# Adversarial Generalized Method of Moments

**Greg Lewis** [* 1]   **Vasilis Syrgkanis** [* 2]

## Abstract

We provide an approach for learning deep neural net representations of models described via conditional moment restrictions. Conditional moment restrictions are widely used, as they are the language by which social scientists describe the assumptions they make to enable causal inference. We formulate the problem of estimating the underling model as a zero-sum game between a modeler and an adversary and apply adversarial training. Our approach is similar in nature to Generative Adversarial Networks (GAN), though here the modeler is learning a representation of a function that satisfies a continuum of moment conditions and the adversary is identifying violating moments. We outline ways of constructing effective adversaries in practice, including kernels centered by k-means clustering, and random forests. We examine the practical performance of our approach in the setting of non-parametric instrumental variable regression.

## 1. Introduction

Understanding how policy choices affect social systems requires an understanding of the causal relationships between those policies and the outcomes of interest. To measure these causal relationships, social scientists look to either field experiments, or quasi-experimental variation in observational data. Most of the observational studies rely on assumptions that can be formalized in moment conditions. This is the basis of the estimation approach known as generalized method of moments (GMM), for which Lars Hansen won a Nobel Prize.

While GMM is an incredibly flexible estimation approach,

---
[*]Equal contribution  [1]Microsoft Research and NBER  [2]Microsoft Research. Correspondence to: Greg Lewis <glewis@microsoft.com>, Vasilis Syrgkanis <vasy@microsoft.com>.

in practice its usage is confined to some special cases. One reason for this is that the underlying independence (randomization) assumptions often imply an infinite number of moment conditions. Imposing all of them is infeasible with finite data, but it is hard to know which ones to select. For some special cases, asymptotic theory provides some guidance, but it is not clear that this guidance translates well when the data is finite and/or the models are non-parametric. Given the increasing availability of data and new machine learning approaches, researchers and data scientists may want to apply adaptive non-parametric learners such as neural nets and trees to these GMM estimation problems, but this requires a way of selecting moment conditions *that are adapted to the hypothesis class of the learner.*

This paper offers a solution. We take the problem of selecting moment conditions and make it part of the overall estimation problem. Specifically, we consider environments in which underlying parameters are partially identified by a set of conditional moment restrictions. We show that the identified set is equivalently characterized as the solution to a zero-sum game between two players: a modeler and an adversary (Theorem 1). The modeler aims to minimize a GMM criterion function over a finite set of unconditional moments characterized by test functions, while the adversary tries to pick those test functions so as to maximize the criterion. We show that under some relatively weak conditions, the solution to this zero-sum game can be found by an adaptive best response process, where in each round the modeler and adversary simultaneously update their parameters and the test functions as a best response to the last round of play. In this way the adversary continually provides feedback to the modeler, similar to the approach taken in Generative Adversarial Networks (GANs). This feedback is adapted to the current weaknesses of the modeler.

Under some conditions, training can be done using stochastic gradient descent, the workhorse tool in training deep neural nets (DNNs). So one possible "deep" architecture is one in which the modeler learns a neural net while the adversary simultaneously learns the bandwidth of a set of Gaussian kernels over the space of instruments, and how to weight them. We prove that such an architecture has good generalization performance, in the sense that it identifies neural nets that with high probability do not substantially violate the moment conditions out of sample (Theorem 2).

We also suggest architectures in which the adversary learns random forests or neural nets, which may be more suitable in high dimensions.

The paper concludes with simulation experiments for the case of non-parametric instrumental variables (IV). Our estimation approach performs better than a number of leading alternatives including a two-stage polynomial basis expansion and the deep IV approach of Hartford et al. (2017).

**Related literature.** There is a large literature on GMM (Hansen, 1982), including asymptotic theory on the efficient (infeasible) instruments in the non-parametric case (Chamberlain, 1987) as well as some practical guidance on how to choose them by use of an auxiliary model (Gallant & Tauchen, 1996). There has been little work combining machine learning with generalized method of moments. A recent exception is the generalized random forest approach of Athey et al. (2016). There, a random forest is used to detect heterogeneity in treatment effects across a covariate set. Those treatment effects themselves are then solved for on each leaf using a "local" GMM estimation equation. Our paper complements their work: we show how to solve such GMM estimation problems when the both the set of available moment conditions and the hypothesis space may be large.

Moreover, our paper has strong connection to the Generative Adversarial Network (GAN) literature (Goodfellow et al., 2014; Goodfellow, 2017; Arjovsky et al., 2017), as we use adversarial training to learn a model of the world. However, unlike GANs our modeler is learning a counterfactual model rather than a distribution of data points, and the adversary is identifying conditional moment restrictions that are violated, rather than trying to distinguish between true and fake samples as in DCGANs (Goodfellow et al., 2014). Our problem is similar in spirit to Wasserstein GANs (Arjovsky et al., 2017), where the adversary is trying to identify moments/functions of the distribution of the generator, that differ from the corresponding moments of the true distribution of data. In our setting, the adversary is trying to identify conditional moment restrictions that are violated by the model of the modeler. For that reason, both the architecture of the modeler and that of the critic are inherently different from the ones used in GANs. Our work also has connections with recent work on CausalGANs for learning graphical models (Kocaoglu et al., 2017). Even though causal GANs also attempt to learn causal relationships about the world, the applications and the architecture is inherently different from our work. In particular, Kocaoglu et al. (2017) look at primarily vision tasks and learning distributions of images that satisfy a probabilistic graphical model structure. Our goal is to enable deep learning for economic applications, where the causal estimation problems are typically formulated as a set of conditional moment restrictions, rather than a probabilistic graphical model constraint.

## 2. The Conditional Moment Problem

We consider the problem of estimating a flexible econometric model $h$ that satisfies a set of conditional moment restrictions:

$$\mathbb{E}[\rho(z;h)|x] = 0 \tag{1}$$

where $x \in \mathcal{X} \subseteq \mathbb{R}^d$, $z \in \mathbb{R}^p$, $h \in \mathcal{H}$ for $\mathcal{H}$ a compact convex hypothesis space and $\rho : \mathbb{R}^p \times \mathcal{H} \to \mathbb{R}^q$. The truth is some model $h_0$ that satisfies all the moment restrictions. We will denote with $\mathcal{H}_I$ the set of all such models (the *identified set*). If $\mathcal{H}$ is small relative to the set of conditional moment restrictions implied by the model, then $h_0$ may be uniquely pinned down by (1). This is the case of *point identification*.

But since the influential work of Manski (1989), there has been increasing interest in the case where the moment restrictions only suffice for partial identification. In this case all that can be learned, even with infinite data, is the set $\mathcal{H}_I$. Our goal is to find a model/hypothesis $h \in \mathcal{H}_I$, with probability going to one as the size of the dataset goes to infinity. Since $h$ is observationally equivalent to $h_0$ — at least with respect to the moments of interest — this seems to be all one can reasonably hope for. It also sidesteps the tricky question of whether the flexible DNN representations we consider are point identified.

### 2.1. Equilibrium Formulation

We start our analysis with the observation that we can replace our conditional moment restrictions with a infinite set of unconditional moments of the form:

$$\mathcal{H}_I \equiv \{h \in \mathcal{H} : \forall f \in \mathcal{F} : \mathbb{E}[\rho(z;h)f(x)] = 0\} \tag{2}$$

where $\mathcal{F}$ is the set of functions $f : \mathcal{X} \to [0, \infty)$. It is easy to see that this is equivalent using as test functions the Dirac functions at each $x \in \mathcal{X}$.

Any hypothesis $h \in \mathcal{H}_I$ must then also be a solution of the following minimax loss minimization problem:

$$\min_{h \in \mathcal{H}} \max_{f \in \mathcal{F}} \left(\mathbb{E}[\rho(z;h)f(x)]\right)^2 \tag{3}$$

The equivalence follows from the fact that any model that is a zero for all test functions must also be so for the worst possible test function; and that the minimum of the squared loss occurs at zero. Conversely, any solution to the minimax problem must attain a zero for every test function, and so is in the identified set.

This characterization of the identified set is impractical. We would like to use the data to learn a model $h$ that has good worst-case generalization performance over test functions

$f$, but the max operator makes this objective function hard to optimize. One of the main insights of the paper is that we can instead turn this into a zero-sum game, where the minimizing player, the *modeler*, wants to choose a model $h \in \mathcal{H}$ so as to minimize the violation of the moment conditions and the maximizing player, the *adversary*, wants to choose a distribution $\sigma$ over test functions $f \in \mathcal{F}$, so as to maximize the expected violation of moments. We can find a Nash equilibrium of this game by an iterative best response algorithm, which in our case will be alternating stochastic gradient descent.

Von Neumann's classic minimax theorem provides conditions under which the Nash equilibrium is in fact a solution to the minimax problem defined by Equation (3). The main requirement is that the modeler's problem is convex and the adversary's is concave. To simplify exposition we will assume the former throughout the paper. On the other hand, the adversary's problem is unlikely to be concave in the original choice set $\mathcal{F}$, and so we allow the adversary to play a mixed strategy $\sigma$. Thus we have reduced the problem of finding an $h$ in the identified set to finding an equilibrium pair $(h^*, \sigma^*)$ of the zero-sum game.

Transforming the problem to a zero-sum game makes training more feasible. Still, the hypothesis class for the adversary is too large to generalize well or estimate from finite data. For example, how should one evaluate a Dirac delta test function at a point $x$ where no data has been observed? We solve this in the usual way: by restricting the hypothesis class to something more reasonable, with hyper-parameters that adapt to the size of the data. To achieve this we introduce the notion of a set of $\gamma$-test functions. In the algorithms that follow below, we consider an adversary that learns test functions based on kernels or decision trees or neural nets.

**Definition 1** ($\gamma$-test functions). *A set of functions $\mathcal{F}$ is a set of $\gamma$-test functions if for any conditional moment $\mathbb{E}[\rho(z; h)|x]$ there exists a function $\bar{f}$ in the convex closure $\bar{\mathcal{F}}$ of $\mathcal{F}$, such that if the violation of the unconditional moment $\mathbb{E}[\rho(z; h)\bar{f}(x)]$ is smaller than $\epsilon$, then the violation of the conditional moment is smaller than $\gamma(\epsilon)$, i.e. $\forall x \in \mathcal{X}, \exists \bar{f} \in \bar{\mathcal{F}}$:*

$$\left|\mathbb{E}[\rho(z; h)\bar{f}(x)]\right| \leq \epsilon \implies |\mathbb{E}[\rho(z; h)|x]| \leq \gamma(\epsilon) \quad (4)$$

For instance, as we show in the supplementary material, if the conditional moments are $\lambda$-Lipschitz in the conditioning variable $x$, then a simple class of test functions are all uniform Kernel test functions of the form $f(x; x_0) = 1\{\|x - x_0\|_\infty \leq h\}$ for $x_0$ being a point on a discretization of the space $\mathcal{X}$ to multiples of some small number $h$. This constitutes a set of $\gamma$-test functions with $\gamma(\epsilon) = h\lambda + \frac{\epsilon}{\mu h^d}$, for $\mu > 0$ a lower bound on the density of $x$. If we have a target $\epsilon$ in mind, we can set the optimal bandwidth $h = (\epsilon/\lambda\mu)^{1/(d+1)}$, to get $\gamma(\epsilon) = 2\lambda^{d/(d+1)} (\epsilon/\mu)^{1/(d+1)}$.

We are now ready to state our first theorem, which argues that if hypotheses that approximately satisfy all the moment conditions are close to the identified set, all $\epsilon$-equilibria of the zero-sum game defined by a class of $\gamma$-test functions have at most $\sqrt{\epsilon}$ violations of the unconditional moments (since the value of the game is the square of the maximum moment violation), at most $\gamma(\sqrt{\epsilon})$ violations of the conditional moments and so are close to the identified set:

**Theorem 1** (Approximate Set Identification). *Suppose that for any $h \in \mathcal{H}$:*

$$|\mathbb{E}[\rho(z; h)|x]| < \epsilon \; \forall x \; \Rightarrow d(h, \mathcal{H}_I) \leq \kappa(\epsilon) \quad (5)$$

*where $d(h, \mathcal{H}_I)$ denotes the distance of $h$ to $\mathcal{H}$ with respect to some metric on space $\mathcal{H}$. Then the Hausdorff distance between the set of hypotheses in the support of any $\epsilon$-equilibrium of the zero-sum game defined by a set $\mathcal{F}$ of $\gamma$-test functions and the identified set defined by Equation (1) is at most $\kappa(\gamma(\sqrt{\epsilon}))$.*

**Remark.** Even though we phrase our approach for conditional moment problems, we can also use the same approach for any over-identified unconditional moment problem, so as to perform optimal moment weighting. Suppose that instead we are given a set of $m$ moments restrictions of the form:

$$\forall i = 1, \ldots, m : \mathbb{E}[\rho_i(z; h)] = 0 \quad (6)$$

In this case there is no need to construct an $\epsilon$-cover of our test functions. Instead, the adversary will best respond by randomizing over the unconditional moment restrictions.

### 2.2. Example: Non-Parametric IV Regression

Before moving to the estimation part of our results, we offer a concrete example of a conditional moment problem and how the approximate identification theorem applies. We consider the case of IV regression where the data $z$ are a tuple $(x, w, y)$, where $x$ is an instrument, $w$ is a treatment and $y$ is an outcome. The structural model of the data that we envision takes the form:

$$y = h(w) + e \quad (7)$$
$$w = f(x) + g(e) + v \quad (8)$$

where $\mathbb{E}[e] = 0$ and $e \perp x$. Notice that $w$ is correlated with $e$ via the treatment equation, so that the treatment variable is endogenous. We have the constraint that $\mathbb{E}[e|x] = 0$. This enables causal inference about the relationship between $w$ and $y$ in the presence of the latent confounder $e$. Our goal is to estimate $h(x)$ based on this set of conditional moment restrictions: $\mathbb{E}[y - h(w)|x] = 0$. Hence, in this setting $\rho(z; h) = y - h(w)$ for $z = (y, w)$.

We investigate whether there is a metric that satisfies our conditions in Theorem 1. Let $h^* \in \mathcal{H}_I$ be a model in the

identified set. Let $h$ be a model for which all conditional moments are violated by at most $\epsilon$. Then this implies the following bound on the difference between $h$ and $h^*$:

$$\begin{aligned}
\epsilon &\geq \quad |\mathbb{E}[y - h(w)|x]| \\
&= \quad |\mathbb{E}[y - h(w)|x]| + |\mathbb{E}[y - h^*(w)|x]| \\
&\geq \quad |\mathbb{E}[h^*(w) - h(w)|x]|
\end{aligned}$$

where the second line follows since $|\mathbb{E}[y - h^*(w)|x]| = 0$ and the third line follows by the triangle inequality. Thus we can define the metric in the $\mathcal{H}$ space as: $\|h - h^*\| = \sup_x |\mathbb{E}[h^*(w) - h(w)|x]|$. If $x$ and $w$ are perfectly correlated and in one-to-one correspondence, then observe that the latter metric corresponds to simply the infinity norm: $\|h - h^*\|_\infty = \sup_w |h(w) - h^*(w)|$. The looser the instrument is correlated with the treatment, then the looser our metric will be. If $x$ is independent of $w$, then our guarantee on being close to the identified set weakens to matching the function $h^*(w)$ in expectation over $w$. For instance, predicting the unconditional mean of $y$'s would be a good such function $h(w)$. Our approach has the nice property that the finite-sample guarantees adapts to the strength of the instrument. The stronger the instrument the better the metric with respect to which we are close to the identified set and the closer we are in finding a function $h$ that is observationally equivalent to some model in the identified set.

## 3. Estimation

We now consider the estimation problem of the hypothesis $h$ from a finite set of data points $\{(z_1, x_1), \ldots, (z_n, x_n)\}$ drawn i.i.d. from the data generating distribution. We will take the standard route to estimation and replace the population problem with the finite sample approximation of it, i.e. for any function let $\mathbb{E}_n$, denote the expectation with respect to the empirical distribution of data, i.e. $\mathbb{E}_n[\rho(z; h)f(x)] = \frac{1}{n}\sum_{i=1}^n \rho(z_i; h)f(x_i)$. Moreover, for simplicity and so as to enable some of the algorithmic approaches that we will invoke in this section, we will restrict ourselves to finite classes of test functions We will then consider the zero-sum game defined by the empirical loss function: for $\theta \in \Theta$ and $\sigma \in \Delta(\mathcal{F})$:

$$L_n(h, \sigma) = \sum_{f \in \mathcal{F}_\epsilon} \sigma_f \left(\mathbb{E}_n[\rho(z; h)f(x)]\right)^2 \qquad (9)$$

To solve this game we will consider flexible parameterizations of the function $h$, denoted by $h_\theta$ for $\theta \in \Theta$, $\Theta$ finite-dimensional. We assume that the parameterization of the modeler is flexible enough such that there always exists $\theta^* \in \Theta$ such that $h_{\theta^*} \in \mathcal{H}_I$. We will then be interested in finding an equilibrium $(\theta^*, \sigma^*)$ of the zero-sum game defined by the empirical loss $L_n(\theta, \sigma)$.

We also assume the latter loss function is convex in $\theta$. For instance the latter would be true if $\rho(z; h_\theta)$ is a linear function

of $\theta$ (e.g. if $\rho(z; h)$ is linear in $h$ - as in the IV regression example - and $h_\theta(w) = \langle \theta, \phi(w) \rangle$, where $\phi(w)$ is some high-dimensional featurization of the input to the model $h$). In the experimental part we will take the implied algorithms and apply them to deep neural network representations of $h$, where the featurization $\phi(w)$ is also learned by the modeler via the first few layers of the neural network. The latter setting does not obey the convexity assumption, however training such deep representations via first order gradient descent methods is known to perform well empirically.

### 3.1. Computation: Simultaneous No-Regret Dynamics

To solve the zero-sum game, we will invoke simultaneous no-regret dynamics. Specifically we can use online projected gradient descent for the modeler and the Hedge algorithm for the critic, generating the following dynamics:

$$\begin{aligned}
\theta_{t+1} &= \Pi_\Theta\left(\theta_t - \eta_m \nabla_\theta L_n(\theta_t, \sigma_t)\right) \\
\sigma_{f,t+1} &\propto \sigma_{f,t} \cdot \exp\left\{\eta_c \left(\mathbb{E}_n[\rho(z; h_{\theta_t})f(x)]\right)^2\right\}
\end{aligned} \qquad (10)$$

where $\Pi_\Theta(\theta) = \arg\min_{\theta^* \in \Theta} \|\theta - \theta^*\|$, denotes the projection of point $\theta$ onto the space $\Theta$ of parameter and $\eta_m, \eta_c$ are hyper-parameters to be set appropriately. Then the pair of the average model $\theta^* = \frac{1}{T}\sum_{t=1}^T \theta_t$ and $\sigma^* = \frac{1}{T}\sum_{t=1}^T \sigma_t$, would correspond to an $\epsilon$-equilibrium of the game. The latter follows by standard arguments on solving zero-sum games using no-regret dynamics (Freund & Schapire, 1999) (see also the recent general formulations of this result (Shalev-Shwartz, 2012; Rakhlin & Sridharan, 2013; Syrgkanis et al., 2015)). Applying these results yields the following theorem:

**Theorem 2.** *Suppose that the loss function $L_n(\theta, \sigma)$ is convex in $\theta$ and for all $\theta \in \Theta$, $\|\theta\| \leq B$ and $\sup_{\sigma \in \Delta(\mathcal{F})} \|\nabla_\theta L_n(\theta, \sigma)\| \leq L$. Moreover, suppose that $\sup_{z,x,\theta \in \Theta, f \in \mathcal{F}} |\rho(z; h_\theta)f(x)| \leq H$. Then by setting $\eta_m = \frac{B}{L\sqrt{2T}}$ and $\eta_c = \frac{\sqrt{\log(d)}}{H^2\sqrt{2T}}$, we are guaranteed that after $T$ iterations of the training dynamics, the average solutions $\theta^* = \frac{1}{T}\sum_{t=1}^T \theta_t$ and $\sigma^* = \frac{1}{T}\sum_{t=1}^T \sigma_t$ are an $\frac{H^2\sqrt{2\log(|\mathcal{F}|)} + BL\sqrt{2}}{\sqrt{T}}$-approximate equilibrium of the zero-sum game defined by $L_n$.*

**Stochastic gradient for the modeler.** In fact the latter would also hold with high-probability if we replaced the gradients in the above equations with unbiased estimates of these gradients. Interestingly, for the loss function $L_n(\theta, \mu)$ we can compute unbiased estimates of the gradient with respect to $\theta$ by using two samples from the distribution of data $z$, since the gradient takes the form:

$$2\sum_f \sigma_f \cdot \mathbb{E}_n[\rho(z; h_\theta)f(x)] \cdot \nabla_\theta \mathbb{E}_n[\rho(z; h_\theta)f(x)]$$

Thus we see that an unbiased estimate of the gradient with respect to $\theta$ takes the form:

$$\hat{\nabla}_{\theta,t} = 2 \sum_f \sigma_f \cdot \rho(z; h_\theta) f(x) f(\tilde{x}) \nabla_\theta \rho(\tilde{z}; h_\theta)$$

where $z$ and $\tilde{z}$ are two independent random samples from the data. The latter unbiased estimate only requires two samples from the empirical distribution and can significantly speed up computation during training time. Similarly, for variance reduction, instead of two independent samples, we can also use two batches of independent samples of some mini-batch size $B$.

**Bounded gradient bias for the critic.** Unfortunately, the update of the critic does not admit a stochastic version since the expectation over the data $z$ is under the square. Hence, an update for the critic requires calculating the moment violation over a large sample of data, so that concentration inequalities kick in and the bias in the gradient is of negligible size. Typically that would be of order $1/\sqrt{B}$ if we use a data set of $B$ samples. Alternatively, we can use the whole set of samples for the updates of the critic. Fortunately, the update of the critic does not involve calculating gradients of a deep neural net, rather it only requires evaluations of a neural net over a large sample set. If the latter is also costly to implement, then we can interleave multiple steps of the modeler in between one step of the critic.

### 3.2. Sample Complexity and Generalization Error

The dynamics of the previous section are guaranteed to find an equilibrium of the game defined by the empirical distribution loss $L_n$. However, we are interested in connecting an equilibrium of the empirical game with the population game and showing that any equilibrium of the empirical game will also satisfy the population moments of Equation (1) to within some vanishing error. We require some conditions on the sample complexity of the hypothesis space $\mathcal{H}_\Theta = \{h_\theta : \theta \in \Theta\}$. Recall the Rademacher complexity of a function class $\mathcal{G}$ and sample $S = \{z_1, \ldots, z_n\}$:

$$R(G, S) = \mathbb{E}_\xi \left[ \sup_{g \in \mathcal{G}} \frac{1}{n} \sum_{i=1}^{n} \xi_i \cdot g(z_i) \right] \quad (11)$$

where $\xi_i$ are Rademacher random variables which take values $\{-1, 1\}$ with equal probability.

**Theorem 3.** *Suppose that the Rademcaher complexity of the function class $\mathcal{A} = \{\rho(\cdot; h_\theta) f(\cdot) : \theta \in \Theta, f \in \mathcal{F}\}$ is upper bounded by $\mathcal{R}$ and that the conditions of Theorem 2 hold. Then with probability $1 - \delta$, the averages $\theta^* = \frac{1}{T} \sum_{t=1}^{T} \theta_t$ and $w^* = \frac{1}{T} \sum_{t=1}^{T} w_t$, correspond to an $\epsilon$-equilibrium of the population game defined by Equation (3) for $\epsilon = O\left( \frac{H^2 \sqrt{2 \log(|\mathcal{F}|)} + BL\sqrt{2}}{\sqrt{T}} + H\mathcal{R} + H^2 \sqrt{\frac{\log(1/\delta)}{n}} \right)$. Therefore: $d(h_{\theta^*}, \mathcal{H}_I) \leq \kappa(\gamma(\sqrt{\epsilon}))$.*

For the case of neural networks, classical results of Anthony & Bartlett (2009) bound their Rademacher complexity as a function of the number of gates and the number of connections used. Moreover, as we show in the supplementary material for the case of Lipschitz moments, applying the latter theorem to the class of uniform Kernel test functions centered at a uniform grid of appropriate size, can lead to convergence rates of the order of $n^{-1/(2(d+1))} \times \max(\sqrt{d \log(d)}, \sqrt{r \log(r)}, \log(1/\delta))$, where $r$ is the number of parameters of the modeler, with probability $1 - \delta$. We remark that an alternative route to connecting the empirical equilibrium with the population one, is to use the fact that our training process is stable in the sense that it employs gradient descent and the hedge algorithm for each player. Both of these algorithms are special cases of the Follow-the-Regularized-Leader algorithms, which are known to be algorithms that are stable with respect to each data point. We defer a more refined analysis of our theoretical results via stability to future work.

