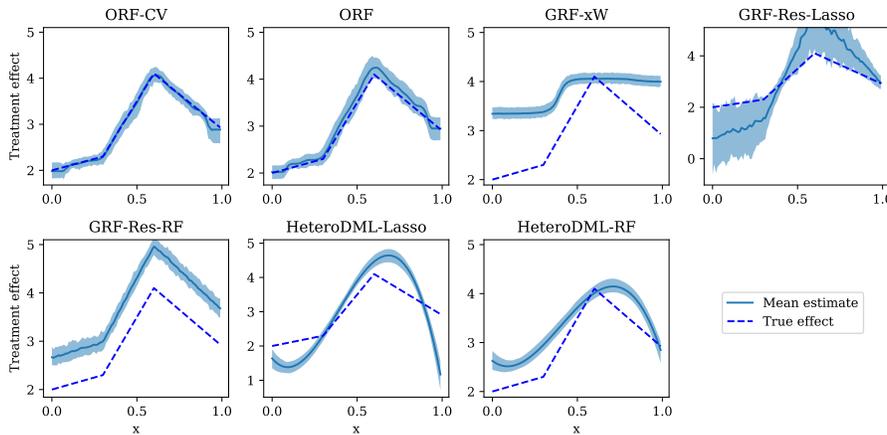


Figure 26: Treatment effect estimations for 100 Monte Carlo experiments with parameters $n = 5000$, $p = 500$, $d = 1$, $\mathbf{k = 15}$, and a piecewise polynomial treatment response. The shaded regions depict the mean and the 5%-95% interval of the 100 experiments.



Figure 27: Treatment effect estimations for 100 Monte Carlo experiments with parameters $n = 5000$, $p = 500$, $d = 1$, $\mathbf{k = 20}$, and a piecewise polynomial treatment response. The shaded regions depict the mean and the 5%-95% interval of the 100 experiments.
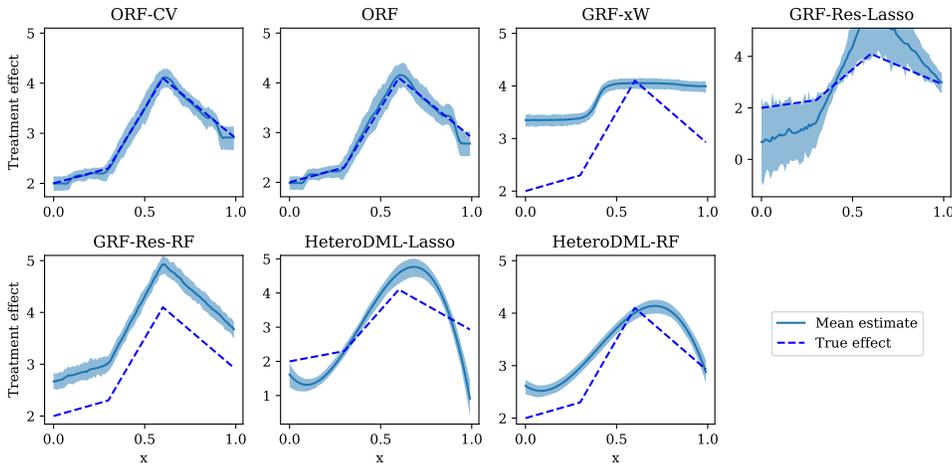
Figure 28: Treatment effect estimations for 100 Monte Carlo experiments with parameters $n = 5000$, $p = 500$, $d = 1$, $\mathbf{k = 25}$, and a piecewise polynomial treatment response. The shaded regions depict the mean and the 5%-95% interval of the 100 experiments.
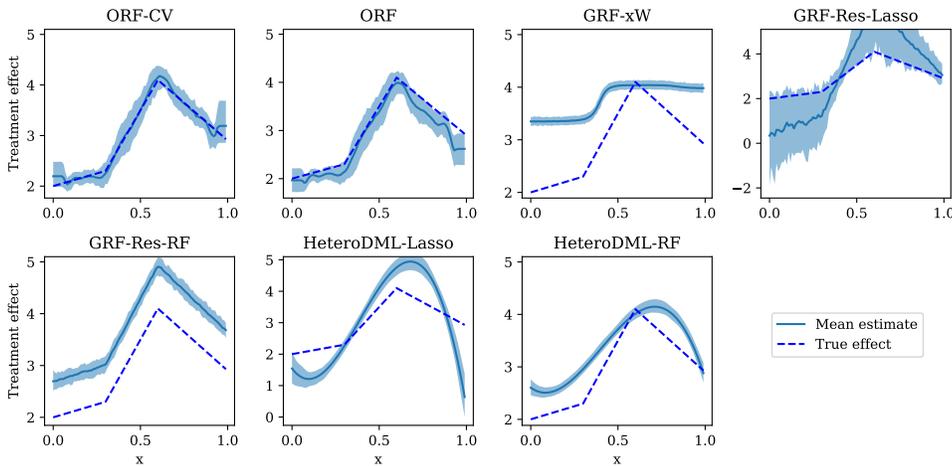