## 4. The Algorithm: Adversarial GMM

Given the theory laid out in the previous sections we are now ready to state the main algorithm that we deployed in practice. In the previous section, we developed formal rates of convergence of the model estimated by our simultaneous training dynamics, as a function of the class of test functions $\mathcal{F}$ used. We also instantiated this class of test functions for the case of Lipschitz moments for the class of uniform kernels centered on a grid of points and derived worst-case concrete rates of convergence. However, in practice, taking a grid of points might be a computationally very heavy approach and might also not use the structure of the data. For this reason we propose that instead of discretizing the whole space of $\mathcal{X}$ we rather cluster the data points into small clusters based on some criterion and then place local kernels around these clusters.

The clustering approach that worked well in the experiments was to construct a K-means clustering of the data points $x_1, \ldots, x_n$ (i.e. the conditioning variables). Then consider the set of functions that correspond to Gaussian kernels around each centroid of the cluster and with standard deviation large enough to encompass the points in the cluster.[1] This gives rise to our main Algorithm 1, with Subprocess 2.

**Alternative Designs of $\mathcal{F}$.** We propose several alternative approaches to constructing the set of test functions $\mathcal{F}$ that might perform better dependend on the structure of the moment problem and the parameterization of the hypothesis

---

[1]Alternatively, a uniform kernel that puts equal weight on the points of each cluster and zero weight outside can also be used. We found experimentally that the hard boundaries of the uniform kernel lead to poorer performance.

**Algorithm 1** AdversarialGMM

1: **Hyperparameters**: Step sizes $\eta_m, \eta_c$, mini-batch sizes $B_m, B_c$, number of models $M$ to use for averaging, a test function generator $G$.
2: **Input**: data $S = \{(z_1, x_1), \ldots, (z_n, x_n)\}$
3: Generate set of test functions $\mathcal{F} = G(S)$
4: **for** $t = 1, \ldots, T$ **do**
5:     Randomly sample with replacement from $S$ three batches of samples $S_1, S_2, S_3$, with sizes $B_m, B_m, B_c$ corrspondingly.
6:     Let $E_S$ denote an expectation with respect to the empirical distribution in a sample $S$. Construct estimates of the gradient of the modeler:

$$\hat{\nabla}_{\theta,t} = 2 \sum_{f \in \mathcal{F}} \sigma_{f,t} \cdot \mathbb{E}_{S_1} \left[\rho(z; h_{\theta_t}) f(x)\right] \cdot$$
$$\mathbb{E}_{S_2} \left[f(x) \nabla_{\theta} \rho(z; h_{\theta})\right]$$

    and the utility of the critic:

$$\hat{U}_{f,t} = \left(\mathbb{E}_{S_3} \left[\rho(z; h_{\theta_t}) f(x)\right]\right)^2, \ \forall f \in \mathcal{F}$$

7:     **Modeler Step.** Take a gradient step for the modeler using projected gradient descent (or any first order algorithm like Adam):

$$\theta_{t+1} = \Pi_{\Theta} \left(\theta_t - \eta_m \hat{\nabla}_{\theta,t}\right)$$

8:     **Critic Step.** Take a Hedge step for the critic:

$$\sigma_{f,t+1} \propto \sigma_{f,t} \cdot \exp\left\{\eta_c \hat{U}_{f,t}\right\}, \ \forall f \in \mathcal{F}$$

9:     **Critic Jitter.** If each $f \in \mathcal{F}$ is parameterized via some parameter $w_f \in W$, then take a gradient step for $w_f$:

$$\nabla_{w_f,t} = 2\mathbb{E}_{S_1} \left[\rho(z; h_{\theta_t}) f(x; w_f)\right] \cdot$$
$$\mathbb{E}_{S_2} \left[\rho(z; h_{\theta}) \nabla_{w_f} f(x; w_f)\right]$$
$$w_{f,t+1} = \Pi_W \left(w_{f,t} + \eta_w \nabla_{w_f,t}\right)$$

10: **end for**
11: Return the average model $h^* = \frac{1}{M} \sum_{t \in I} h_t$, where $I$ is a set of $M$ randomly chosen steps during training.

**Algorithm 2** KMeans based test function generation

1: **Hyperparameters**: Number of clusters $K$, minimum radius size $r$.
2: **Input**: data $S = \{(z_1, x_1), \ldots, (z_n, x_n)\}$
3: Cluster the data-points $(x_1, \ldots, x_n)$ into $K$ clusters using $K$-means.
4: For each cluster $i \in [K]$ let $f_i : \mathcal{X} \to R$ be a Gaussian Kernel:

$$f_i(x) = \frac{1}{(2\pi\sigma^2)^{1/2d}} \exp\left\{-\frac{\|x - x_i\|_W^2}{2\sigma^2}\right\} \quad (12)$$

with standard deviation $\sigma$ equals to twice the distance to the $r$-th closest data point to the centroid of cluster $i$, with respect to a $W$-matrix weight norm:

$$\|x\|_W^2 = x^T W x \quad (13)$$

for some positive definite matrix $W = VV^T$.
5: Return $\mathcal{F} = \{f_i : i \in [K]\}$

space:

- **Random Data-Points.** Instead of constructing a $K$-means clustering, one could simply take $K$ random data points as centroids and then construct gaussian or other kernels around these data points so as to cover a minimum sized neighborhood of $r$ other points

- **Random Forest.** To take advantage of the structure of the problem we build a random forest by regressing the variables $z$ on the variables $x$. This will essentially uncover the dimensions of $x$ that really have any effect on $z$ and hence are useful in identifying the model $h$. Subsequently, we construct a test function for each leaf of the random forest which corresponds to a local kernel at the center of the leaf and covering the data-points of the leaf. One could either use a Gaussian kernel constructed based on the minimum enclosing ellipsoid of the points in the leaf or a uniform kernel on the leaf points.

- **Sieve-based test functions.** An alternative to the local test functions, is to construct test functions that correspond to a sieve approximation of any function. One such sieve is the polynomial sieve where the test functions are simply $\mathcal{F} = \{x^d : d \in \{0, \ldots, K\}\}$ for some upper bound degree $K$ (potentially dividing by a normalizing constant so that each test function takes values in the same range). The latter corresponds to the standard approach of semi-parametric estimation in the econometrics literature (see e.g. the recent work on two-step sieve GMM estimation by (Chen & Liao, 2015))

- **Neural net based test functions.** Yet another alternative is to be fully flexible in the creation of the test functions, by representing them as neural networks $f(\cdot; w_1), \ldots, f(\cdot; w_m)$. We want these different test functions to provide a good cover of the test function space. Hence, we can create a secondary game between the neural network critics, e.g. by adding a term in the loss of the critics that represents some metric of similarity with other critics. The latter would incentivize the critics to approximate different combinations of moments. Subsequently a meta-critic will optimize the weight placed on each critic using Hedge. The latter is closer in spirit to Wasserstein GANs (Arjovsky et al., 2017), which also approximate the class of all 1-Lipschitz functions with a neural net representation.

## 5. Experimental Evaluation

We applied our AdversarialGMM algorithm to the problem of non-parametric instrumental variable regression from Section 2.2.

**Data Generating Processes.** In each experiment we generated 1000 data points. We analyzed the experimental performance for the following two data generating processes:

**DGP** 1:

$$
\begin{aligned}
w &= (1 - \gamma)x_1 + \gamma e + \zeta \\
y &= h_0(w) + e + \delta \\
e &\sim \mathcal{N}(0, 2), \quad x \sim \mathcal{N}(0, 2I_d) \\
\zeta &\sim \mathcal{N}(0, 0.1), \quad \delta \sim \mathcal{N}(0, 0.1)
\end{aligned}
$$

**DGP** 2:

$$
\begin{aligned}
w &= (1 - \gamma)\left(x_1 \cdot 1\{x_1 > 0\} + x_2 \cdot 1\{x_2 < 0\}\right) + \gamma e + \zeta \\
y &= h_0(w) + e + \delta \\
e &\sim \mathcal{N}(0, 2), \quad x \sim \mathcal{N}(0, 2I_d) \\
\zeta &\sim \mathcal{N}(0, 0.1), \quad \delta \sim \mathcal{N}(0, 0.1)
\end{aligned}
$$

In other words, we generate $d$ instruments drawn from independent normal distributions. However, in DGP 1 only the first instrument has any effect on the treatment, while in DGP 2 only the first and the second instruments have any effect on the treatment. In DGP 1 the effect of the first instrument is linear. In DGP 2, the first instrument affects the treatment when it is positive and the second instrument when it is negative. The constant $\gamma$ captures the *strength of the instrument*, since it controls the relative weight between the instrument and the confounding factor $e$. The goal of the multiple instruments is to see if the adversarial GMM algorithm will identify the relevant instruments through learning the appropriate weighting matrix $W = VV^T$ of the Kernel

norm. In particular, whether the matrix $V$ will project to the dimensions of the relevant instruments and ignore the rest.

We considered different types of functions for the true model:

- 2-d polynomial (2dpoly): $h_0(w) = -1.5w + .9w^2$

- 3-d polynomial (3dpoly): $h_0(w) = -1.5w + .9w^2 + w^3$

- Absolute value function (abs): $h_0(w) = |w|$

- Identity function (linear): $h_0(w) = w$

- Sigmoid function (sigmoid): $h_0(w) = \frac{2}{1+e^{-2x}}$

- Sin function (sin): $h_0(w) = \sin(x)$

- Step function (step): $h_0(w) = 1 + 1.5 \cdot 1\{x \geq 0\}$

- Random piece-wise linear function (rand_pw): $h_0(w) = \sum_{i=1}^{5} (a_i w + b_i) 1\{x_i \in [\tau_{i-1}, \tau_i]\}$, where $\tau_0 = -2$ and $\tau_i$ are chosen uniformly at random over a grid of step 0.1 in $[-2, 2]$, $a_i$ are chosen uniformly at random in $[-4, 4]$, $b_1$ is chosen uniformly at random in $[-1, 1]$ and $b_i$ for $i > 1$ are chosen appropriately so that the whole function is continuous.

**Training Hyper-Parameters.** We evaluated each of these functions with the same training process without further fine-tuning so that we can argue about their adaptivity to different types of true models. We applied the adversarial GMM algorithm with $T = 400$ steps of the training dynamics, where the modeler was trained with a batch-size of $B = 100$ randomly sampled data points (with replacement) per gradient step, while the critic's loss was evaluated on the whole sample set of $n = 1000$ points. We updated the critic after every iteration of the modeler and we also updated the weight matrix of the Kernel of the critic functions after every iteration. For the matrix $W = VV^T$, we used an $2 \times d$ matrix $V$, i.e. we projected the $d$ instruments into a two-dimensional space that was learned via a gradient method (unless $d = 1$ in which case we used $W = I$). Figure 1 gives an example of the gaussian Kernels generated by the algorithm. For stability and improved performance, we used the ADAM algorithm for all the gradient optimizations (modeler and critic jitter), rather than stochastic gradient descent.

The modeler was represented as deep neural net with three hidden layers, each with width of $r = 1000$ RELU (rectified linear units) gates. The critic was randomizing over the set of gaussian kernel functions centered on the centroids of a $K$-means clustering with $K = 50$ performed on the $d$-dimensional $\mathcal{X}$ space. The learning rate of the modeler and critic were $\eta_m = \eta_w = 0.007 \approx 1/\sqrt{rT}$ and $\eta_c = 0.11 \approx 1/\sqrt{\log(K)T}$, correspondingly.
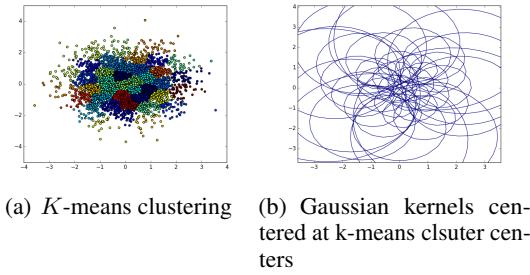
(a) $K$-means clustering    (b) Gaussian kernels centered at k-means cluster centers

*Figure 1.* $K$-means clustering of the two-dimensional space $\mathcal{X}$.

**Benchmarks.** We compared our approach to a polynomial two-stage least squares (2SLSPoly) approach, where we represented $h_\theta(w)$ as a third degree polynomial and then estimated $w, w^2, w^3$ in separate first stage regressions on all interactions of $x_1$ and $x_2$ up to degree 3, regularizing to remove instability. We also compared it to running a direct polynomial regression (DirectPoly), direct neural net regression (DirectNN), standard linear two-stage least squares (2SLS), and DeepIV.

Figure 2 shows qualitative examples of fitted functions using each of these methods and for each true functional model $h_0$. We show three different outputs from our DNN: the average across all steps of the dynamics (avg), the final model (final), and the model that performed best in-sample moment violation (best).

**Performance Metrics.** We repeated each experiment $M = 100$ times on fresh samples of the data and we calculated the $R^2$ of each method in each of the runs on a set of test treatment points. More concretely, we generated a set of test points $\{w_1, \ldots, w_K\}$ and evaluated, the mean squared error of each estimator $\hat{h}$

$$\text{MSE}(\hat{h}) = \frac{1}{K} \sum_{i=1}^{K} \left( h_0(w_i) - \hat{h}(w_i) \right)^2 \qquad (14)$$

and the R-squared, which captures how much of the variance of the function was explained by our model:

$$R^2 = 1 - \frac{\text{MSE}(\hat{h})}{\frac{1}{K} \sum_{i=1}^{K} \left( h_0(w_i) - \frac{1}{K} \sum_{t=1}^{K} h_0(w_t) \right)^2} \qquad (15)$$

We created two types of test treatment points: i) in the first case we generated treatment points from the marginal distribution of the original data generating process, ii) in the second case we created a uniform grid of 100 points between the 10 and 90 percentile of the treatment variable in the original set of sample points on which we trained our model.

In Figures 3 and 4, we give the median $R^2$ and the $5 - 95$ percentiles of the $R^2$ across the $M = 100$ experiments for each of the two DGPs. Moreover, in Figure 5, shows how the $R^2$ performance of each method varies as the number of instruments $d$ increases. For this figure we merged together the experiments from all the different true models $h_0$ (except rand_pw) and computed median $R^2$ across the $M = 700$ total experiments. Finally in Figure 6 we plot the median $R^2$ and the $10 - 90$ percentiles specifically for the random true model rand_pw as the strength and dimension of the instruments varies.

In the Appendix we provide more detailed performance for each true model separately as well as a plethora of more detailed experimental results and figures.

**Conclusion.** We find that the Adversarial GMM algorithm is consistently either the best performing method or within statistically insignificant different from the best performing method. It is not the best performing approach in models where a linear model approximates the truth well, in which case linear 2SLS outperforms it. Even then it several times outperforms 2SLS or is very close to its performance. The performance does degrade as the number of instruments $d$ increases, as expected. However, even up to $d = 10$ it maintains great performance. We conclude that Adversarial GMM is an experimentally proven method for non-linear IV estimation, which performs well in a multitude of different underlying models without the need for any fine-tuning for the model at hand.

## References

Anthony, Martin and Bartlett, Peter L. *Neural Network Learning: Theoretical Foundations*. Cambridge University Press, New York, NY, USA, 1st edition, 2009. ISBN 052111862X, 9780521118620.

Arjovsky, Martin, Chintala, Soumith, and Bottou, Léon. Wasserstein gan. *arXiv preprint arXiv:1701.07875*, 2017.

Athey, Susan, Tibshirani, Julie, and Wager, Stefan. Generalized random forests. 2016.

Bartlett, Peter L. and Mendelson, Shahar. Rademacher and gaussian complexities: Risk bounds and structural results. *J. Mach. Learn. Res.*, 3:463–482, March 2003. ISSN 1532-4435.

Chamberlain, Gary. Asymptotic efficiency in estimation with conditional moment restrictions. *Journal of Econometrics*, 34(3):305 – 334, 1987. ISSN 0304-4076. doi: https://doi.org/10.1016/0304-4076(87)90015-7.

Chen, Xiaohong and Liao, Zhipeng. Sieve semiparametric two-step gmm under weak dependence. Cowles Foun-

(a) 2dpoly

(b) 3dpoly

(c) abs

(d) linear
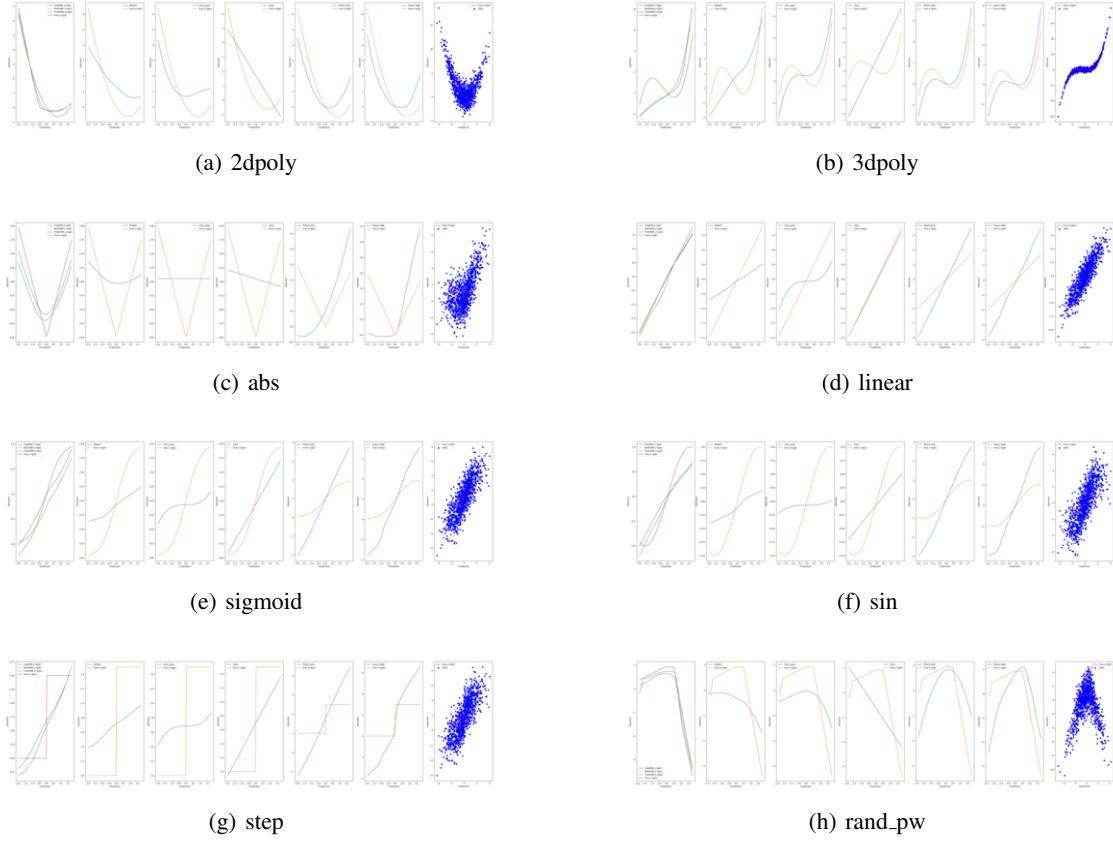
(e) sigmoid

(f) sin

(g) step

(h) rand_pw

*Figure 2.* Examples of fitted functions via each method (in order from left to right: AGMM (best, final, avg), DeepIV, 2SLSPoly, 2SLS, DirectPoly, DirectNN, data points). Data Generating Process: DGP 1, Instrument strength: $\gamma = 0.5$, Number of instruments: $d = 1$, Number of samples: $n = 1000$, Training steps: $T = 400$, Number of critics: $K = 50$, Kernel radius: $r = 50$ data points, Critic jitter: yes.

| function | AGMM Avg | DeepIV | 2SLSpoly | 2SLS | DirectPoly | DirectNN |
|---|---|---|---|---|---|---|
| abs | **0.79** **(0.49, 0.93)** | 0.11 (-0.09, 0.24) | -0.16 (-0.34, -0.02) | -0.18 (-0.38, -0.06) | -3.25 (-3.77, -2.65) | -2.90 (-3.78, -2.15) |
| 2dpoly | **0.97** **(0.91, 0.99)** | 0.49 (0.32, 0.60) | 0.47 (0.28, 0.59) | 0.56 (0.42, 0.65) | 0.67 (0.57, 0.72) | 0.66 (0.55, 0.74) |
| sigmoid | **0.90** **(0.60, 0.98)** | 0.53 (0.30, 0.71) | 0.23 (0.16, 0.33) | 0.89 (0.79, 0.96) | -1.16 (-1.43, -0.82) | -1.26 (-1.91, -0.75) |
| step | **0.67** **(0.41, 0.79)** | 0.40 (0.22, 0.55) | 0.14 (0.08, 0.21) | 0.66 (0.58, 0.73) | -1.05 (-1.35, -0.74) | -1.11 (-1.63, -0.59) |
| 3dpoly | 0.21 (-1.30, 0.78) | -0.34 (-1.31, -0.00) | 0.03 (-0.59, 0.40) | -9.57 (-15.19, -5.14) | **0.44** **(0.20, 0.65)** | 0.35 (-0.16, 0.65) |
| sin | **0.89** **(0.58, 0.98)** | 0.41 (0.25, 0.58) | 0.12 (0.04, 0.19) | 0.74 (0.59, 0.85) | -0.67 (-0.96, -0.40) | -0.94 (-1.39, -0.55) |
| linear | 0.97 (0.90, 0.99) | 0.67 (0.50, 0.79) | 0.43 (0.35, 0.52) | **1.00** **(0.98, 1.00)** | 0.00 (-0.14, 0.14) | 0.00 (-0.27, 0.25) |
| rand_pw | **0.86** **(0.04, 0.98)** | 0.41 (-0.14, 0.71) | 0.26 (-0.15, 0.48) | 0.61 (-0.69, 0.99) | 0.26 (-3.73, 0.88) | 0.38 (-3.61, 0.88) |

*Figure 3.* Median $R^2$ and the $5 - 95$ percentiles of the $R^2$ of each method for a grid of 100 treatment points used as test set. Boldface portrays the best performing method in terms of median $R^2$. Data Generating Process: DGP 1, Instrument strength: $\gamma = 0.5$, Number of instruments: $d = 1$, Number of samples: $n = 1000$, Number of experiments: $M = 100$, Training steps: $T = 400$, Number of critics: $K = 50$, Kernel radius: $r = 50$ data points, Critic jitter: yes.

| function | AGMM Avg | DeepIV | 2SLSpoly | 2SLS | DirectPoly | DirectNN |
|---|---|---|---|---|---|---|
| abs | **0.67** (0.17, 0.86) | 0.08 (-0.12, 0.22) | -0.16 (-0.41, -0.05) | -0.20 (-0.45, -0.06) | -5.43 (-6.51, -4.71) | -5.33 (-6.78, -4.27) |
| 2dpoly | **0.95** (0.82, 0.98) | 0.33 (0.17, 0.44) | 0.33 (0.09, 0.45) | 0.60 (0.46, 0.68) | 0.45 (0.33, 0.55) | 0.42 (0.25, 0.56) |
| sigmoid | 0.86 (0.43, 0.96) | 0.39 (0.21, 0.54) | 0.20 (0.09, 0.30) | **0.90** (0.71, 0.96) | -2.02 (-2.46, -1.53) | -2.38 (-3.27, -1.66) |
| step | 0.59 (0.15, 0.74) | 0.29 (0.17, 0.44) | 0.11 (0.03, 0.17) | **0.66** (0.47, 0.73) | -1.52 (-1.93, -1.12) | -1.74 (-2.31, -1.19) |
| 3dpoly | -1.91 (-4.67, -0.15) | -0.68 (-1.90, -0.04) | -1.11 (-2.02, -0.18) | -13.88 (-24.32, -8.76) | **-0.36** (-1.02, 0.19) | -0.54 (-1.47, 0.13) |
| sin | **0.87** (0.42, 0.96) | 0.30 (0.16, 0.47) | 0.11 (-0.00, 0.18) | 0.76 (0.51, 0.89) | -1.41 (-1.81, -1.00) | -2.00 (-2.60, -1.41) |
| linear | 0.96 (0.79, 1.00) | 0.50 (0.32, 0.66) | 0.40 (0.31, 0.48) | **0.99** (0.96, 1.00) | -0.57 (-0.75, -0.37) | -0.65 (-0.98, -0.32) |
| rand_pw | **0.73** (-0.53, 0.98) | 0.28 (-0.27, 0.58) | 0.19 (-0.37, 0.44) | 0.67 (-0.77, 0.99) | -0.14 (-6.99, 0.84) | -0.00 (-6.85, 0.83) |

*Figure 4.* Median $R^2$ and the $5 - 95$ percentiles of the $R^2$ of each method for a grid of 100 treatment points used as test set. Boldface portrays the best performing method in terms of median $R^2$. Data Generating Process: DGP 2, Instrument strength: $\gamma = 0.5$, Number of instruments: $d = 2$, Number of samples: $n = 1000$, Number of experiments: $M = 100$, Training steps: $T = 400$, Number of critics: $K = 50$, Kernel radius: $r = 50$ data points, Critic jitter: yes.

dation Discussion Papers 2012, Cowles Foundation for Research in Economics, Yale University, 2015.

Freund, Yoav and Schapire, Robert E. Adaptive game playing using multiplicative weights. *Games and Economic Behavior*, 29(1):79 – 103, 1999. ISSN 0899-8256. doi: https://doi.org/10.1006/game.1999.0738.

Gallant, A. and Tauchen, George. Which moments to match? *Econometric Theory*, 12(04):657–681, 1996.

Goodfellow, Ian, Pouget-Abadie, Jean, Mirza, Mehdi, Xu, Bing, Warde-Farley, David, Ozair, Sherjil, Courville, Aaron, and Bengio, Yoshua. Generative adversarial nets. In Ghahramani, Z., Welling, M., Cortes, C., Lawrence, N. D., and Weinberger, K. Q. (eds.), *Advances in Neural Information Processing Systems 27*, pp. 2672–2680. Curran Associates, Inc., 2014.

Goodfellow, Ian J. NIPS 2016 tutorial: Generative adversarial networks. *CoRR*, abs/1701.00160, 2017.

Hansen, Lars Peter. Large sample properties of generalized method of moments estimators. *Econometrica*, 50(4): 1029–1054, 1982. ISSN 00129682, 14680262.

Hartford, Jason, Lewis, Greg, Leyton-Brown, Kevin, and Taddy, Matt. Deep IV: A flexible approach for counterfactual prediction. In Precup, Doina and Teh, Yee Whye (eds.), *Proceedings of the 34th International Conference on Machine Learning*, volume 70 of *Proceedings of Machine Learning Research*, pp. 1414–1423, International Convention Centre, Sydney, Australia, 06–11 Aug 2017. PMLR.

Kocaoglu, M., Snyder, C., Dimakis, A. G., and Vishwanath, S. CausalGAN: Learning Causal Implicit Generative Models with Adversarial Training. *ArXiv e-prints*, September 2017.

Koltchinskii, Vladimir and Panchenko, Dmitriy. Rademacher processes and bounding the risk of

function learning. In Giné, Evarist, Mason, David M., and Wellner, Jon A. (eds.), *High Dimensional Probability II*, pp. 443–457, Boston, MA, 2000. Birkhäuser Boston. ISBN 978-1-4612-1358-1.

Manski, Charles F. Anatomy of the selection problem. *The Journal of Human Resources*, 24(3):343–360, 1989. ISSN 0022166X.
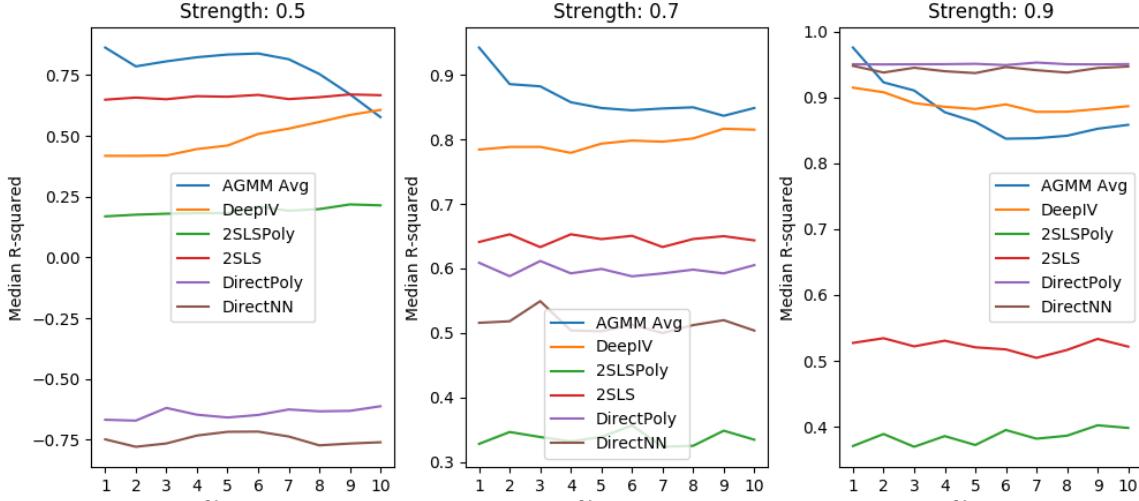
Rakhlin, Alexander and Sridharan, Karthik. Optimization, learning, and games with predictable sequences. In *Proceedings of the 26th International Conference on Neural Information Processing Systems - Volume 2*, NIPS'13, pp. 3066–3074, USA, 2013. Curran Associates Inc.

Shalev-Shwartz, Shai. Online learning and online convex optimization. *Found. Trends Mach. Learn.*, 4(2):107–194, February 2012. ISSN 1935-8237. doi: 10.1561/2200000018.
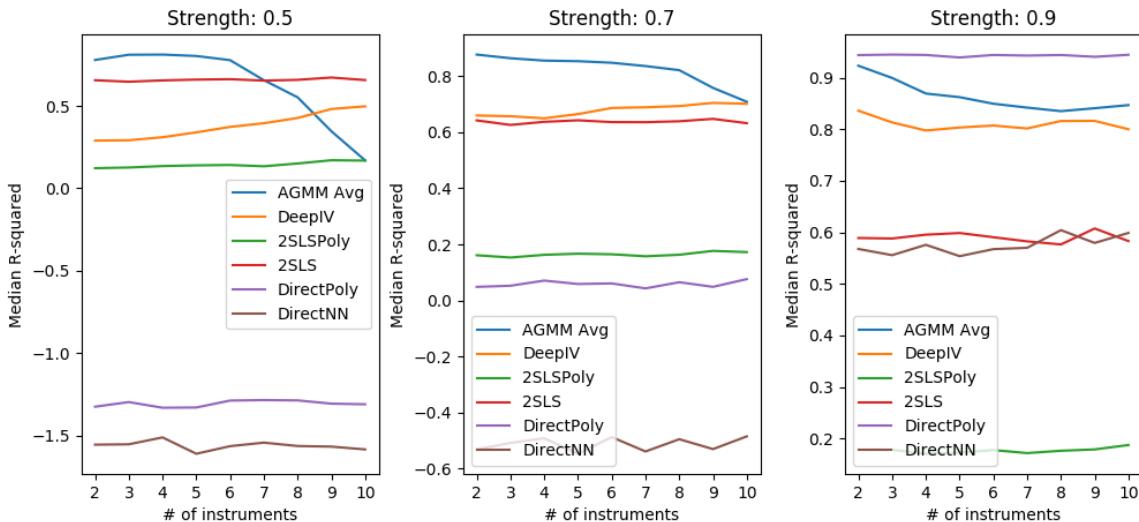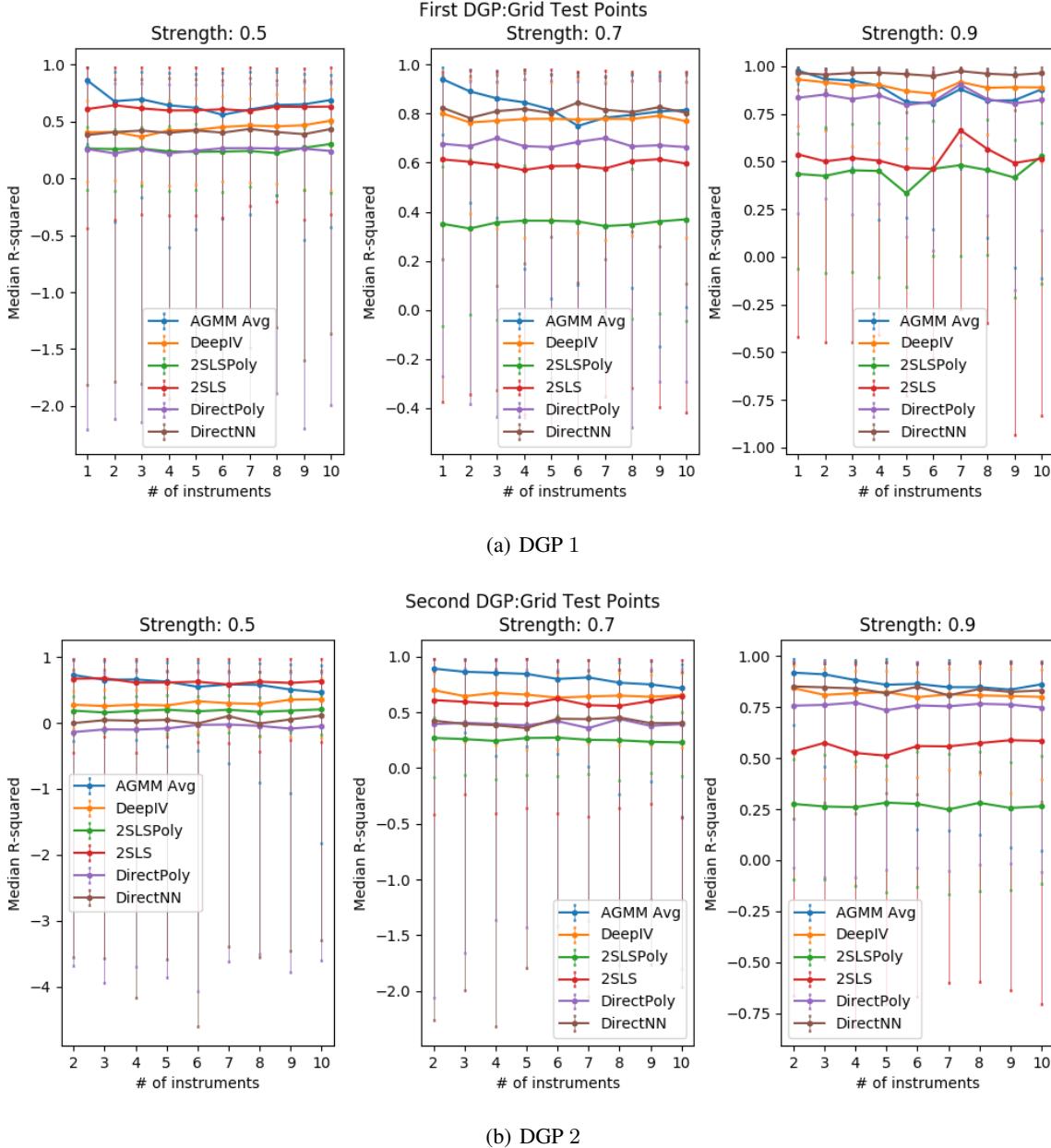
Shalev-Shwartz, Shai and Ben-David, Shai. *Understanding Machine Learning: From Theory to Algorithms*. Cambridge University Press, New York, NY, USA, 2014. ISBN 1107057132, 9781107057135.

Stone, Charles J. Optimal global rates of convergence for nonparametric regression. *Ann. Statist.*, 10(4):1040–1053, 12 1982. doi: 10.1214/aos/1176345969.

Syrgkanis, Vasilis, Agarwal, Alekh, Luo, Haipeng, and Schapire, Robert E. Fast convergence of regularized learning in games. In *Advances in Neural Information Processing Systems*, pp. 2989–2997, 2015.

(a) DGP 1



(b) DGP 2

*Figure 5.* Median $R^2$ of each method for a grid of 100 treatment points used as test set and for all the models $h_0$ put together (except model rand_pw). Performance is portrayed as a function of the number of instruments $d \in \{1, \ldots, 10\}$ and for different values of instrument strength $\gamma \in \{.5, .7, .9\}$. Number of samples: $n = 1000$, Number of experiments: $M = 700$, Training steps: $T = 400$, Number of critics: $K = 50$, Kernel radius: $r = 50$ data points, Critic jitter: yes.

(a) DGP 1



(b) DGP 2

*Figure 6.* Median $R^2$ and the $10 - 90$ percentiles of the $R^2$ of each method for a grid of 100 treatment points used as test set and for true model rand_pw. Performance is portrayed as a function of the number of instruments $d \in \{1, \ldots, 10\}$ and for different values of instrument strength $\gamma \in \{.5, .7, .9\}$. Number of samples: $n = 1000$, Number of experiments: $M = 100$, Training steps: $T = 400$, Number of critics: $K = 50$, Kernel radius: $r = 50$ data points, Critic jitter: yes.

# A. Omitted Proofs

## A.1. Lipschitz Conditional Moments

Let $\mathcal{X} = [0,1]^d$ and suppose that the function $\Psi(h, x) = \mathbb{E}[\rho(z; h)|x]$ is $\lambda$-Lipschitz with respect to $x$ and with respect to the $\|\cdot\|_\infty$ norm and that the density of the distribution of $x$ is lower bounded by $\mu > 0$. Then consider the set of test functions corresponding to uniform Kernels around a set of grid points, i.e.: discretize the space $\mathcal{X}$ to a grid of multiples of some number $h$ (the bandwidth) and let $\mathcal{X}_h$ denote the discretized set of points. Then for each $x_0 \in \mathcal{X}_h$, let $f(x; x_0) = 1\{\|x - x_0\|_\infty \le h\}$. For any point $x^*$, let $x_h^*$ denote its closest grid point. Then observe that if $|\mathbb{E}[\rho(z; h)f(x; x_h^*)]| \le \epsilon$, then:

$$
\begin{aligned}
\epsilon &\ge |\mathbb{E}[\rho(z; h)f(x; x_h^*)]| \\
&\ge |\mathbb{E}[\mathbb{E}[\rho(z; h)|x]1\{\|x - x_h^*\|_\infty \le h\}]| \\
&= |\mathbb{E}[\mathbb{E}[\rho(z; h)|x_h^*]1\{\|x - x_h^*\|_\infty \le h\}] + \mathbb{E}[(\mathbb{E}[\rho(z; h)|x] - \mathbb{E}[\rho(z; h)|x_h^*])1\{\|x - x_h^*\|_\infty \le h\}]| \\
&\ge |\mathbb{E}[\rho(z; h)|x_h^*]|\Pr[\|x - x_h^*\|_\infty \le h] - |\mathbb{E}[(\mathbb{E}[\rho(z; h)|x] - \mathbb{E}[\rho(z; h)|x_h^*])1\{\|x - x_h^*\|_\infty \le h\}]| \\
&\ge |\mathbb{E}[\rho(z; h)|x_h^*]|\Pr[\|x - x_h^*\|_\infty \le h] - h\lambda\Pr[\|x - x_h^*\|_\infty \le h]
\end{aligned}
$$

Hence, by re-arranging:

$$
|\mathbb{E}[\rho(z; h)|x_h^*]| \le h\lambda + \frac{\epsilon}{\Pr[\|x - x_h^*\|_\infty \le h]} \le h\lambda + \frac{\epsilon}{\mu h^d}
$$

Thus the class of test functions $\mathcal{F} = \{f(x; x_0) : x_0 \in \mathcal{X}_h\}$ is a set of $\gamma$-test functions with $\gamma(\epsilon) = h\lambda + \frac{\epsilon}{\mu h^d}$. If we have a target $\epsilon$ in mind, we can set the optimal bandwidth $h = (\epsilon/\lambda\mu)^{1/(d+1)}$, to get $\gamma(\epsilon) = 2\lambda^{d/(d+1)}(\epsilon/\mu)^{1/(d+1)}$. The latter slow convergence with respect to $d$ is a typical rate in non-parametric regression problems in $d$-dimensions (Stone, 1982).

Furthermore, we can show that this class also has good generalization properties. Specifically, observe that $\mathcal{A} = \{\rho(\cdot; h_\theta)1\{\|\cdot - x_0\|_\infty \le h\} : \theta \in \Theta, x_0 \in \mathcal{X}_h\}$. Further assume that $\rho(\cdot; h_\theta) \in [0,1]$ is a $\lambda$-Lipschitz function of $\theta$ and $\theta \in [0,1]^r$. Then by standard covering number arguments, the latter class $\mathcal{A}$ can be approximated to within $\epsilon$ by a finite class of functions of size $N = O((1/\epsilon)^r(1/h)^d)$. Therefore the Rademacher complexity is bounded by $O\left(\sqrt{\frac{r\log(r) + d\log(1/h)}{n}}\right)$ (see e.g. Lemma 27.5 of (Shalev-Shwartz & Ben-David, 2014)). Moreover, $|\mathcal{F}| = (1/h)^d$ and we remind that $\gamma(\epsilon) = O(h + \frac{\epsilon}{h^d})$. Combining all the above we get that if we run the dynamics for $T = O(n)$ iterations, then we are guaranteed that the averages are an $\epsilon$-equilibrium of the population game for $\epsilon = O\left(\frac{\sqrt{d\log(1/h)} + \sqrt{r\log(r)} + \sqrt{\log(1/\delta)}}{\sqrt{n}}\right)$ with probability $1 - \delta$ and that the conditional moment violations are upper bounded by $\gamma(\epsilon)$. Balancing $h$ appropriately, we get an error rate of the order of $n^{-1/(2(d+1))}$.

## A.2. Proof of Theorem 1

For $h \in \mathcal{H}$ and $\sigma \in \Delta(\mathcal{F})$ let $L(h, \sigma) = \mathbb{E}_{f \sim \sigma}\left[(\mathbb{E}[\rho(z; h)f(x)])^2\right]$ denote the loss of the zero-sum game. Observe that if the modeler chooses an $h \in \mathcal{H}_I$, then he is guaranteed zero loss. Hence, the value of the game is zero. Subsequently, it is easy to see that at any $\epsilon$-equilibrium $(h^*, w^*)$: $\sup_{f \in \mathcal{F}}(\mathbb{E}[\rho(z; h)f(x)])^2 \le \epsilon \Rightarrow \sup_{f \in \mathcal{F}}|\mathbb{E}[\rho(z; h)f(x)]| \le \sqrt{\epsilon}$. The latter also implies that this inequality holds for any convex combination of functions in $\mathcal{F}$, i.e. for any $\bar{f} \in \bar{\mathcal{F}}$: $|\mathbb{E}[\rho(z; h)\bar{f}(x)]| \le \sqrt{\epsilon}$. Hence, by the triangle inequality and the property of $\gamma$-test functions, for any $x \in \mathcal{X}$: $|\mathbb{E}[\rho(z; h)|x]| \le \gamma(\sqrt{\epsilon})$. The theorem then follows by Equation (5).

## A.3. Proof of Theorem 2

By Corollary 2.14 and 2.17 of (Shalev-Shwartz, 2012) we have small regret. In particular:

$$
\frac{1}{T}\sum_{t=1}^T L_n(\theta_t, \sigma_t) \le \inf_{\theta \in \Theta} \frac{1}{T}\sum_{t=1}^T L_n(\theta, \sigma_t) + \epsilon_1(T)
$$

$$
\frac{1}{T}\sum_{t=1}^T L_n(\theta_t, \sigma_t) \ge \sup_{\sigma \in \Delta(\mathcal{F})} \frac{1}{T}\sum_{t=1}^T L_n(\theta_t, \sigma) - \epsilon_2(T)
$$

for $\epsilon_1(T) = \frac{BL\sqrt{2}}{\sqrt{T}}$ and $\epsilon_2(T) = \frac{H^2\sqrt{2\log(|\mathcal{F}|)}}{\sqrt{T}}$. Subsequently, since the game is convex in $\theta$ and concave in $\sigma$, we can use the well-known fact that the average solutions $\theta^*$ and $\sigma^*$ is an $\epsilon_1(T) + \epsilon_2(T)$-approximate equilibrium (see e.g. (Rakhlin & Sridharan, 2013)).

### A.4. Proof of Theorem 3

By the result of (Koltchinskii & Panchenko, 2000; Bartlett & Mendelson, 2003) connecting Rademacher complexity and uniform convergence of empirical processes, we have with probability $1 - \delta$:

$$\sup_{\theta \in \Theta, f \in \mathcal{F}} |\mathbb{E}_n\left[\rho(z; h_\theta)f(x)\right] - \mathbb{E}\left[\rho(z; h_\theta)f(x)\right]| \leq \Delta \triangleq O\left(\mathcal{R} + H\sqrt{\frac{\log(1/\delta)}{n}}\right)$$

Thus since $(\theta^*, \sigma^*)$ is an $\epsilon$-equilibrium of the empirical game, where $\epsilon$ is given by Theorem 2, we get that it is also an $O(\epsilon + \Delta)$ equilibrium of the population game. More formally, let $L(\theta, \sigma) = \mathbb{E}_{f \sim \sigma}\left[\left(\mathbb{E}\left[\rho(z; h_\theta)f(x)\right]\right)^2\right]$ be the loss of the population game. Then by the uniform convergence property and the boundedness of the moment by $H$, we have with probability $1 - \delta$:

$$\sup_{\theta, \sigma} |L_n(\theta, \sigma) - L(\theta, \sigma)| \leq 2H\Delta$$

Subsequently, we can check that $\theta^*$ and $\sigma^*$ satisfy the approximate equilibrium conditions

$$\begin{aligned}
L(\theta^*, \sigma^*) &\leq L_n(\theta^*, \sigma^*) + 2H\Delta \\
&\leq \inf_{\theta \in \Theta} L_n(\theta, \sigma^*) + \epsilon + 2H\Delta \\
&\leq \inf_{\theta \in \Theta} L(\theta, \sigma^*) + \epsilon + 4H\Delta
\end{aligned}$$

Similarly for the maximizing player.

## B. Examples of Local Kernels

We simulated a data generating process with two instruments. Instrument 1 affects the treatment only when it is negative and instrument 2 affects the treatment only when it is positive. We wanted to see the types of discretizations $\mathcal{F}_\epsilon$, i.e. local kernels that the different algorithms we propose would have created. Below we depict the kernels based on, i) choosing random data points as centers and putting local gaussians around them with standard deviation equal to the distance of this point to the $k$-th nearest neighbor, where $k$ is a hyperparameter, ii) a random forest, where each leaf defines a local kernel, iii) a random forest where we put local gaussians at the centers of each leaf and with standard deviation equal to the maximum distance from the center of any data point in the leaf. The forest was trained by regressing the treatment on the instrument. We see that the forest correctly picks up that the treatment is more sensitive to the instruments in the lower quadrant and hence puts more local kernels around that area.