Figure 29: Treatment effect estimations for 100 Monte Carlo experiments with parameters $n = 5000$, $p = 500$, $d = 1$, $\mathbf{k = 30}$, and a piecewise polynomial treatment response. The shaded regions depict the mean and the 5%-95% interval of the 100 experiments.

## H.4. Experimental results for larger control support

We present experimental results for a piecewise linear treatment response $\theta_0$, with $n = 5000$ samples and large support $k \in \{50, 75, 100, 150, 200\}$. Figures 30-35 illustrate that the behavior of the ORF-CV algorithm, with parameters set in accordance our theoretical results, is consistent up until fairly large support sizes. Our method performs well with respect to the chosen evaluation metrics and outperform other estimators for larger support sizes.



Figure 30: Bias, variance and RMSE as a function of support size. The solid lines represent the mean of the metrics across test points, averaged over the Monte Carlo experiments, and the filled regions depict the standard deviation, scaled down by a factor of 3 for clarity.
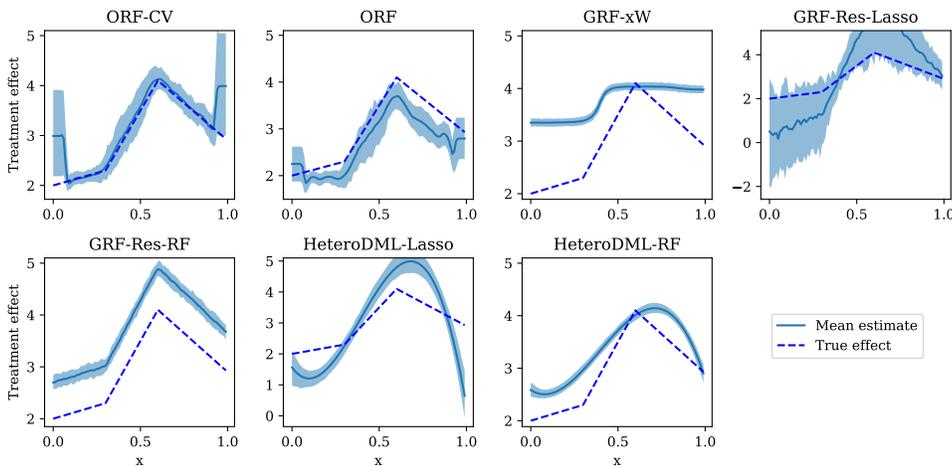


Figure 31: Treatment effect estimations for 100 Monte Carlo experiments with parameters $n = 5000$, $p = 500$, $d = 1$, $\mathbf{k = 50}$, and a piecewise linear treatment response. The shaded regions depict the mean and the $5\%$-$95\%$ interval of the 100 experiments.



Figure 32: Treatment effect estimations for 100 Monte Carlo experiments with parameters $n = 5000$, $p = 500$, $d = 1$, $\mathbf{k = 75}$, and a piecewise linear treatment response. The shaded regions depict the mean and the $5\%$-$95\%$ interval of the 100 experiments.

Figure 33: Treatment effect estimations for 100 Monte Carlo experiments with parameters $n = 5000$, $p = 500$, $d = 1$, $\mathbf{k = 100}$, and a piecewise linear treatment response. The shaded regions depict the mean and the 5%-95% interval of the 100 experiments.

Figure 34: Treatment effect estimations for 100 Monte Carlo experiments with parameters $n = 5000$, $p = 500$, $d = 1$, $\mathbf{k = 150}$, and a piecewise linear treatment response. The shaded regions depict the mean and the 5%-95% interval of the 100 experiments.

Figure 35: Treatment effect estimations for 100 Monte Carlo experiments with parameters $n = 5000$, $p = 500$, $d = 1$, $\mathbf{k = 200}$, and a piecewise linear treatment response. The shaded regions depict the mean and the 5%-95% interval of the 100 experiments.

## H.5. Experimental results for two-dimensional heterogeneity

We introduce experimental results for a two-dimensional $x$ and corresponding $\theta_0$ given by:

$$\theta_0(x_1, x_2) = \theta_{\text{piecewise linear}}(x_1)\mathbb{I}_{x_2=0} + \theta_{\text{piecewise constant}}(x_1)\mathbb{I}_{x_2=1}$$

where $x_1 \sim U[0,1]$ and $x_2 \sim Bern(0.5)$. In Figures 36-44, we examine the overall behavior of the ORF-CV and ORF estimators, as well as the behavior across the slices $x_2 = 0$ and $x_2 = 1$. We compare the performance of the ORF-CV and ORF estimators with alternative methods for $n = 5000$ and $k \in \{1, 5, 10, 15, 20, 25, 30\}$. We conclude that the ORF-CV algorithm yields a better performance for all support sizes and evaluation metrics.



Figure 36: Overall bias, variance and RMSE as a function of support size for $n = 5000$, $p = 500$, $d = 2$. The solid lines represent the mean of the metrics across test points, averaged over the Monte Carlo experiments, and the filled regions depict the standard deviation, scaled down by a factor of 3 for clarity.