(a) Gaussians centered at random points (b) Decision tree trained by regression of $w$ on $x$ (c) Gaussian kernel centered on leafs of decision tree

# C. Further Figures of Experimental Results

## C.1. First DGP
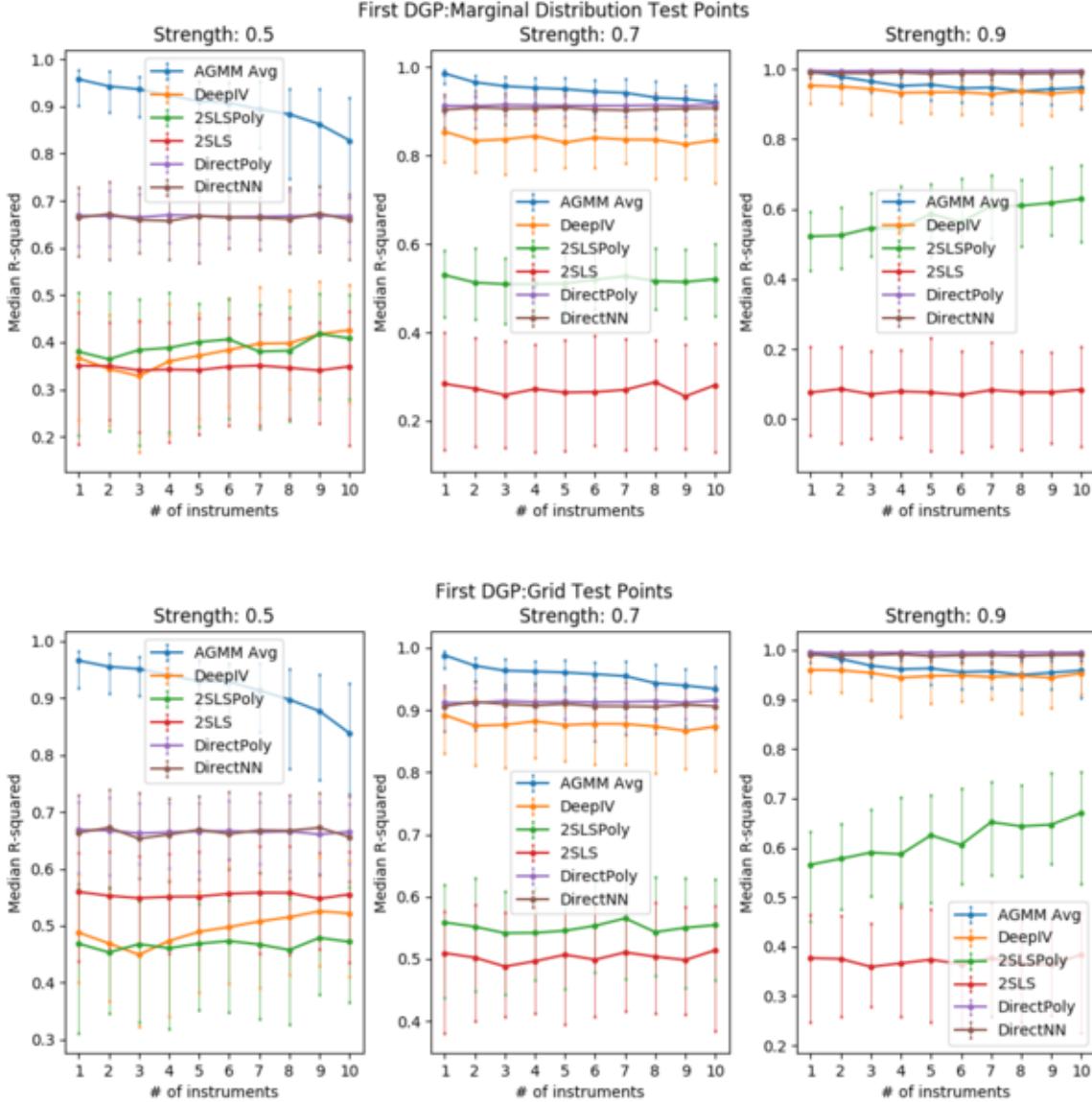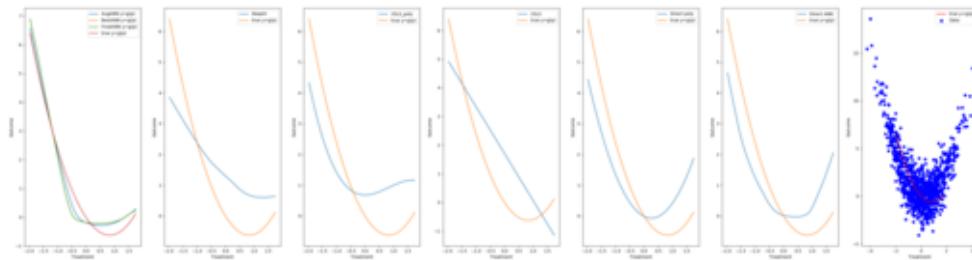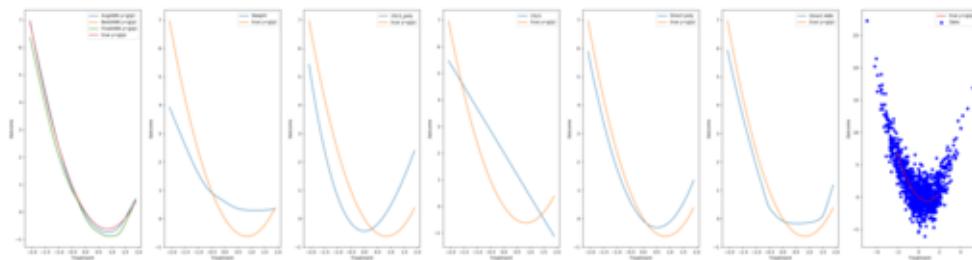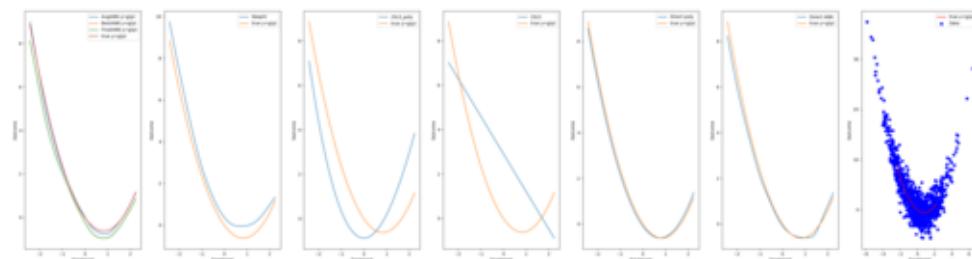
### C.1.1. SECOND DEGREE POLYNOMIAL



*Figure 7.* Median and $10 - 90$ percentiles of $R^2$ across 100 experiments as a function of number of instruments and instrument strength $\gamma$. $h_0(w) = -1.5w + .9w^2$. Number of samples: 1000, Training steps: 400, Number of critics: 50, Kernel radius: 50 data points, Critic jitter: yes
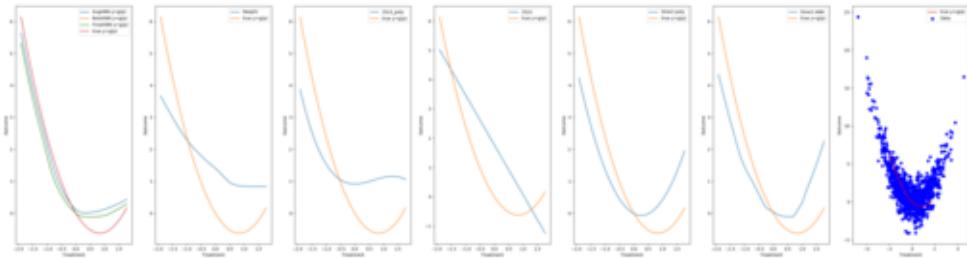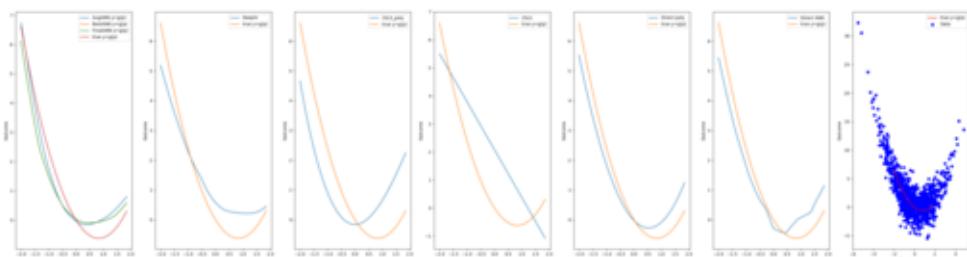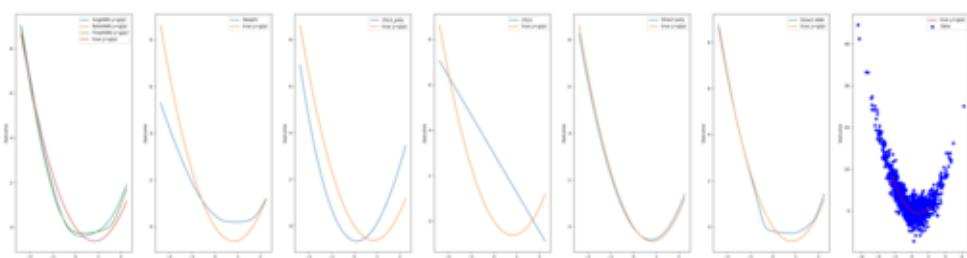
(a) Strength 0.5



(b) Strength 0.7



(c) Strength 0.9

*Figure 8.* Examples of fitted functions via each method (in order from left to right: AGMM (best, final, avg), DeepIV, 2SLSPoly, 2SLS, DirectPoly, DirectNN, data points). $h_0(w) = -1.5w + .9w^2$, Number of instruments: 1, Number of samples: 1000, Training steps: 400, Number of critics: 50, Kernel radius: 50 data points, Critic jitter: yes

(a) Strength 0.5



(b) Strength 0.7



(c) Strength 0.9

*Figure 9.* Examples of fitted functions via each method (in order from left to right: AGMM (best, final, avg), DeepIV, 2SLSPoly, 2SLS, DirectPoly, DirectNN, data points). $h_0(w) = -1.5w + .9w^2$, Number of instruments: 5, Number of samples: 1000, Training steps: 400, Number of critics: 50, Kernel radius: 50 data points, Critic jitter: yes
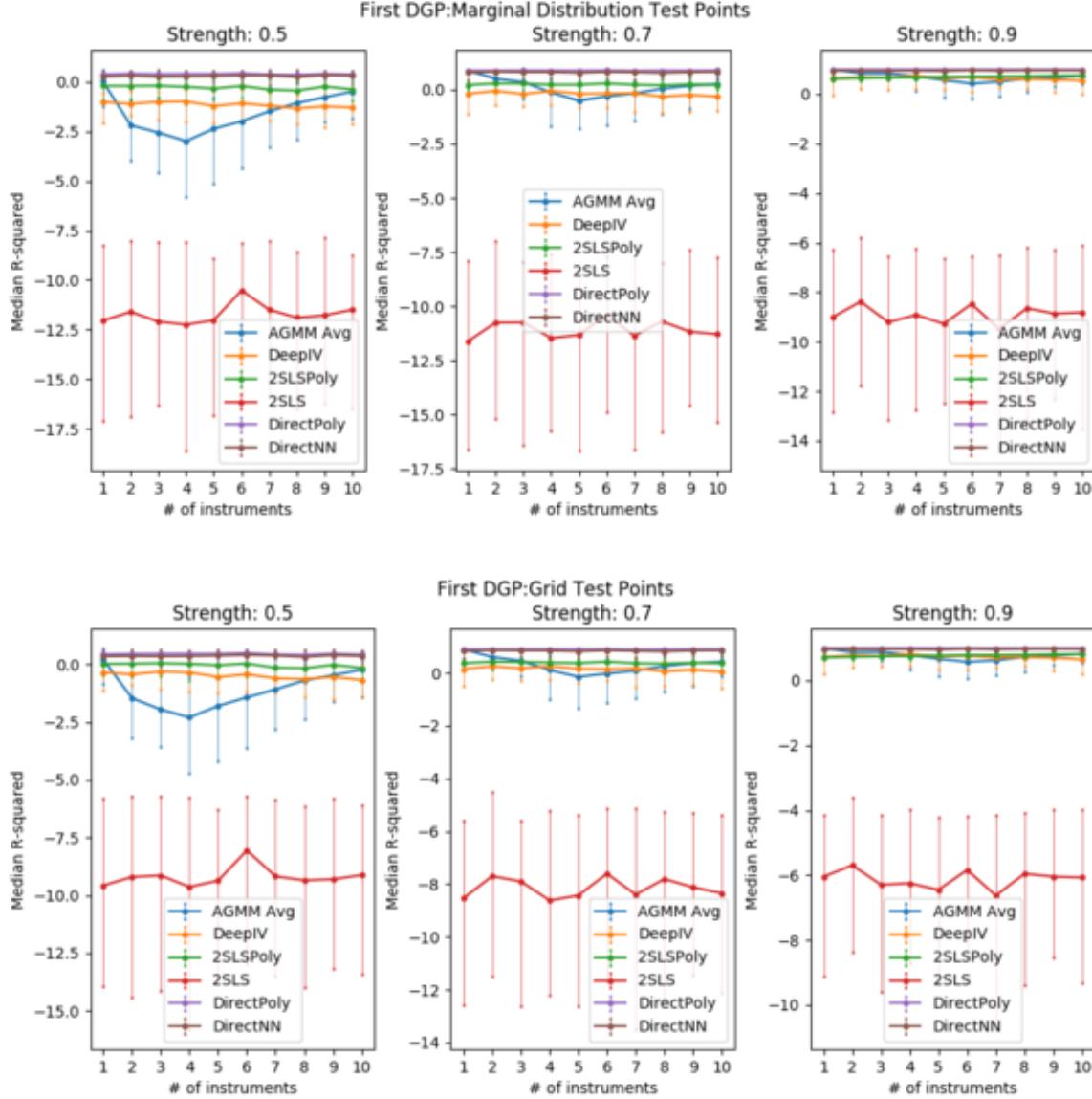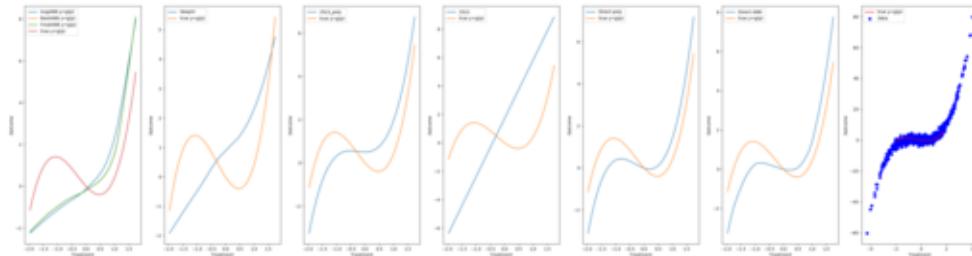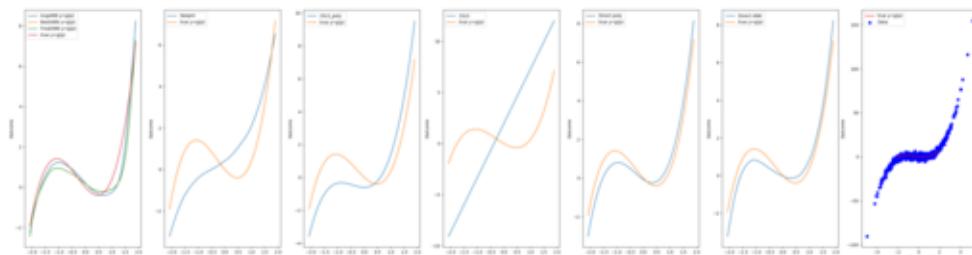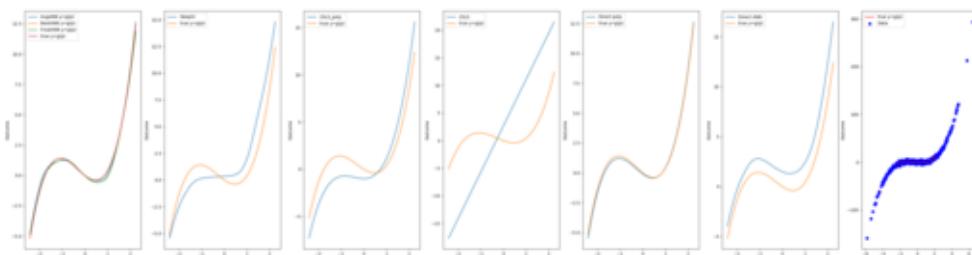
## C.1.2. THIRD DEGREE POLYNOMIAL



*Figure 10.* Median and $10 - 90$ percentiles of $R^2$ across 100 experiments as a function of number of instruments and instrument strength $\gamma$. $h_0(w) = -1.5w + .9w^2 + w^3$. Number of samples: 1000, Training steps: 400, Number of critics: 50, Kernel radius: 50 data points, Critic jitter: yes
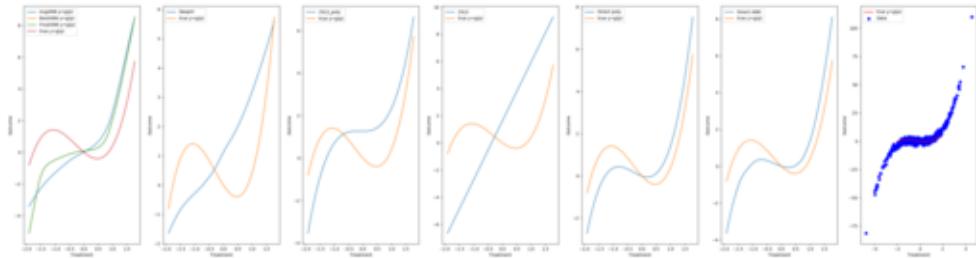
(a) Strength 0.5



(b) Strength 0.7



(c) Strength 0.9

*Figure 11.* Examples of fitted functions via each method (in order from left to right: AGMM (best, final, avg), DeepIV, 2SLSPoly, 2SLS, DirectPoly, DirectNN, data points). $h_0(w) = -1.5w + .9w^2 + w^3$, Number of instruments: 1, Number of samples: 1000, Training steps: 400, Number of critics: 50, Kernel radius: 50 data points, Critic jitter: yes

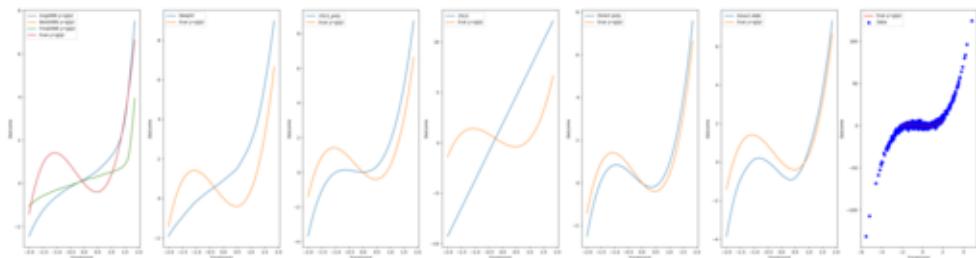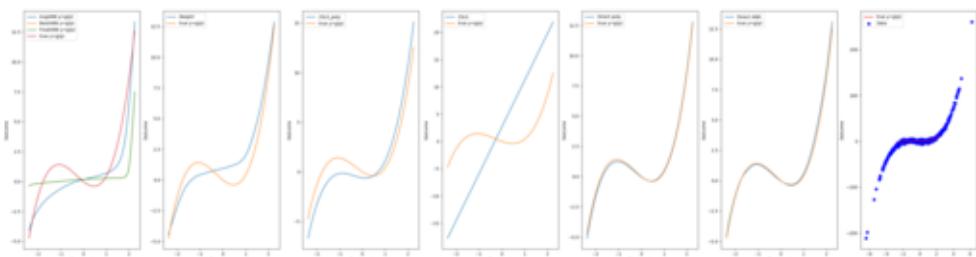(a) Strength 0.5



(b) Strength 0.7



(c) Strength 0.9

*Figure 12.* Examples of fitted functions via each method (in order from left to right: AGMM (best, final, avg), DeepIV, 2SLSPoly, 2SLS, DirectPoly, DirectNN, data points). $h_0(w) = -1.5w + .9w^2 + w^3$, Number of instruments: 5, Number of samples: 1000, Training steps: 400, Number of critics: 50, Kernel radius: 50 data points, Critic jitter: yes

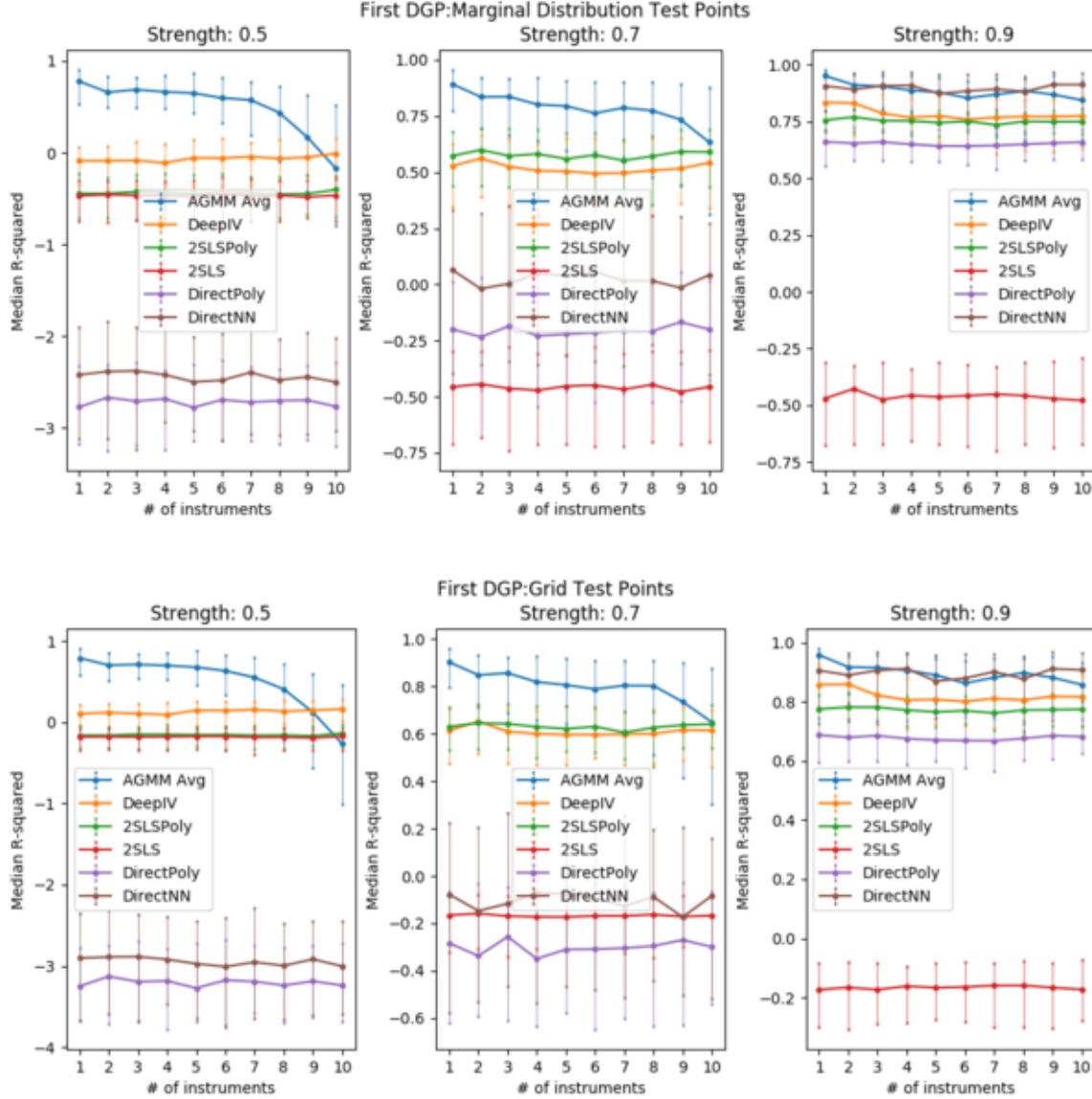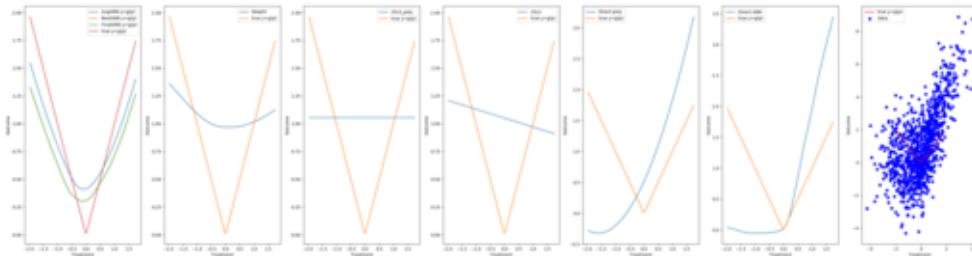## C.1.3. ABSOLUTE VALUE



*Figure 13.* Median $R^2$ as a function of number of instruments and instrument strength. $h_0(w) = |w|$. Number of samples: 1000, Training steps: 400, Number of critics: 50, Kernel radius: 50 data points, Critic jitter: yes

(a) Strength 0.5



(b) Strength 0.7



(c) Strength 0.9

*Figure 14.* Examples of fitted functions via each method (in order from left to right: AGMM (best, final, avg), DeepIV, 2SLSPoly, 2SLS, DirectPoly, DirectNN, data points). $h_0(w) = |w|$, Number of instruments: 1, Number of samples: 1000, Training steps: 400, Number of critics: 50, Kernel radius: 50 data points, Critic jitter: yes

(a) Strength 0.5



(b) Strength 0.7



(c) Strength 0.9

*Figure 15.* Examples of fitted functions via each method (in order from left to right: AGMM (best, final, avg), DeepIV, 2SLSPoly, 2SLS, DirectPoly, DirectNN, data points). $h_0(w) = |w|$, Number of instruments: 5, Number of samples: 1000, Training steps: 400, Number of critics: 50, Kernel radius: 50 data points, Critic jitter: yes

## C.1.4. IDENTIFY FUNCTION
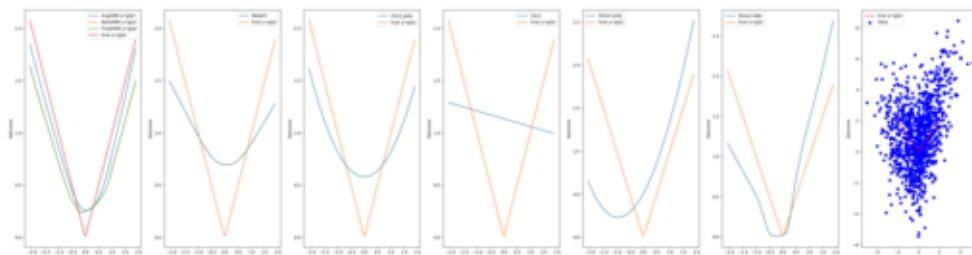


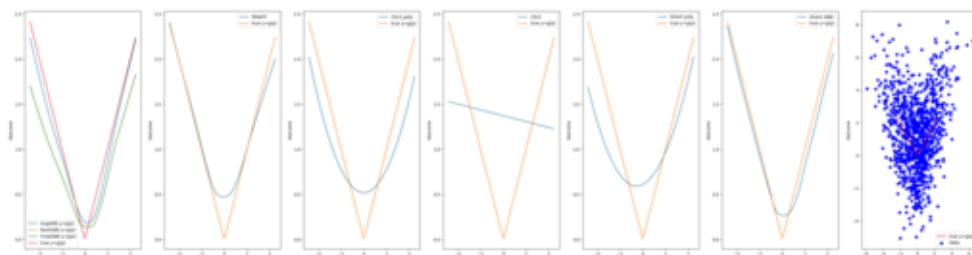*Figure 16.* Median and $10 - 90$ percentiles of $R^2$ across 100 experiments as a function of number of instruments and instrument strength $\gamma$. $h_0(w) = x$. Number of samples: 1000, Training steps: 400, Number of critics: 50, Kernel radius: 50 data points, Critic jitter: yes
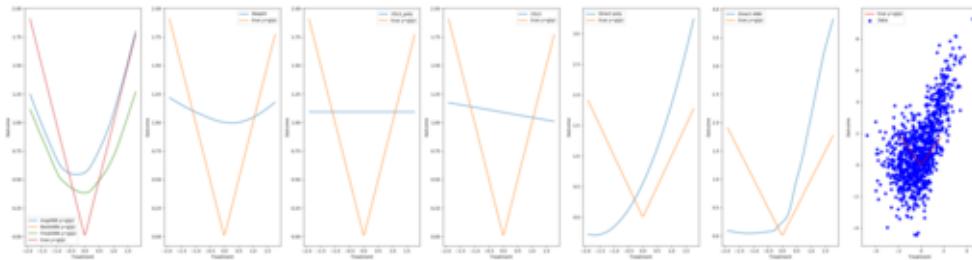
(a) Strength 0.5