Figure 37: Bias, variance and RMSE as a function of support for $n = 5000$, $p = 500$, $d = 2$ and slices $\mathbf{x_2} = \mathbf{0}$ and $\mathbf{x_2} = \mathbf{1}$, respectively. The solid lines represent the mean of the metrics across test points, averaged over the Monte Carlo experiments, and the filled regions depict the standard deviation, scaled down by a factor of 3 for clarity.
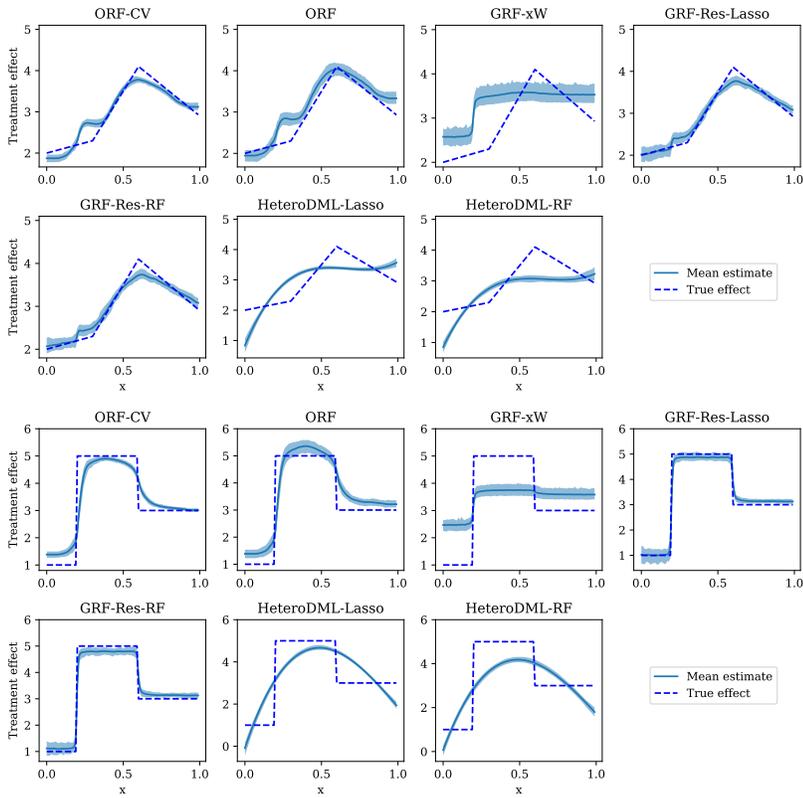
Figure 38: Treatment effect estimations for 100 Monte Carlo experiments with parameters $n = 5000$, $p = 500$, $d = 2$, $\mathbf{k = 1}$, and slices $\mathbf{x_2 = 0}$ and $\mathbf{x_2 = 1}$, respectively. The shaded regions depict the mean and the 5%-95% interval of the 100 experiments.
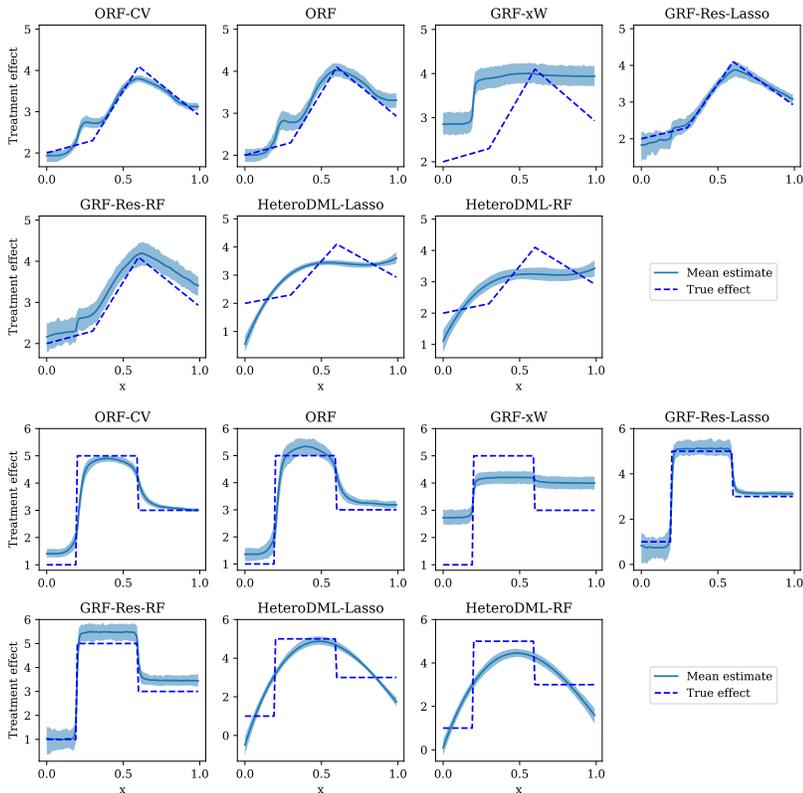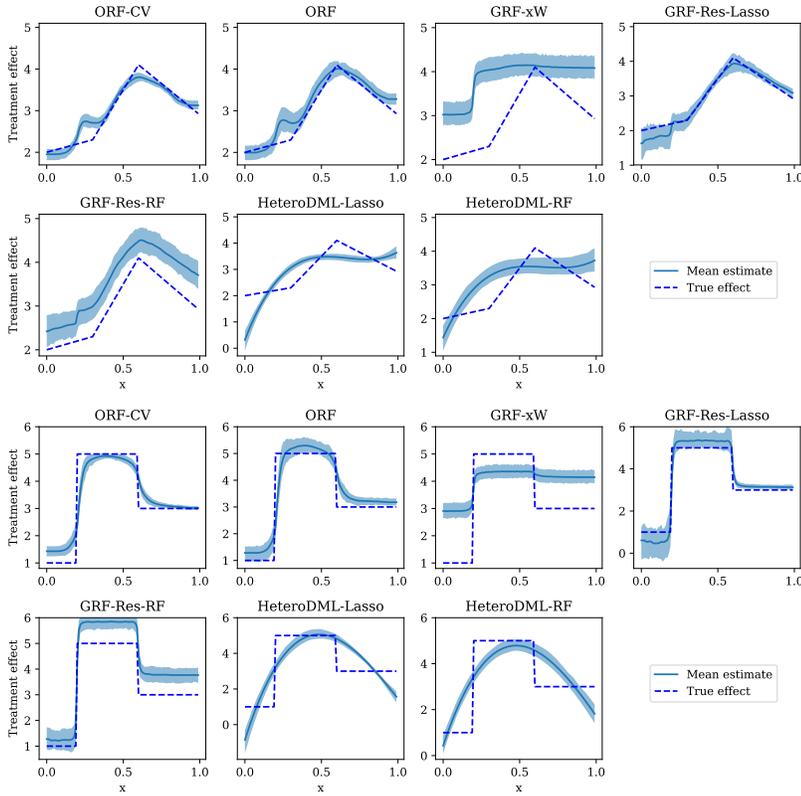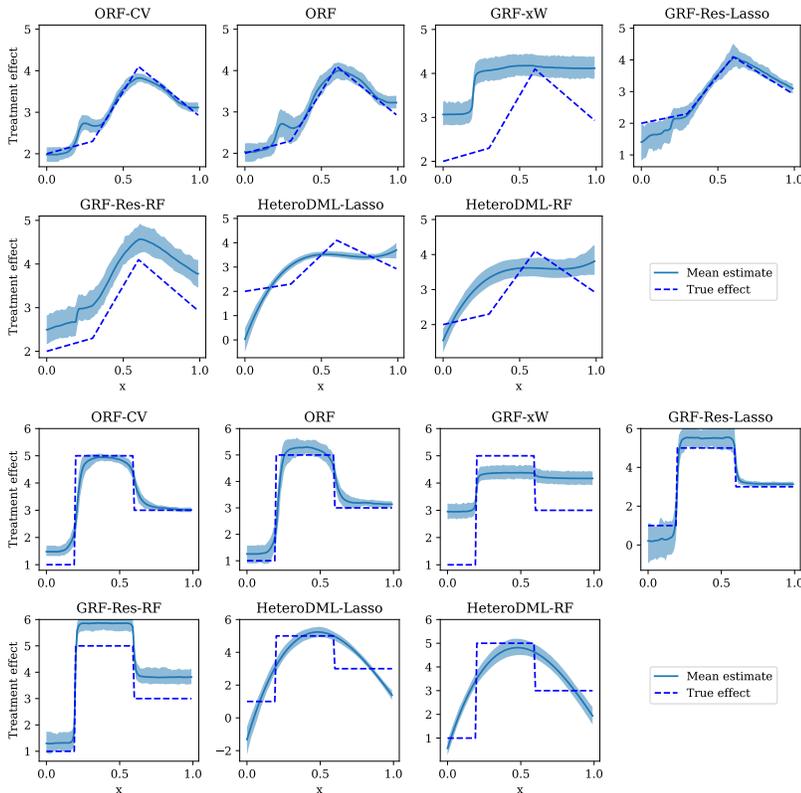