(b) Strength 0.7



(c) Strength 0.9

*Figure 17.* Examples of fitted functions via each method (in order from left to right: AGMM (best, final, avg), DeepIV, 2SLSPoly, 2SLS, DirectPoly, DirectNN, data points). $h_0(w) = w$, Number of instruments: 1, Number of samples: 1000, Training steps: 400, Number of critics: 50, Kernel radius: 50 data points, Critic jitter: yes
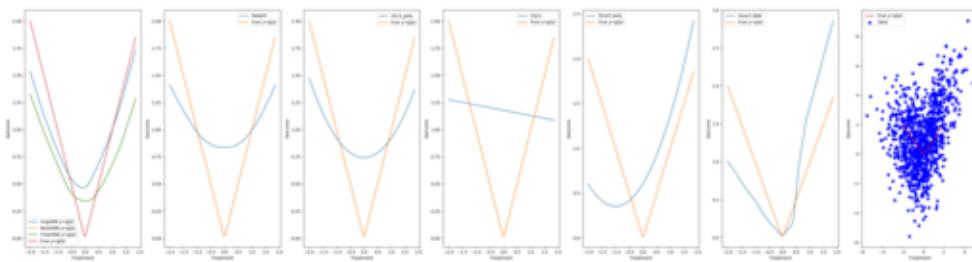
(a) Strength 0.5



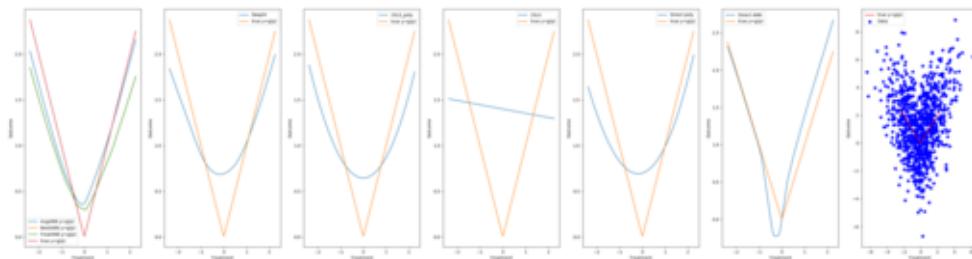(b) Strength 0.7



(c) Strength 0.9

*Figure 18.* Examples of fitted functions via each method (in order from left to right: AGMM (best, final, avg), DeepIV, 2SLSPoly, 2SLS, DirectPoly, DirectNN, data points). $h_0(w) = w$, Number of instruments: 5, Number of samples: 1000, Training steps: 400, Number of critics: 50, Kernel radius: 50 data points, Critic jitter: yes

## C.1.5. SIGMOID FUNCTION



*Figure 19.* Median and $10 - 90$ percentiles of $R^2$ across 100 experiments as a function of number of instruments and instrument strength $\gamma$. $h_0(w) = \frac{2}{1+e^{-2x}}$. Number of samples: 1000, Training steps: 400, Number of critics: 50, Kernel radius: 50 data points, Critic jitter: yes

(a) Strength 0.5



(b) Strength 0.7



(c) Strength 0.9

*Figure 20.* Examples of fitted functions via each method (in order from left to right: AGMM (best, final, avg), DeepIV, 2SLSPoly, 2SLS, DirectPoly, DirectNN, data points). $h_0(w) = \frac{2}{1+e^{-2x}}$, Number of instruments: 1, Number of samples: 1000, Training steps: 400, Number of critics: 50, Kernel radius: 50 data points, Critic jitter: yes

(a) Strength 0.5



(b) Strength 0.7



(c) Strength 0.9

*Figure 21.* Examples of fitted functions via each method (in order from left to right: AGMM (best, final, avg), DeepIV, 2SLSPoly, 2SLS, DirectPoly, DirectNN, data points). $h_0(w) = \frac{2}{1+e^{-2x}}$, Number of instruments: 5, Number of samples: 1000, Training steps: 400, Number of critics: 50, Kernel radius: 50 data points, Critic jitter: yes
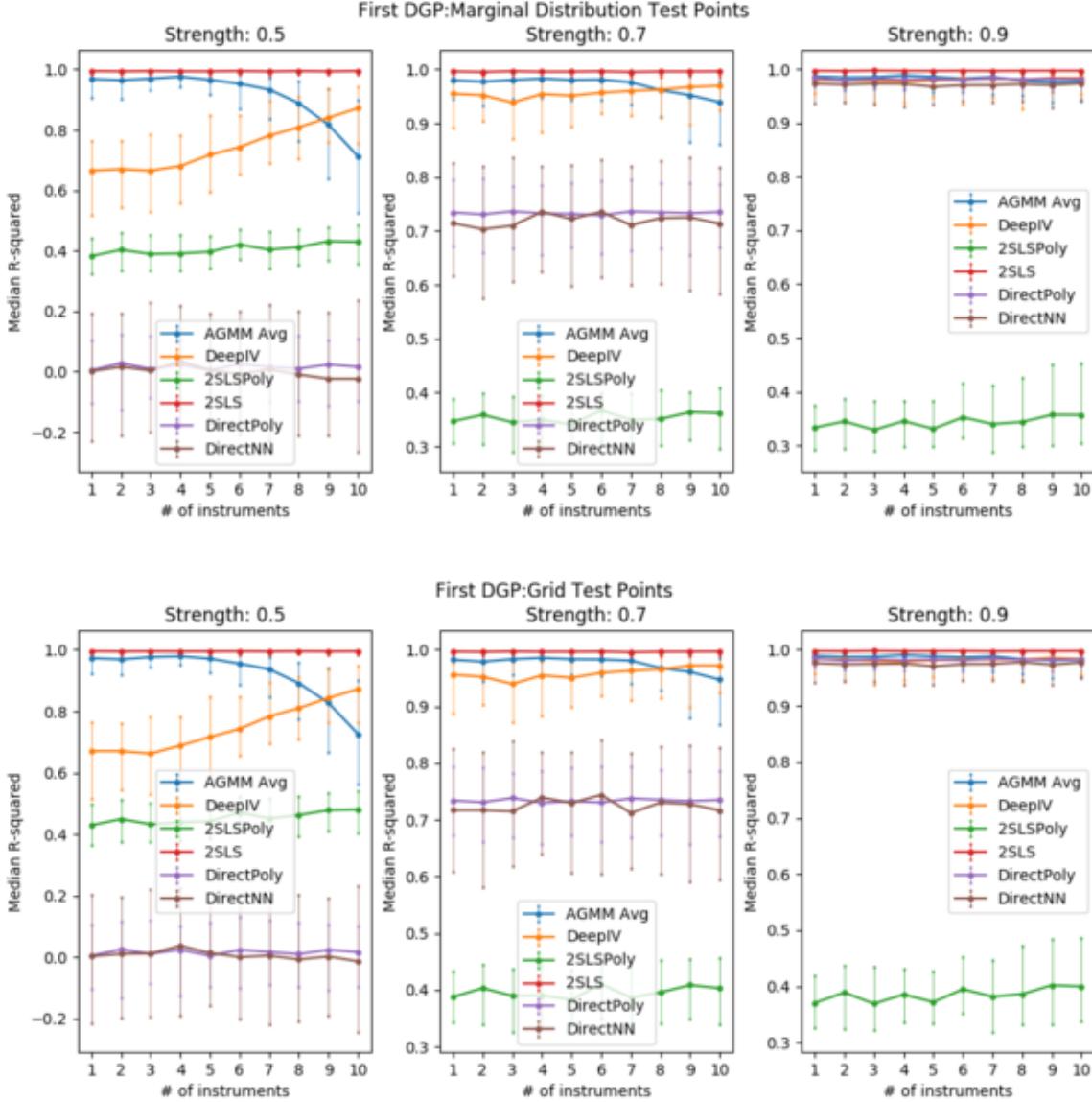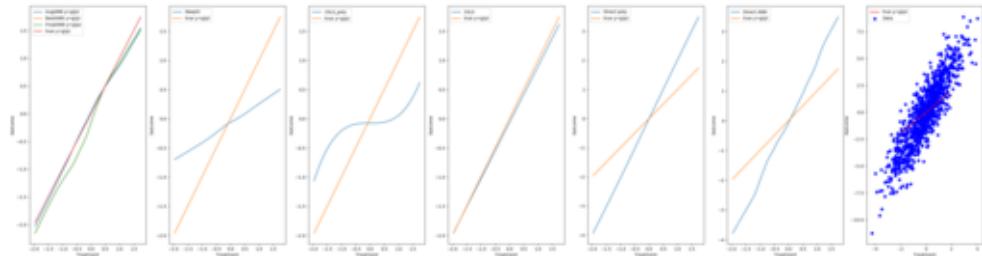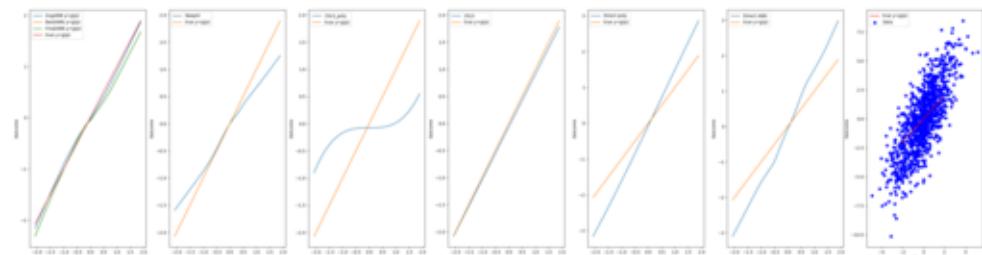
## C.1.6. SIN FUNCTION



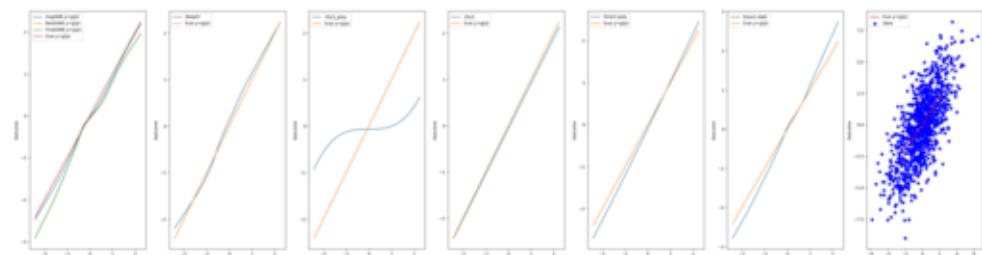*Figure 22.* Median and $10 - 90$ percentiles of $R^2$ across 100 experiments as a function of number of instruments and instrument strength $\gamma$. $h_0(w) = \sin(x)$. Number of samples: 1000, Training steps: 400, Number of critics: 50, Kernel radius: 50 data points, Critic jitter: yes
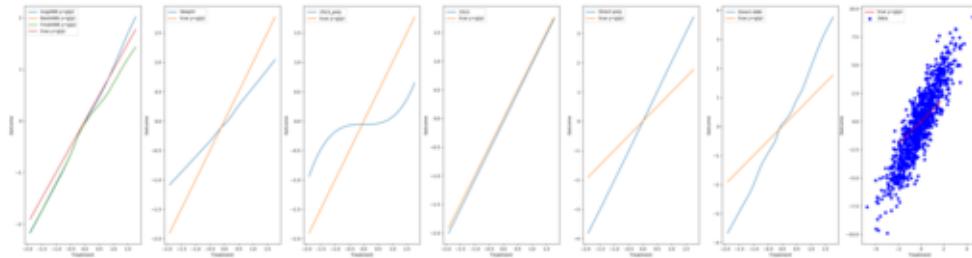
(a) Strength 0.5



(b) Strength 0.7



(c) Strength 0.9

*Figure 23.* Examples of fitted functions via each method (in order from left to right: AGMM (best, final, avg), DeepIV, 2SLSPoly, 2SLS, DirectPoly, DirectNN, data points). $h_0(w) = sin(x)$, Number of instruments: 1, Number of samples: 1000, Training steps: 400, Number of critics: 50, Kernel radius: 50 data points, Critic jitter: yes
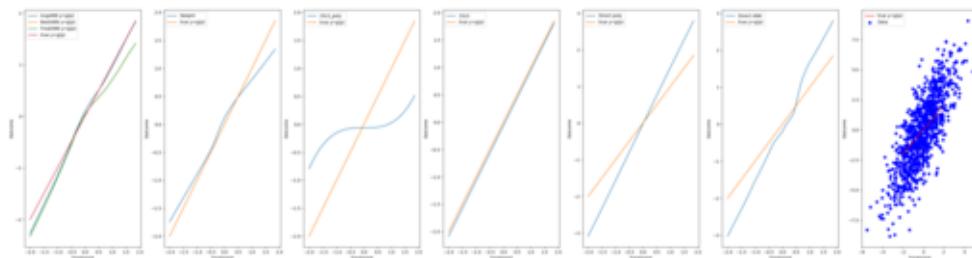
(a) Strength 0.5
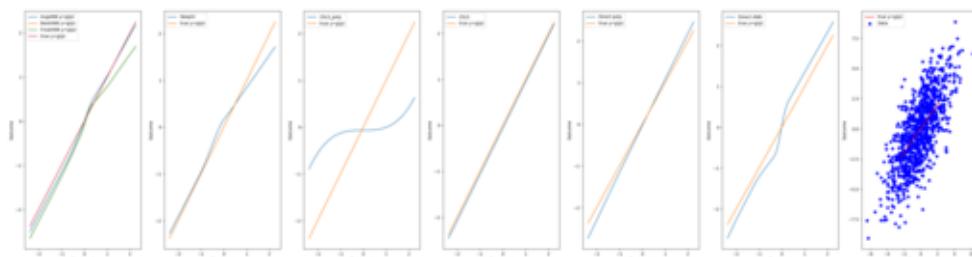


(b) Strength 0.7



(c) Strength 0.9

*Figure 24.* Examples of fitted functions via each method (in order from left to right: AGMM (best, final, avg), DeepIV, 2SLSPoly, 2SLS, DirectPoly, DirectNN, data points). $h_0(w) = sin(x)$, Number of instruments: 5, Number of samples: 1000, Training steps: 400, Number of critics: 50, Kernel radius: 50 data points, Critic jitter: yes

## C.1.7. STEP FUNCTION



*Figure 25.* Median and $10 - 90$ percentiles of $R^2$ across 100 experiments as a function of number of instruments and instrument strength $\gamma$. $h_0(w) = 1\{x < 0\} + 2.51\{x \geq 0\}$. Number of samples: 1000, Training steps: 400, Number of critics: 50, Kernel radius: 50 data points, Critic jitter: yes

(a) Strength 0.5



(b) Strength 0.7



(c) Strength 0.9

*Figure 26.* Examples of fitted functions via each method (in order from left to right: AGMM (best, final, avg), DeepIV, 2SLSPoly, 2SLS, DirectPoly, DirectNN, data points). $h_0(w) = 1\{x < 0\} + 2.51\{x \geq 0\}$, Number of instruments: 1, Number of samples: 1000, Training steps: 400, Number of critics: 50, Kernel radius: 50 data points, Critic jitter: yes

(a) Strength 0.5



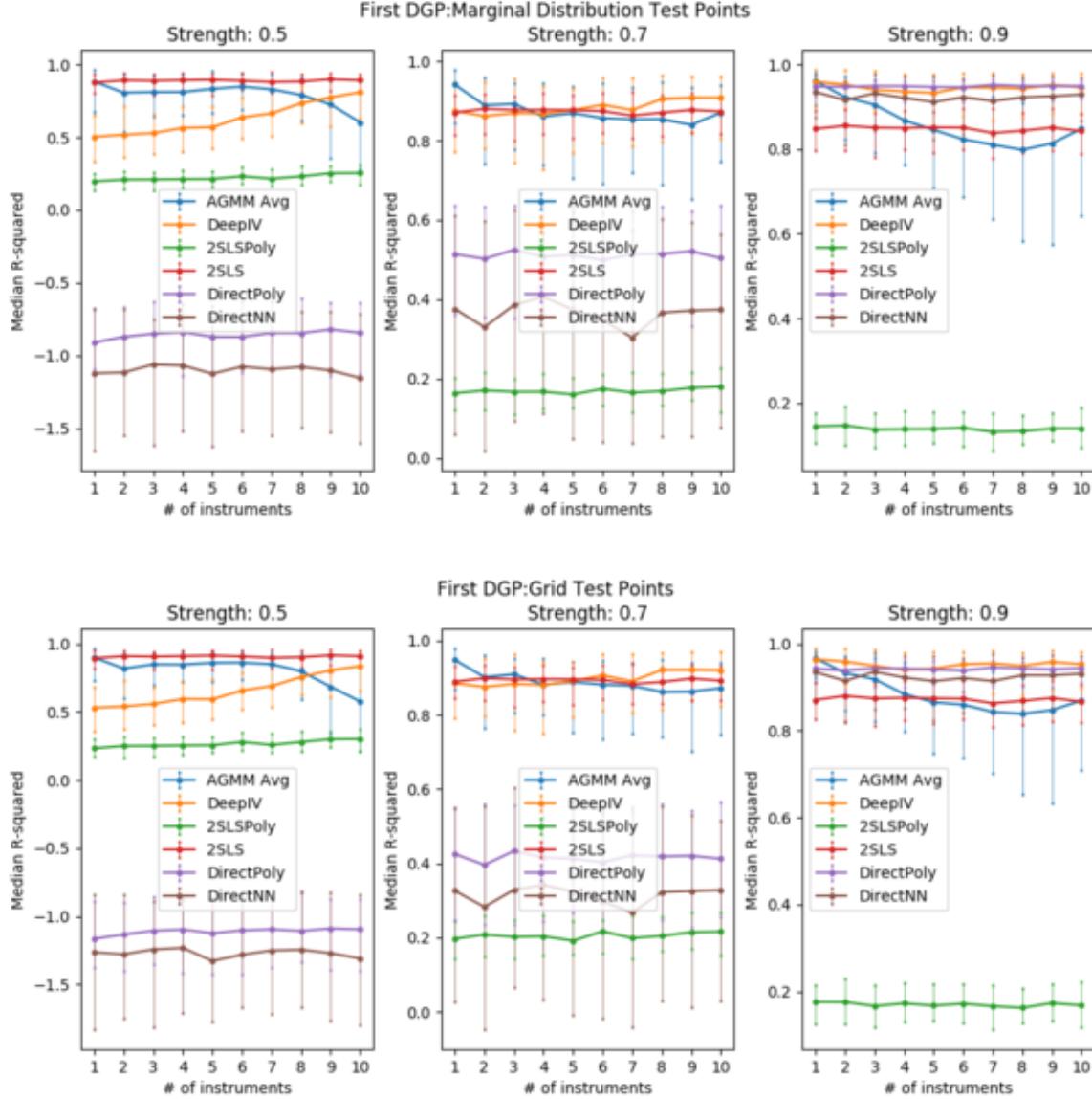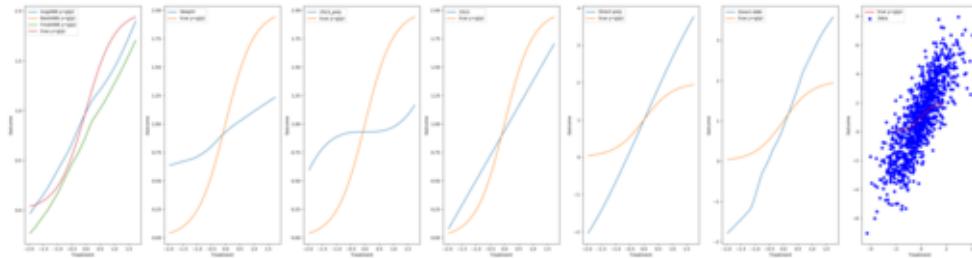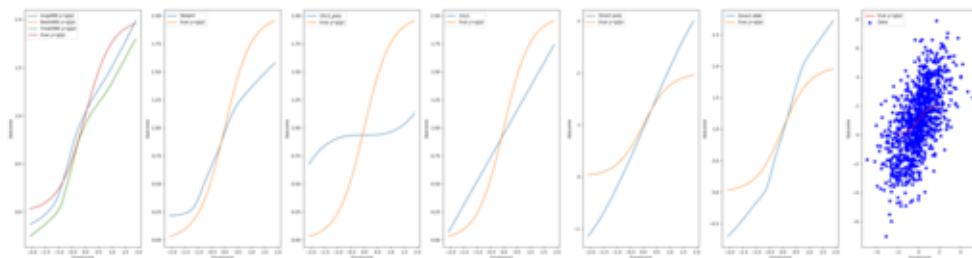(b) Strength 0.7



(c) Strength 0.9

*Figure 27.* Examples of fitted functions via each method (in order from left to right: AGMM (best, final, avg), DeepIV, 2SLSPoly, 2SLS, DirectPoly, DirectNN, data points). $h_0(w) = 1\{x < 0\} + 2.51\{x \geq 0\}$, Number of instruments: 5, Number of samples: 1000, Training steps: 400, Number of critics: 50, Kernel radius: 50 data points, Critic jitter: yes
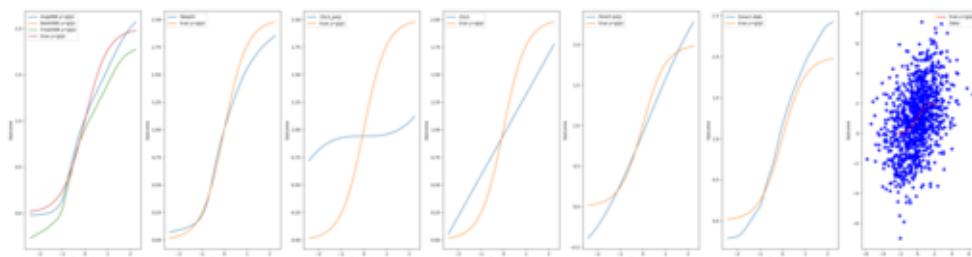
## C.2. Second DGP

### C.2.1. SECOND DEGREE POLYNOMIAL



*Figure 28.* Median and $10-90$ percentiles of $R^2$ across 100 experiments as a function of number of instruments and instrument strength $\gamma$. $h_0(w) = -1.5w + .9w^2$. Number of samples: 1000, Training steps: 400, Number of critics: 50, Kernel radius: 50 data points, Critic jitter: yes
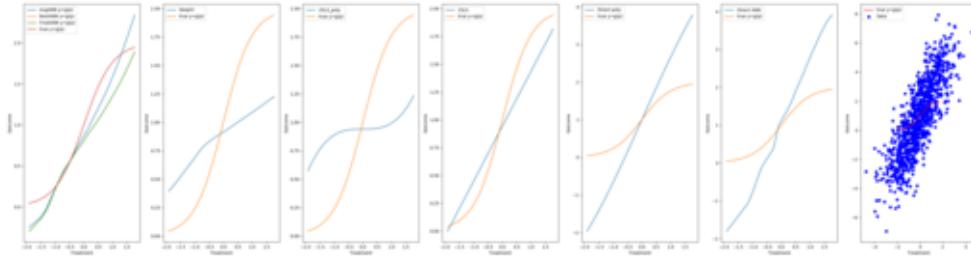
(a) Strength 0.5