Figure 39: Treatment effect estimations for 100 Monte Carlo experiments with parameters $n = 5000$, $p = 500$, $d = 2$, $\mathbf{k = 5}$, and slices $\mathbf{x_2 = 0}$ and $\mathbf{x_2 = 1}$, respectively. The shaded regions depict the mean and the 5%-95% interval of the 100 experiments.

Figure 40: Treatment effect estimations for 100 Monte Carlo experiments with parameters $n = 5000$, $p = 500$, $d = 2$, $\mathbf{k = 10}$, and slices $\mathbf{x_2 = 0}$ and $\mathbf{x_2 = 1}$, respectively. The shaded regions depict the mean and the 5%-95% interval of the 100 experiments.
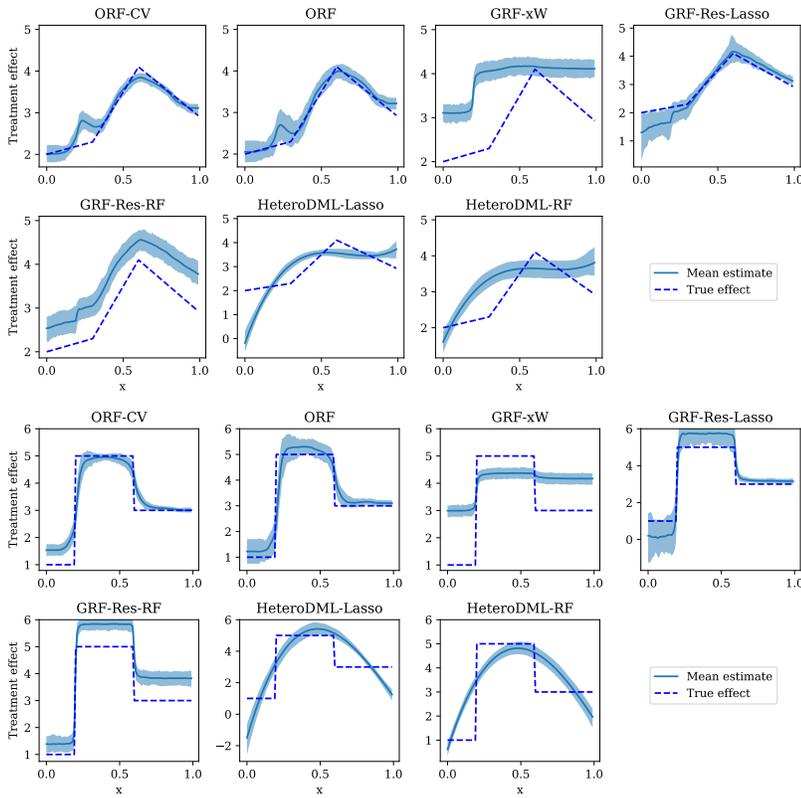
Figure 41: Treatment effect estimations for 100 Monte Carlo experiments with parameters $n = 5000$, $p = 500$, $d = 2$, $\mathbf{k = 15}$, and slices $\mathbf{x_2 = 0}$ and $\mathbf{x_2 = 1}$, respectively. The shaded regions depict the mean and the 5%-95% interval of the 100 experiments.