(b) Strength 0.7



(c) Strength 0.9

*Figure 29.* Examples of fitted functions via each method (in order from left to right: AGMM (best, final, avg), DeepIV, 2SLSPoly, 2SLS, DirectPoly, DirectNN, data points). $h_0(w) = -1.5w + .9w^2$, Number of instruments: 2, Number of samples: 1000, Training steps: 400, Number of critics: 50, Kernel radius: 50 data points, Critic jitter: yes
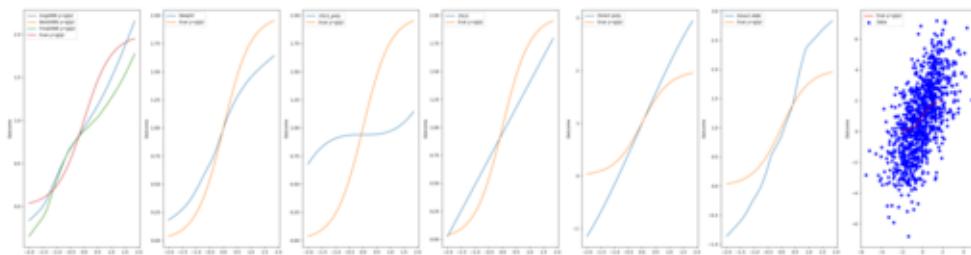
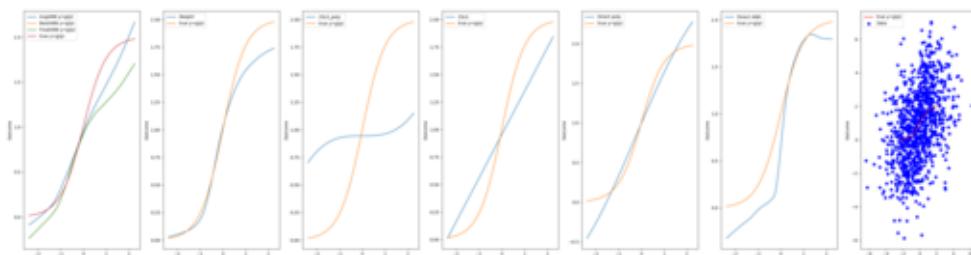(a) Strength 0.5



(b) Strength 0.7



(c) Strength 0.9

*Figure 30.* Examples of fitted functions via each method (in order from left to right: AGMM (best, final, avg), DeepIV, 2SLSPoly, 2SLS, DirectPoly, DirectNN, data points). $h_0(w) = -1.5w + .9w^2$, Number of instruments: 5, Number of samples: 1000, Training steps: 400, Number of critics: 50, Kernel radius: 50 data points, Critic jitter: yes

## C.2.2. THIRD DEGREE POLYNOMIAL



*Figure 31.* Median and $10 - 90$ percentiles of $R^2$ across 100 experiments as a function of number of instruments and instrument strength $\gamma$. $h_0(w) = -1.5w + .9w^2 + w^3$. Number of samples: 1000, Training steps: 400, Number of critics: 50, Kernel radius: 50 data points, Critic jitter: yes

(a) Strength 0.5



(b) Strength 0.7



(c) Strength 0.9

*Figure 32.* Examples of fitted functions via each method (in order from left to right: AGMM (best, final, avg), DeepIV, 2SLSPoly, 2SLS, DirectPoly, DirectNN, data points). $h_0(w) = -1.5w + .9w^2 + w^3$, Number of instruments: 2, Number of samples: 1000, Training steps: 400, Number of critics: 50, Kernel radius: 50 data points, Critic jitter: yes
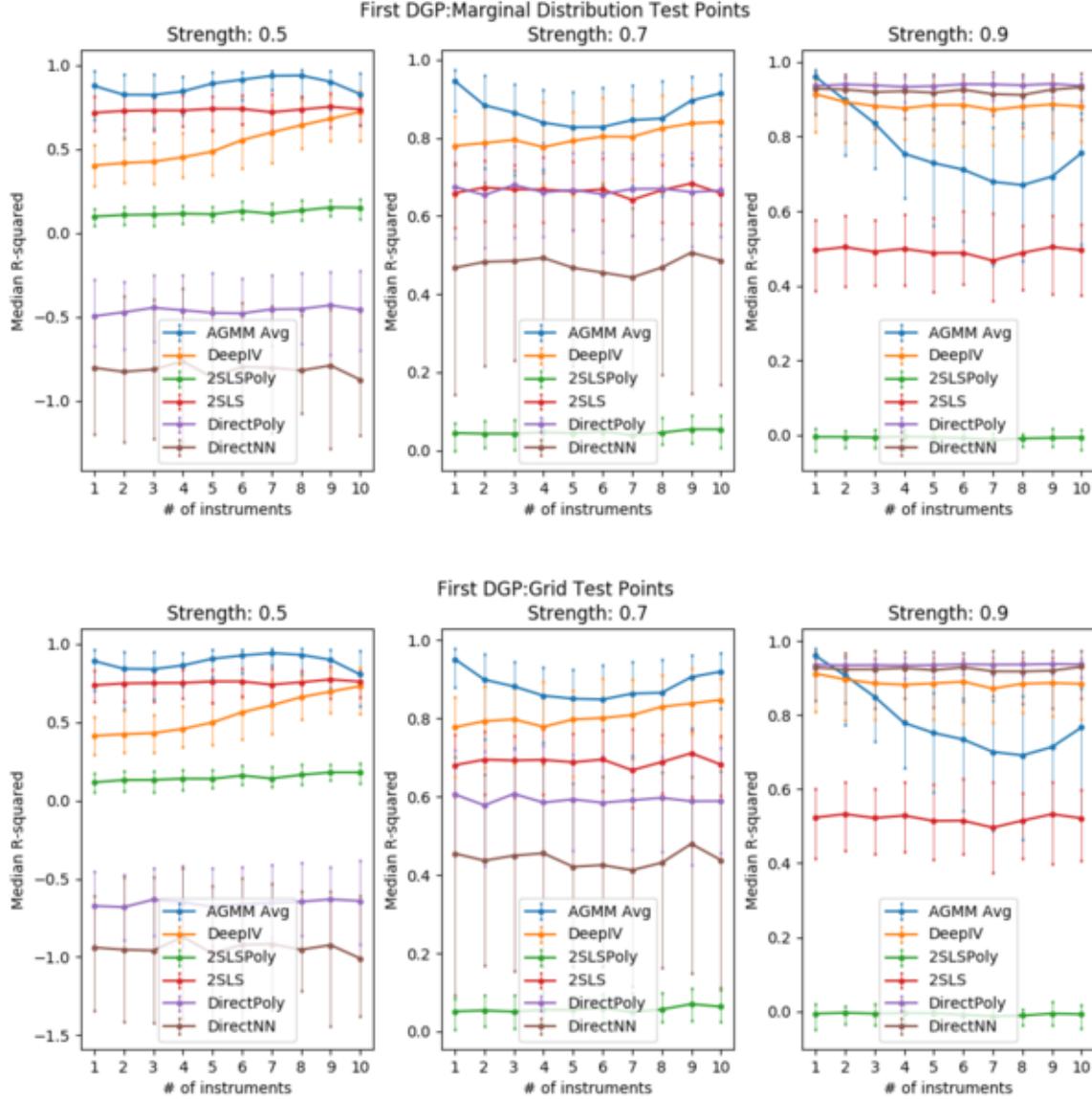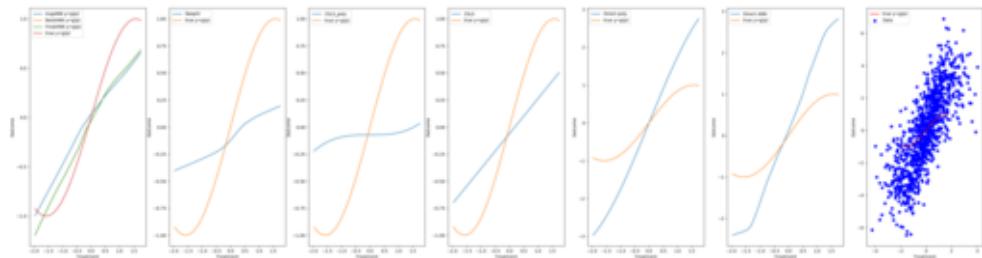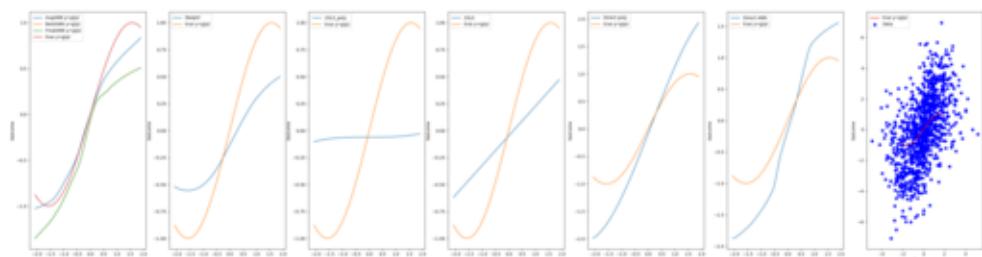
(a) Strength 0.5



(b) Strength 0.7



(c) Strength 0.9

*Figure 33.* Examples of fitted functions via each method (in order from left to right: AGMM (best, final, avg), DeepIV, 2SLSPoly, 2SLS, DirectPoly, DirectNN, data points). $h_0(w) = -1.5w + .9w^2 + w^3$, Number of instruments: 5, Number of samples: 1000, Training steps: 400, Number of critics: 50, Kernel radius: 50 data points, Critic jitter: yes

## C.2.3. ABSOLUTE VALUE



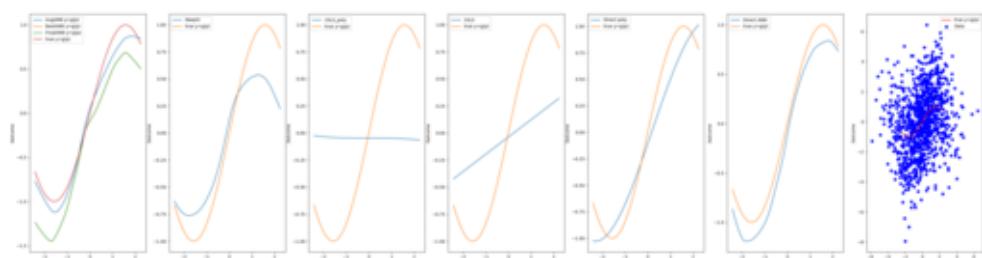*Figure 34.* Median and $10 - 90$ percentiles of $R^2$ across 100 experiments as a function of number of instruments and instrument strength $\gamma$. $h_0(w) = |w|$. Number of samples: 1000, Training steps: 400, Number of critics: 50, Kernel radius: 50 data points, Critic jitter: yes

(a) Strength 0.5



(b) Strength 0.7



(c) Strength 0.9

*Figure 35.* Examples of fitted functions via each method (in order from left to right: AGMM (best, final, avg), DeepIV, 2SLSPoly, 2SLS, DirectPoly, DirectNN, data points). $h_0(w) = |w|$, Number of instruments: 2, Number of samples: 1000, Training steps: 400, Number of critics: 50, Kernel radius: 50 data points, Critic jitter: yes
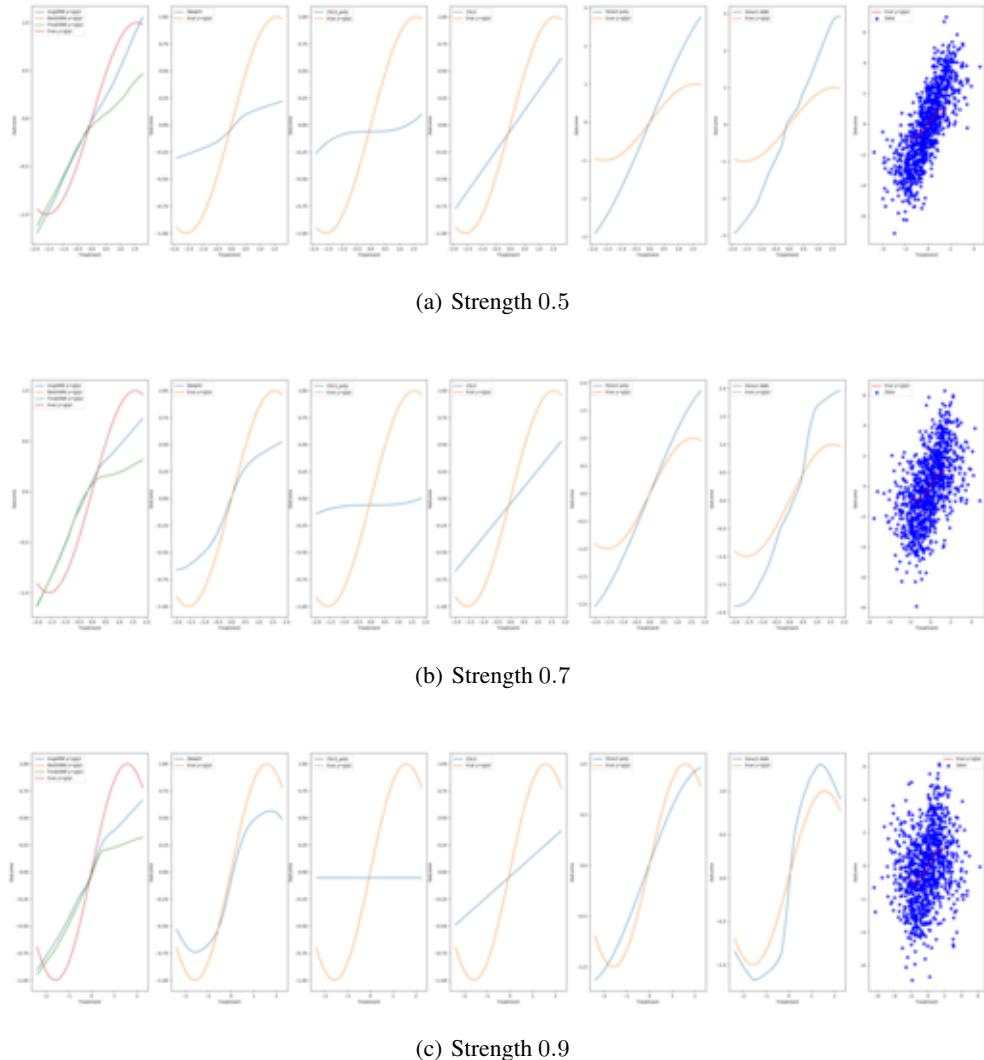
(a) Strength 0.5



(b) Strength 0.7



(c) Strength 0.9

*Figure 36.* Examples of fitted functions via each method (in order from left to right: AGMM (best, final, avg), DeepIV, 2SLSPoly, 2SLS, DirectPoly, DirectNN, data points). $h_0(w) = |w|$, Number of instruments: 5, Number of samples: 1000, Training steps: 400, Number of critics: 50, Kernel radius: 50 data points, Critic jitter: yes
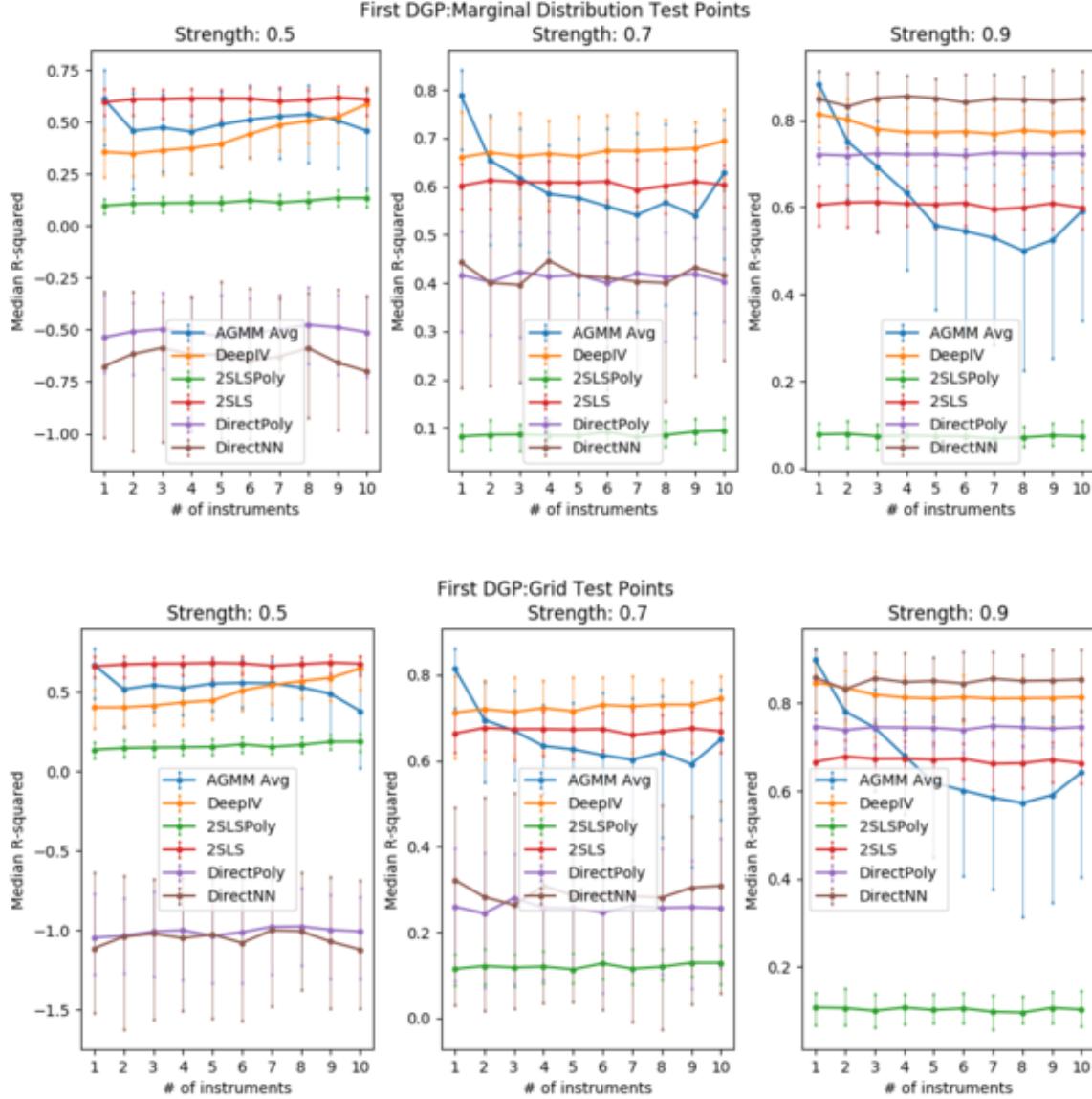
## C.2.4. IDENTIFY FUNCTION
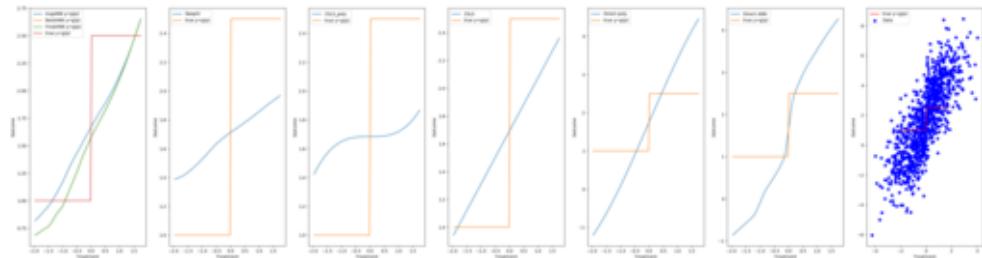


*Figure 37.* Median and $10 - 90$ percentiles of $R^2$ across 100 experiments as a function of number of instruments and instrument strength $\gamma$. $h_0(w) = x$. Number of samples: 1000, Training steps: 400, Number of critics: 50, Kernel radius: 50 data points, Critic jitter: yes

(a) Strength 0.5



(b) Strength 0.7



(c) Strength 0.9

*Figure 38.* Examples of fitted functions via each method (in order from left to right: AGMM (best, final, avg), DeepIV, 2SLSPoly, 2SLS, DirectPoly, DirectNN, data points). $h_0(w) = w$, Number of instruments: 2, Number of samples: 1000, Training steps: 400, Number of critics: 50, Kernel radius: 50 data points, Critic jitter: yes

(a) Strength 0.5



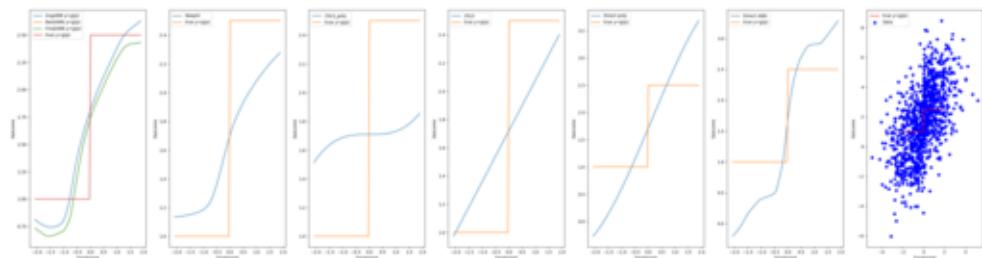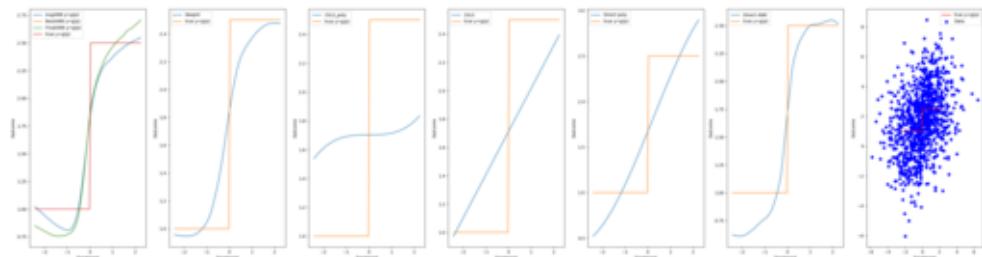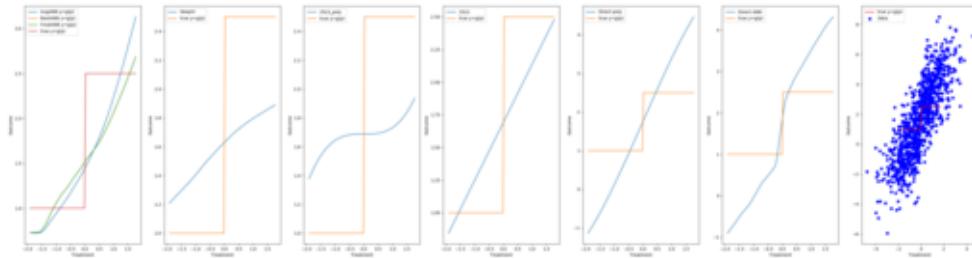(b) Strength 0.7



(c) Strength 0.9

*Figure 39.* Examples of fitted functions via each method (in order from left to right: AGMM (best, final, avg), DeepIV, 2SLSPoly, 2SLS, DirectPoly, DirectNN, data points). $h_0(w) = w$, Number of instruments: 5, Number of samples: 1000, Training steps: 400, Number of critics: 50, Kernel radius: 50 data points, Critic jitter: yes

## C.2.5. SIGMOID FUNCTION



*Figure 40.* Median and $10 - 90$ percentiles of $R^2$ across 100 experiments as a function of number of instruments and instrument strength $\gamma$. $h_0(w) = \frac{2}{1+e^{-2x}}$. Number of samples: 1000, Training steps: 400, Number of critics: 50, Kernel radius: 50 data points, Critic jitter: yes
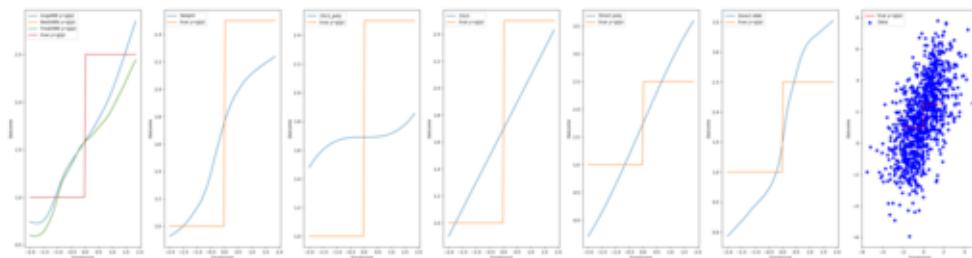
(a) Strength 0.5

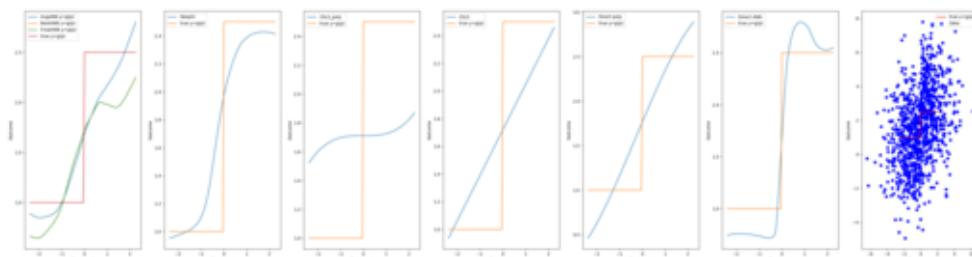

(b) Strength 0.7



(c) Strength 0.9

*Figure 41.* Examples of fitted functions via each method (in order from left to right: AGMM (best, final, avg), DeepIV, 2SLSPoly, 2SLS, DirectPoly, DirectNN, data points). $h_0(w) = \frac{2}{1+e^{-2x}}$, Number of instruments: 2, Number of samples: 1000, Training steps: 400, Number of critics: 50, Kernel radius: 50 data points, Critic jitter: yes