Figure 42: Treatment effect estimations for 100 Monte Carlo experiments with parameters $n = 5000$, $p = 500$, $d = 2$, $\mathbf{k} = \mathbf{20}$, and slices $\mathbf{x_2} = \mathbf{0}$ and $\mathbf{x_2} = \mathbf{1}$, respectively. The shaded regions depict the mean and the 5%-95% interval of the 100 experiments.
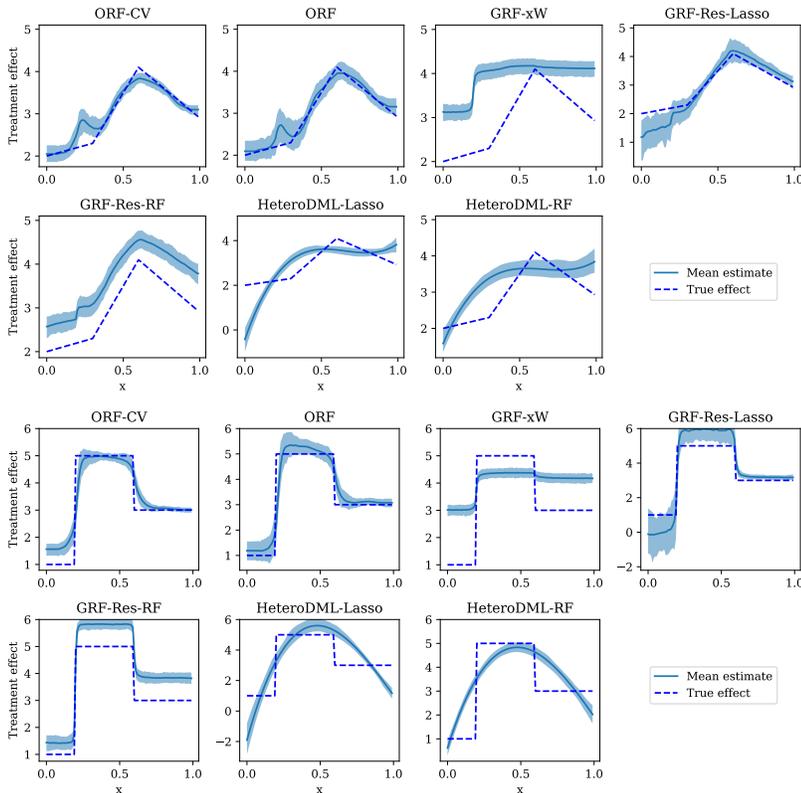
Figure 43: Treatment effect estimations for 100 Monte Carlo experiments with parameters $n = 5000$, $p = 500$, $d = 2$, $\mathbf{k} = \mathbf{25}$, and slices $\mathbf{x_2} = \mathbf{0}$ and $\mathbf{x_2} = \mathbf{1}$, respectively. The shaded regions depict the mean and the 5%-95% interval of the 100 experiments.
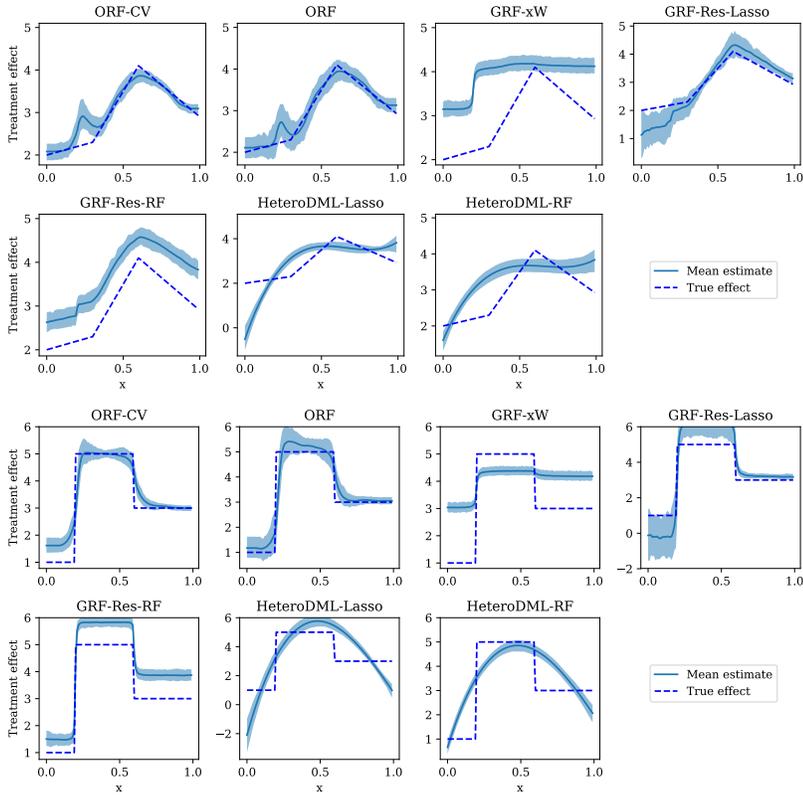
Figure 44: Treatment effect estimations for 100 Monte Carlo experiments with parameters $n = 5000$, $p = 500$, $d = 2$, $\mathbf{k = 30}$, and slices $\mathbf{x_2 = 0}$ and $\mathbf{x_2 = 1}$, respectively. The shaded regions depict the mean and the 5%-95% interval of the 100 experiments.