(a) Strength 0.5



(b) Strength 0.7



(c) Strength 0.9

*Figure 42.* Examples of fitted functions via each method (in order from left to right: AGMM (best, final, avg), DeepIV, 2SLSPoly, 2SLS, DirectPoly, DirectNN, data points). $h_0(w) = \frac{2}{1+e^{-2x}}$, Number of instruments: 5, Number of samples: 1000, Training steps: 400, Number of critics: 50, Kernel radius: 50 data points, Critic jitter: yes
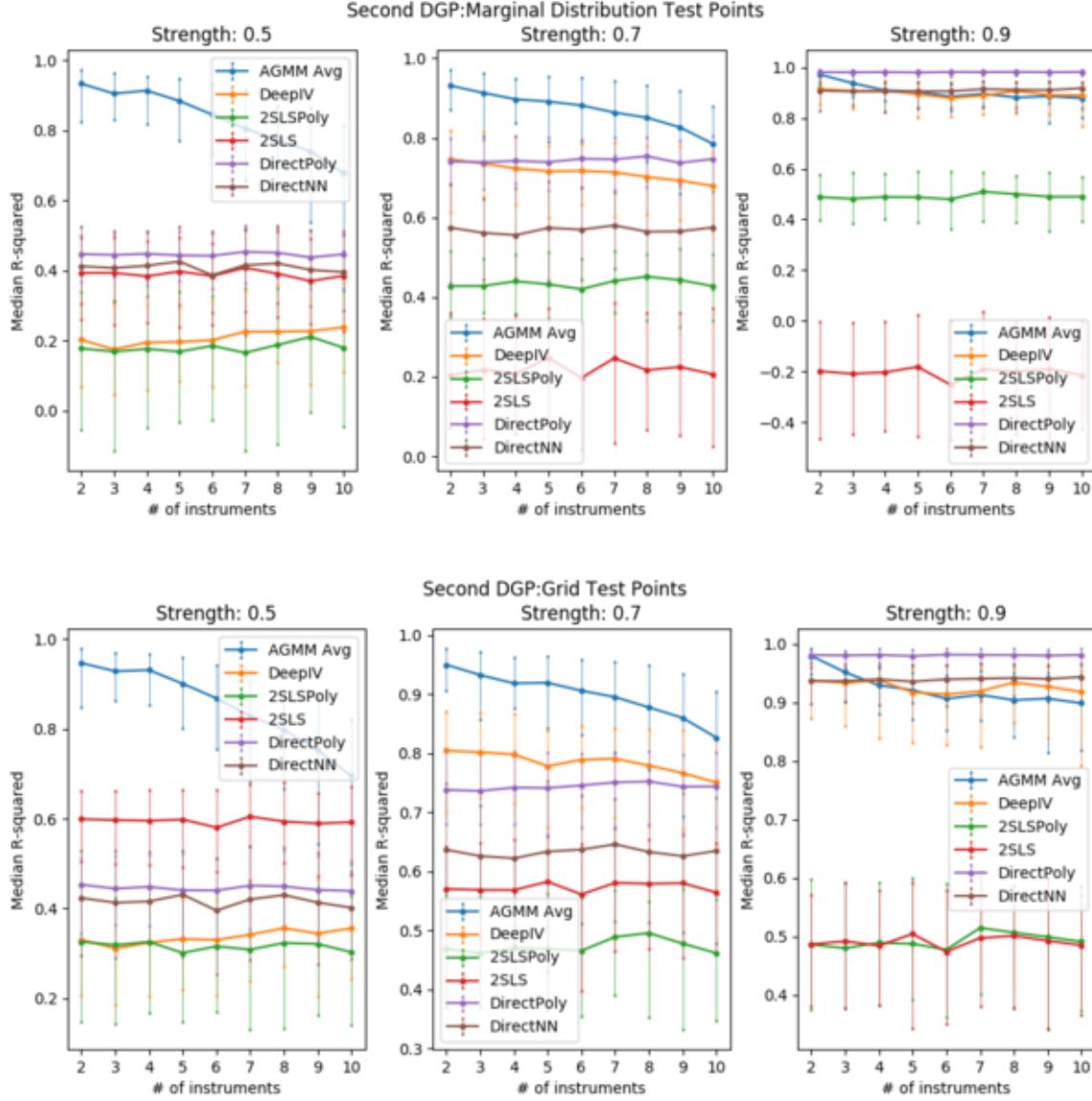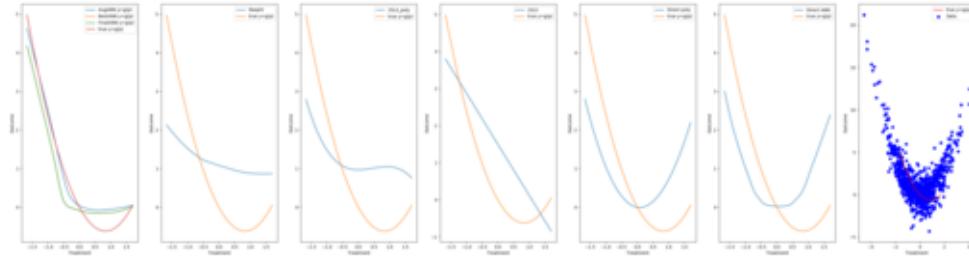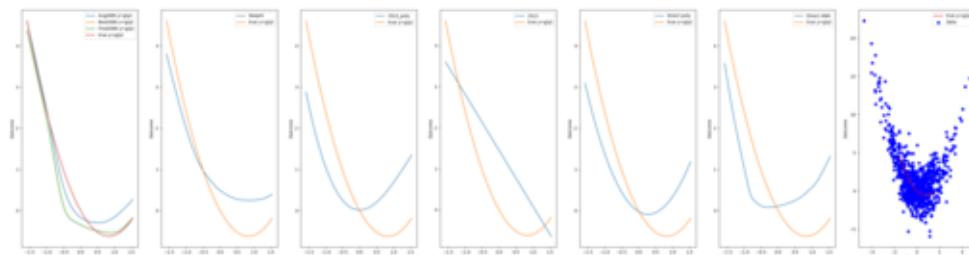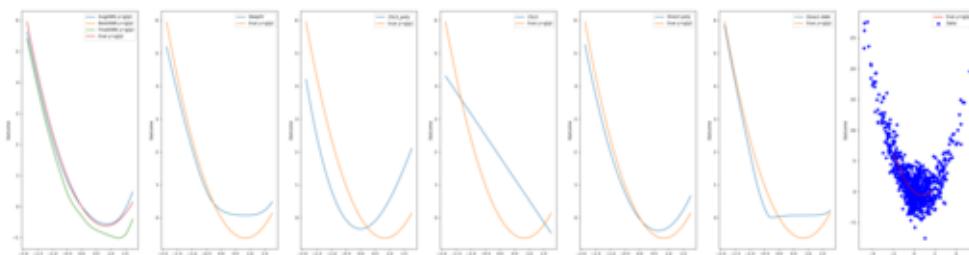
## C.2.6. Sin Function



*Figure 43.* Median and $10 - 90$ percentiles of $R^2$ across 100 experiments as a function of number of instruments and instrument strength $\gamma$. $h_0(w) = \sin(x)$. Number of samples: 1000, Training steps: 400, Number of critics: 50, Kernel radius: 50 data points, Critic jitter: yes

(a) Strength 0.5



(b) Strength 0.7



(c) Strength 0.9

*Figure 44.* Examples of fitted functions via each method (in order from left to right: AGMM (best, final, avg), DeepIV, 2SLSPoly, 2SLS, DirectPoly, DirectNN, data points). $h_0(w) = sin(x)$, Number of instruments: 2, Number of samples: 1000, Training steps: 400, Number of critics: 50, Kernel radius: 50 data points, Critic jitter: yes
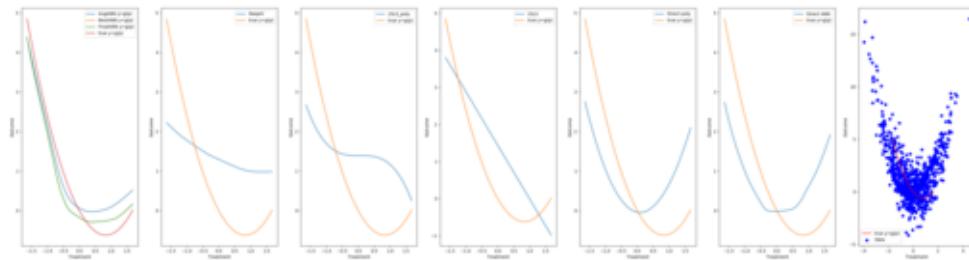
(a) Strength 0.5



(b) Strength 0.7



(c) Strength 0.9

*Figure 45.* Examples of fitted functions via each method (in order from left to right: AGMM (best, final, avg), DeepIV, 2SLSPoly, 2SLS, DirectPoly, DirectNN, data points). $h_0(w) = sin(x)$, Number of instruments: 5, Number of samples: 1000, Training steps: 400, Number of critics: 50, Kernel radius: 50 data points, Critic jitter: yes

C.2.7. STEP FUNCTION



*Figure 46.* Median and $10 - 90$ percentiles of $R^2$ across 100 experiments as a function of number of instruments and instrument strength $\gamma$. $h_0(w) = 1\{x < 0\} + 2.5 1\{x \geq 0\}$. Number of samples: 1000, Training steps: 400, Number of critics: 50, Kernel radius: 50 data points, Critic jitter: yes
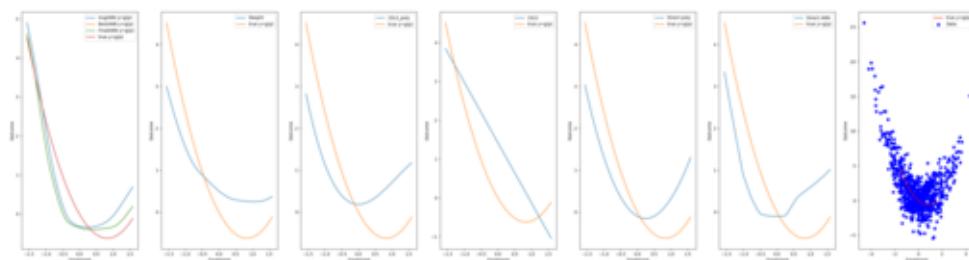
(a) Strength 0.5
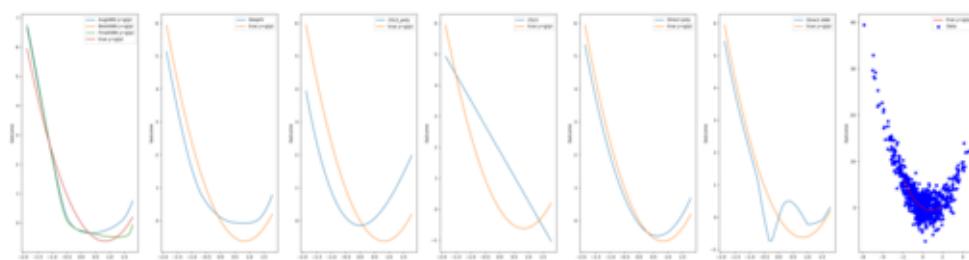


(b) Strength 0.7



(c) Strength 0.9

*Figure 47.* Examples of fitted functions via each method (in order from left to right: AGMM (best, final, avg), DeepIV, 2SLSPoly, 2SLS, DirectPoly, DirectNN, data points). $h_0(w) = 1\{x < 0\} + 2.51\{x \geq 0\}$, Number of instruments: 2, Number of samples: 1000, Training steps: 400, Number of critics: 50, Kernel radius: 50 data points, Critic jitter: yes

(a) Strength 0.5



(b) Strength 0.7



(c) Strength 0.9

*Figure 48.* Examples of fitted functions via each method (in order from left to right: AGMM (best, final, avg), DeepIV, 2SLSPoly, 2SLS, DirectPoly, DirectNN, data points). $h_0(w) = 1\{x < 0\} + 2.51\{x \geq 0\}$, Number of instruments: 5, Number of samples: 1000, Training steps: 400, Number of critics: 50, Kernel radius: 50 data points, Critic jitter: yes
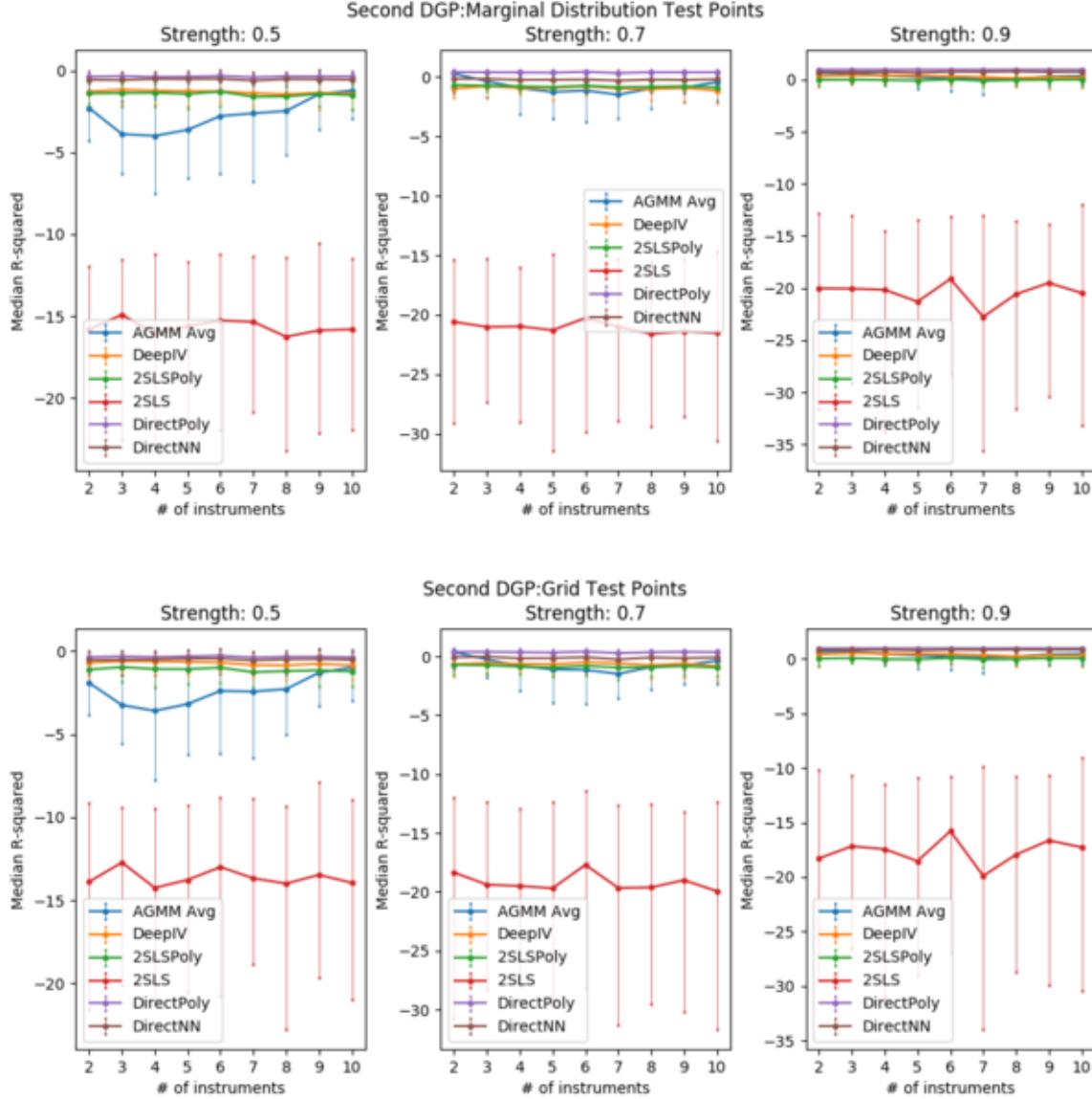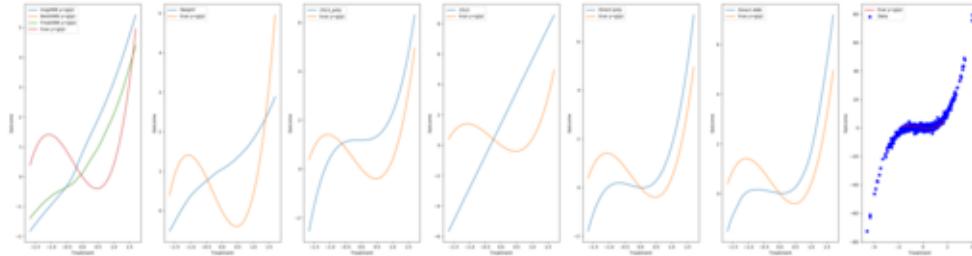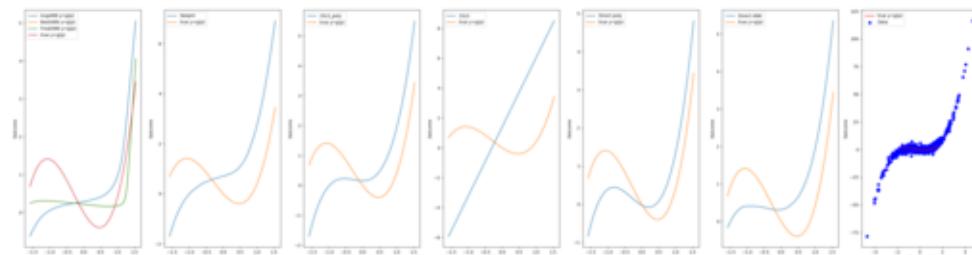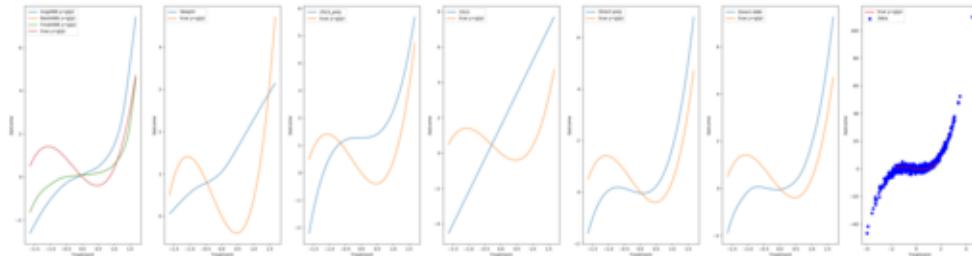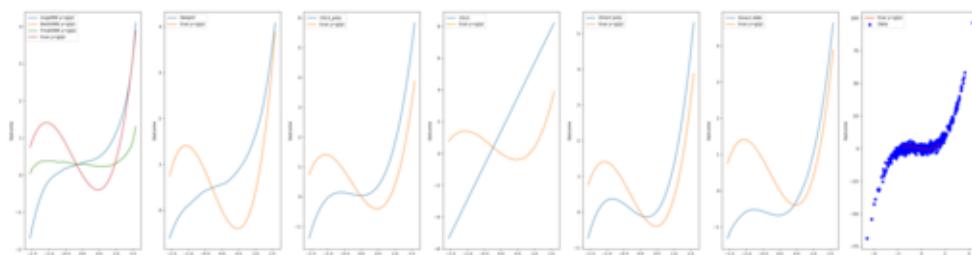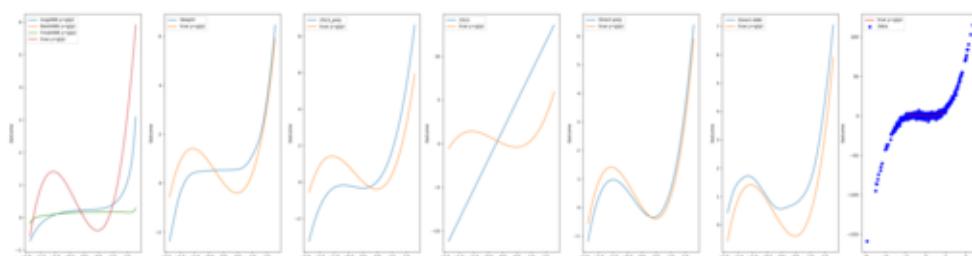
# Further Tables of Experimental Results
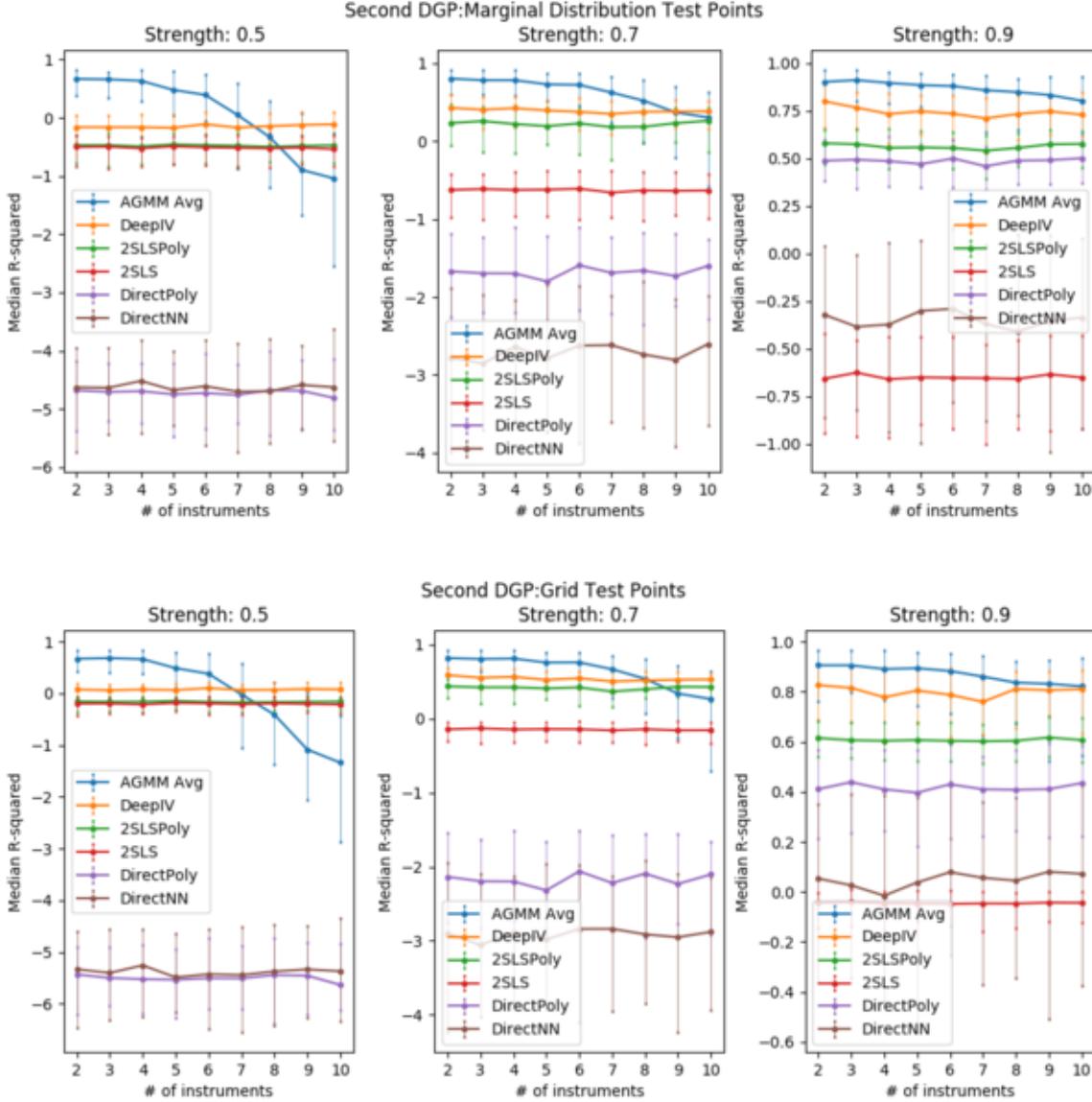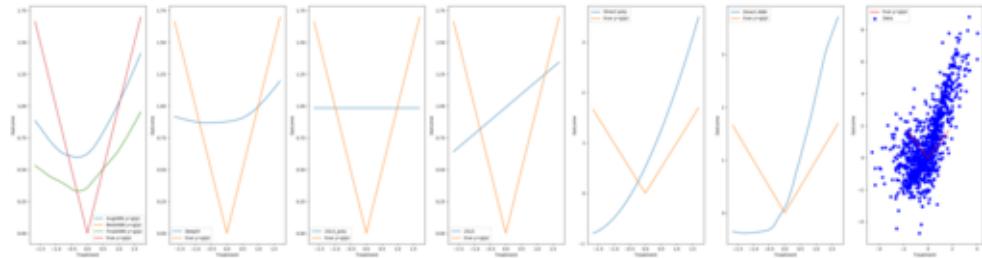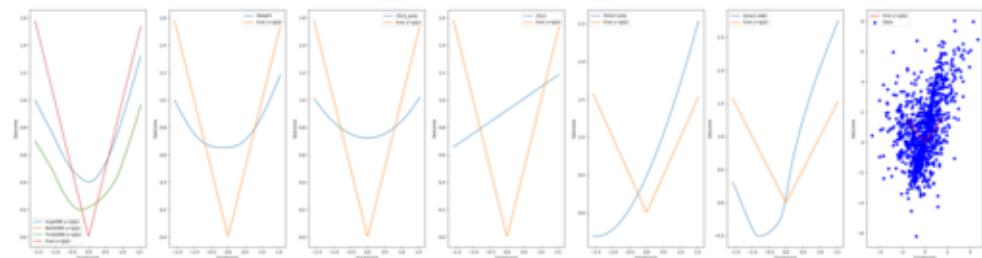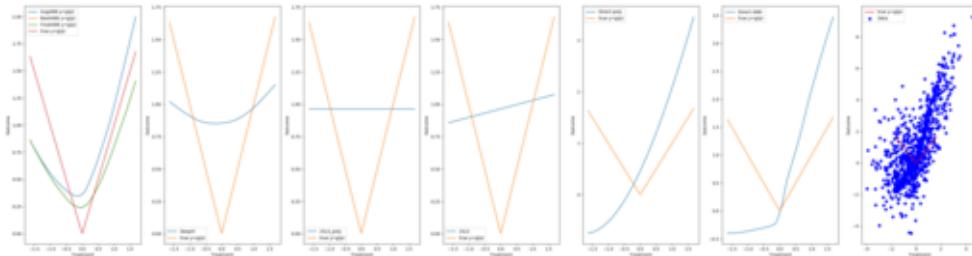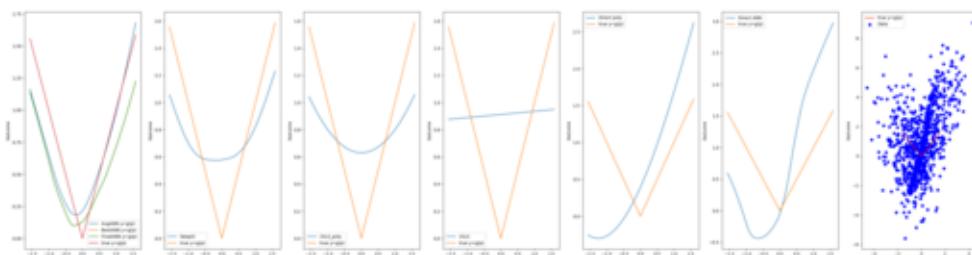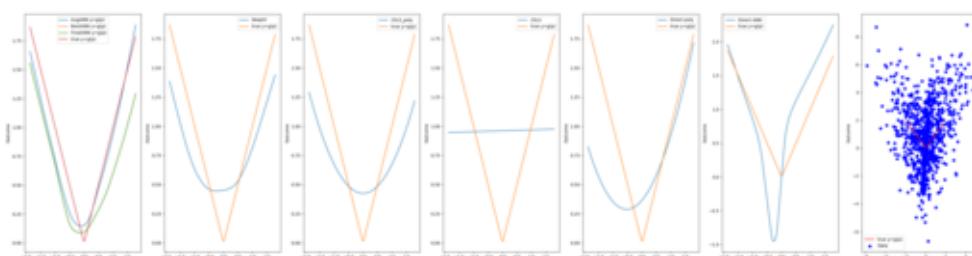
# First DGP

## D. Marginal Distribution Test Points

### D.1. Number of Instruments:1

#### D.1.1. INSTRUMENT STRENGTH: 0.5

| function | AGMM Avg | DeepIV | 2SLSpoly | 2SLS | DirectPoly | DirectNN |
|---|---|---|---|---|---|---|
| abs | **0.78** (0.46, 0.92) | -0.09 (-0.44, 0.11) | -0.44 (-0.78, -0.20) | -0.47 (-0.81, -0.23) | -2.77 (-3.27, -2.20) | -2.42 (-3.36, -1.66) |
| 2dpoly | **0.96** (0.89, 0.98) | 0.37 (0.18, 0.51) | 0.38 (0.18, 0.52) | 0.35 (0.15, 0.49) | 0.67 (0.58, 0.72) | 0.66 (0.56, 0.74) |
| sigmoid | **0.88** (0.56, 0.97) | 0.50 (0.28, 0.69) | 0.20 (0.13, 0.28) | 0.88 (0.77, 0.95) | -0.91 (-1.15, -0.59) | -1.12 (-1.73, -0.62) |
| step | **0.61** (0.34, 0.76) | 0.36 (0.19, 0.49) | 0.10 (0.05, 0.15) | 0.60 (0.52, 0.67) | -0.54 (-0.75, -0.30) | -0.68 (-1.11, -0.27) |
| 3dpoly | 0.02 (-1.86, 0.71) | -1.02 (-2.27, -0.43) | -0.21 (-1.01, 0.22) | -12.04 (-18.23, -7.07) | **0.38** (0.15, 0.59) | 0.26 (-0.27, 0.58) |
| sin | **0.88** (0.56, 0.98) | 0.40 (0.24, 0.57) | 0.10 (0.03, 0.16) | 0.72 (0.58, 0.82) | -0.50 (-0.73, -0.23) | -0.80 (-1.28, -0.41) |
| linear | 0.97 (0.88, 0.99) | 0.67 (0.49, 0.79) | 0.38 (0.31, 0.47) | **0.99** (0.98, 1.00) | 0.00 (-0.14, 0.14) | 0.00 (-0.29, 0.24) |
| rand_pw | **0.84** (-0.01, 0.98) | 0.34 (-0.36, 0.70) | 0.20 (-0.43, 0.43) | 0.54 (-0.93, 0.98) | 0.26 (-4.29, 0.88) | 0.43 (-3.39, 0.90) |

*Figure 49.* Instrument strength: 0.5, Number of instruments: 1, Number of samples: 1000, Number of experiments: 100, Training steps: 400, Number of critics: 50, Kernel radius: 50 data points, Critic jitter: 1, Test points: Distribution

#### D.1.2. INSTRUMENT STRENGTH: 0.7

| function | AGMM Avg | DeepIV | 2SLSpoly | 2SLS | DirectPoly | DirectNN |
|---|---|---|---|---|---|---|
| abs | **0.89** (0.72, 0.96) | 0.53 (0.25, 0.68) | 0.57 (0.38, 0.69) | -0.46 (-0.77, -0.24) | -0.20 (-0.56, 0.08) | 0.07 (-0.55, 0.36) |
| 2dpoly | **0.99** (0.95, 1.00) | 0.85 (0.76, 0.91) | 0.53 (0.40, 0.59) | 0.28 (0.08, 0.42) | 0.91 (0.89, 0.94) | 0.90 (0.84, 0.95) |
| sigmoid | **0.94** (0.82, 0.98) | 0.87 (0.73, 0.96) | 0.16 (0.11, 0.21) | 0.87 (0.80, 0.93) | 0.51 (0.33, 0.67) | 0.38 (-0.00, 0.66) |
| step | **0.79** (0.65, 0.87) | 0.66 (0.53, 0.79) | 0.08 (0.05, 0.11) | 0.60 (0.54, 0.66) | 0.42 (0.26, 0.53) | 0.44 (0.12, 0.62) |
| 3dpoly | 0.87 (0.75, 0.95) | -0.19 (-1.33, 0.43) | 0.22 (-0.10, 0.48) | -11.61 (-18.61, -6.98) | **0.88** (0.81, 0.93) | 0.79 (0.33, 0.93) |
| sin | **0.95** (0.84, 0.98) | 0.78 (0.64, 0.91) | 0.04 (-0.01, 0.08) | 0.66 (0.55, 0.74) | 0.67 (0.52, 0.78) | 0.47 (0.09, 0.77) |
| linear | 0.98 (0.93, 0.99) | 0.95 (0.87, 0.98) | 0.35 (0.29, 0.40) | **1.00** (0.98, 1.00) | 0.73 (0.65, 0.81) | 0.71 (0.57, 0.84) |
| rand_pw | **0.93** (0.55, 0.99) | 0.77 (0.16, 0.96) | 0.29 (-0.12, 0.64) | 0.54 (-0.89, 0.98) | 0.66 (-0.78, 0.95) | 0.83 (-0.04, 0.97) |

*Figure 50.* Instrument strength: 0.7, Number of instruments: 1, Number of samples: 1000, Number of experiments: 100, Training steps: 400, Number of critics: 50, Kernel radius: 50 data points, Critic jitter: 1, Test points: Distribution

#### D.1.3. INSTRUMENT STRENGTH: 0.9

| function | AGMM Avg | DeepIV | 2SLSpoly | 2SLS | DirectPoly | DirectNN |
|---|---|---|---|---|---|---|
| abs | **0.95** (0.87, 0.98) | 0.83 (0.68, 0.93) | 0.76 (0.69, 0.81) | -0.47 (-0.70, -0.28) | 0.66 (0.53, 0.73) | 0.91 (0.78, 0.97) |
| 2dpoly | 0.99 (0.99, 1.00) | 0.95 (0.88, 0.98) | 0.52 (0.40, 0.60) | 0.08 (-0.11, 0.25) | **1.00** (0.99, 1.00) | 0.99 (0.97, 1.00) |
| sigmoid | **0.96** (0.90, 0.99) | 0.96 (0.86, 0.99) | 0.14 (0.09, 0.19) | 0.85 (0.78, 0.90) | 0.95 (0.91, 0.96) | 0.94 (0.83, 0.98) |
| step | **0.88** (0.83, 0.92) | 0.81 (0.71, 0.87) | 0.08 (0.04, 0.11) | 0.61 (0.54, 0.67) | 0.72 (0.69, 0.74) | 0.85 (0.76, 0.93) |
| 3dpoly | 0.99 (0.97, 1.00) | 0.63 (-0.24, 0.88) | 0.63 (0.49, 0.72) | -9.00 (-13.68, -5.62) | **1.00** (0.99, 1.00) | 0.95 (0.62, 0.99) |
| sin | **0.96** (0.89, 0.98) | 0.91 (0.79, 0.96) | -0.01 (-0.05, 0.02) | 0.49 (0.37, 0.61) | 0.94 (0.89, 0.96) | 0.93 (0.81, 0.97) |
| linear | 0.99 (0.96, 0.99) | 0.98 (0.94, 0.99) | 0.33 (0.28, 0.39) | **1.00** (0.99, 1.00) | 0.98 (0.96, 1.00) | 0.97 (0.93, 0.99) |
| rand_pw | **0.97** (0.85, 1.00) | 0.92 (0.59, 0.98) | 0.39 (-0.25, 0.72) | 0.49 (-1.11, 0.98) | 0.81 (-0.01, 0.99) | 0.96 (0.81, 0.99) |

*Figure 51.* Instrument strength: 0.9, Number of instruments: 1, Number of samples: 1000, Number of experiments: 100, Training steps: 400, Number of critics: 50, Kernel radius: 50 data points, Critic jitter: 1, Test points: Distribution

## D.2. Number of Instruments:2

### D.2.1. INSTRUMENT STRENGTH: 0.5

| function | AGMM Avg | DeepIV | 2SLSpoly | 2SLS | DirectPoly | DirectNN |
|---|---|---|---|---|---|---|
| abs | **0.66 (0.31, 0.87)** | -0.09 (-0.43, 0.11) | -0.45 (-0.79, -0.17) | -0.46 (-0.81, -0.27) | -2.67 (-3.44, -2.23) | -2.39 (-3.25, -1.59) |
| 2dpoly | **0.94 (0.87, 0.98)** | 0.34 (0.20, 0.49) | 0.36 (0.11, 0.53) | 0.35 (0.18, 0.46) | 0.67 (0.58, 0.73) | 0.67 (0.56, 0.76) |
| sigmoid | 0.81 (0.48, 0.94) | 0.52 (0.30, 0.74) | 0.21 (0.12, 0.27) | **0.89 (0.78, 0.95)** | -0.87 (-1.17, -0.63) | -1.11 (-1.76, -0.60) |
| step | 0.46 (0.11, 0.66) | 0.35 (0.21, 0.48) | 0.11 (0.05, 0.15) | **0.61 (0.50, 0.68)** | -0.51 (-0.75, -0.33) | -0.62 (-1.21, -0.27) |
| 3dpoly | -2.19 (-4.30, -0.41) | -1.11 (-2.28, -0.37) | -0.21 (-1.24, 0.23) | -11.60 (-18.00, -7.49) | **0.41 (0.03, 0.60)** | 0.31 (-0.23, 0.63) |
| sin | **0.82 (0.45, 0.95)** | 0.42 (0.23, 0.60) | 0.11 (0.03, 0.16) | 0.73 (0.57, 0.83) | -0.47 (-0.72, -0.27) | -0.83 (-1.34, -0.25) |
| linear | 0.96 (0.88, 0.99) | 0.67 (0.52, 0.81) | 0.40 (0.31, 0.47) | **0.99 (0.98, 1.00)** | 0.03 (-0.16, 0.14) | 0.02 (-0.26, 0.25) |
| rand_pw | **0.61 (-0.74, 0.94)** | 0.31 (-0.33, 0.68) | 0.19 (-0.37, 0.44) | 0.58 (-0.89, 0.98) | 0.24 (-4.09, 0.89) | 0.43 (-3.71, 0.89) |

*Figure 52.* Instrument strength: 0.5, Number of instruments: 2, Number of samples: 1000, Number of experiments: 100, Training steps: 400, Number of critics: 50, Kernel radius: 50 data points, Critic jitter: 1, Test points: Distribution

### D.2.2. INSTRUMENT STRENGTH: 0.7

| function | AGMM Avg | DeepIV | 2SLSpoly | 2SLS | DirectPoly | DirectNN |
|---|---|---|---|---|---|---|
| abs | **0.84 (0.64, 0.94)** | 0.56 (0.32, 0.68) | 0.60 (0.38, 0.71) | -0.44 (-0.72, -0.28) | -0.23 (-0.59, 0.13) | -0.02 (-0.46, 0.38) |
| 2dpoly | **0.97 (0.92, 0.98)** | 0.83 (0.69, 0.91) | 0.51 (0.40, 0.61) | 0.27 (0.10, 0.41) | 0.91 (0.88, 0.94) | 0.91 (0.85, 0.95) |
| sigmoid | **0.89 (0.71, 0.97)** | 0.86 (0.70, 0.97) | 0.17 (0.11, 0.23) | 0.88 (0.80, 0.93) | 0.50 (0.31, 0.67) | 0.33 (-0.05, 0.63) |
| step | 0.65 (0.46, 0.77) | **0.67 (0.46, 0.75)** | 0.09 (0.05, 0.12) | 0.61 (0.54, 0.66) | 0.40 (0.26, 0.52) | 0.40 (0.16, 0.64) |
| 3dpoly | 0.51 (-0.22, 0.82) | -0.06 (-0.98, 0.36) | 0.28 (-0.14, 0.52) | -10.75 (-16.15, -6.68) | **0.88 (0.81, 0.93)** | 0.83 (0.34, 0.94) |
| sin | **0.88 (0.67, 0.97)** | 0.79 (0.62, 0.89) | 0.04 (-0.00, 0.08) | 0.67 (0.56, 0.76) | 0.65 (0.49, 0.80) | 0.48 (0.12, 0.72) |
| linear | 0.98 (0.93, 1.00) | 0.95 (0.88, 0.98) | 0.36 (0.29, 0.41) | **1.00 (0.98, 1.00)** | 0.73 (0.64, 0.81) | 0.70 (0.53, 0.83) |
| rand_pw | **0.87 (0.24, 0.98)** | 0.72 (0.23, 0.95) | 0.29 (-0.12, 0.69) | 0.55 (-1.02, 0.98) | 0.66 (-0.86, 0.96) | 0.80 (0.11, 0.96) |

*Figure 53.* Instrument strength: 0.7, Number of instruments: 2, Number of samples: 1000, Number of experiments: 100, Training steps: 400, Number of critics: 50, Kernel radius: 50 data points, Critic jitter: 1, Test points: Distribution

### D.2.3. INSTRUMENT STRENGTH: 0.9

| function | AGMM Avg | DeepIV | 2SLSpoly | 2SLS | DirectPoly | DirectNN |
|---|---|---|---|---|---|---|
| abs | **0.91 (0.79, 0.97)** | 0.83 (0.62, 0.92) | 0.77 (0.68, 0.82) | -0.43 (-0.74, -0.27) | 0.65 (0.55, 0.74) | 0.89 (0.79, 0.96) |
| 2dpoly | 0.98 (0.95, 0.99) | 0.95 (0.82, 0.98) | 0.52 (0.39, 0.62) | 0.09 (-0.11, 0.25) | **0.99 (0.99, 1.00)** | 0.99 (0.97, 1.00) |
| sigmoid | 0.92 (0.80, 0.98) | **0.95 (0.88, 0.99)** | 0.15 (0.09, 0.20) | 0.86 (0.79, 0.90) | 0.95 (0.91, 0.96) | 0.92 (0.79, 0.98) |
| step | 0.75 (0.60, 0.82) | 0.80 (0.69, 0.85) | 0.08 (0.04, 0.11) | 0.61 (0.54, 0.66) | 0.72 (0.69, 0.74) | **0.83 (0.71, 0.91)** |
| 3dpoly | 0.85 (0.51, 0.91) | 0.69 (-0.06, 0.88) | 0.65 (0.47, 0.76) | -8.39 (-12.35, -5.19) | **1.00 (0.99, 1.00)** | 0.96 (0.72, 1.00) |
| sin | 0.90 (0.72, 0.96) | 0.89 (0.74, 0.96) | -0.01 (-0.04, 0.02) | 0.50 (0.36, 0.60) | **0.94 (0.89, 0.96)** | 0.93 (0.82, 0.98) |
| linear | 0.99 (0.95, 1.00) | 0.98 (0.91, 0.99) | 0.34 (0.28, 0.39) | **1.00 (0.99, 1.00)** | 0.98 (0.96, 1.00) | 0.97 (0.93, 0.99) |
| rand_pw | 0.93 (0.42, 0.99) | 0.89 (0.52, 0.99) | 0.38 (-0.19, 0.75) | 0.43 (-1.04, 0.98) | 0.83 (0.00, 0.99) | **0.95 (0.84, 0.99)** |

*Figure 54.* Instrument strength: 0.9, Number of instruments: 2, Number of samples: 1000, Number of experiments: 100, Training steps: 400, Number of critics: 50, Kernel radius: 50 data points, Critic jitter: 1, Test points: Distribution

## D.3. Number of Instruments:3

### D.3.1. INSTRUMENT STRENGTH: 0.5

| function | AGMM Avg | DeepIV | 2SLSpoly | 2SLS | DirectPoly | DirectNN |
|---|---|---|---|---|---|---|
| abs | **0.69 (0.45, 0.86)** | -0.08 (-0.49, 0.15) | -0.43 (-0.82, -0.18) | -0.47 (-0.83, -0.25) | -2.71 (-3.36, -2.16) | -2.38 (-3.30, -1.79) |
| 2dpoly | **0.94 (0.85, 0.97)** | 0.33 (0.14, 0.47) | 0.38 (0.08, 0.51) | 0.34 (0.17, 0.46) | 0.66 (0.59, 0.73) | 0.66 (0.57, 0.74) |
| sigmoid | 0.81 (0.58, 0.93) | 0.53 (0.35, 0.72) | 0.21 (0.13, 0.27) | **0.89 (0.76, 0.94)** | -0.85 (-1.14, -0.55) | -1.06 (-1.88, -0.63) |
| step | 0.47 (0.17, 0.66) | 0.36 (0.22, 0.50) | 0.11 (0.05, 0.14) | **0.61 (0.50, 0.67)** | -0.50 (-0.73, -0.29) | -0.59 (-1.13, -0.27) |
| 3dpoly | -2.57 (-4.99, -0.49) | -1.01 (-2.43, -0.36) | -0.20 (-1.06, 0.14) | -12.09 (-17.91, -7.42) | **0.39 (0.09, 0.59)** | 0.26 (-0.47, 0.62) |
| sin | **0.82 (0.58, 0.96)** | 0.42 (0.25, 0.56) | 0.11 (0.03, 0.16) | 0.73 (0.56, 0.82) | -0.45 (-0.72, -0.20) | -0.81 (-1.33, -0.36) |
| linear | 0.97 (0.92, 0.99) | 0.67 (0.48, 0.80) | 0.39 (0.32, 0.46) | **0.99 (0.98, 1.00)** | 0.01 (-0.14, 0.14) | 0.00 (-0.23, 0.28) |
| rand_pw | **0.62 (-0.69, 0.93)** | 0.34 (-0.26, 0.69) | 0.18 (-0.24, 0.45) | 0.55 (-0.87, 0.98) | 0.25 (-3.86, 0.88) | 0.47 (-3.21, 0.89) |

*Figure 55.* Instrument strength: 0.5, Number of instruments: 3, Number of samples: 1000, Number of experiments: 100, Training steps: 400, Number of critics: 50, Kernel radius: 50 data points, Critic jitter: 1, Test points: Distribution

### D.3.2. INSTRUMENT STRENGTH: 0.7

| function | AGMM Avg | DeepIV | 2SLSpoly | 2SLS | DirectPoly | DirectNN |
|---|---|---|---|---|---|---|
| abs | **0.84 (0.62, 0.92)** | 0.53 (0.31, 0.68) | 0.57 (0.39, 0.71) | -0.46 (-0.83, -0.25) | -0.19 (-0.55, 0.07) | 0.00 (-0.39, 0.38) |
| 2dpoly | **0.96 (0.90, 0.98)** | 0.84 (0.73, 0.90) | 0.51 (0.39, 0.58) | 0.26 (0.11, 0.39) | 0.91 (0.87, 0.94) | 0.91 (0.86, 0.94) |
| sigmoid | **0.89 (0.69, 0.95)** | 0.87 (0.71, 0.97) | 0.17 (0.10, 0.22) | 0.88 (0.78, 0.92) | 0.52 (0.33, 0.67) | 0.39 (-0.03, 0.69) |
| step | 0.62 (0.43, 0.76) | **0.66 (0.53, 0.76)** | 0.09 (0.05, 0.12) | 0.61 (0.53, 0.66) | 0.42 (0.27, 0.53) | 0.40 (0.13, 0.64) |
| 3dpoly | 0.35 (-1.10, 0.78) | -0.21 (-0.97, 0.31) | 0.28 (-0.11, 0.46) | -10.75 (-17.88, -7.58) | **0.89 (0.82, 0.93)** | 0.81 (0.35, 0.93) |
| sin | **0.86 (0.68, 0.95)** | 0.80 (0.66, 0.91) | 0.04 (-0.01, 0.08) | 0.67 (0.52, 0.75) | 0.68 (0.51, 0.81) | 0.48 (0.12, 0.78) |
| linear | 0.98 (0.93, 1.00) | 0.94 (0.85, 0.98) | 0.35 (0.26, 0.40) | **1.00 (0.98, 1.00)** | 0.74 (0.65, 0.81) | 0.71 (0.58, 0.85) |
| rand_pw | **0.84 (0.02, 0.98)** | 0.73 (0.14, 0.95) | 0.31 (-0.12, 0.64) | 0.51 (-0.82, 0.98) | 0.70 (-0.67, 0.95) | 0.82 (-0.20, 0.97) |

*Figure 56.* Instrument strength: 0.7, Number of instruments: 3, Number of samples: 1000, Number of experiments: 100, Training steps: 400, Number of critics: 50, Kernel radius: 50 data points, Critic jitter: 1, Test points: Distribution

### D.3.3. INSTRUMENT STRENGTH: 0.9

| function | AGMM Avg | DeepIV | 2SLSpoly | 2SLS | DirectPoly | DirectNN |
|---|---|---|---|---|---|---|
| abs | 0.90 (0.75, 0.96) | 0.79 (0.62, 0.88) | 0.75 (0.70, 0.81) | -0.47 (-0.75, -0.28) | 0.66 (0.56, 0.72) | **0.91 (0.78, 0.98)** |
| 2dpoly | 0.96 (0.92, 0.99) | 0.94 (0.83, 0.97) | 0.55 (0.43, 0.69) | 0.07 (-0.08, 0.21) | **0.99 (0.99, 1.00)** | 0.99 (0.97, 1.00) |
| sigmoid | 0.90 (0.75, 0.97) | 0.94 (0.88, 0.99) | 0.14 (0.08, 0.18) | 0.85 (0.77, 0.90) | **0.95 (0.91, 0.96)** | 0.93 (0.83, 0.98) |
| step | 0.69 (0.51, 0.81) | 0.78 (0.65, 0.85) | 0.07 (0.04, 0.10) | 0.61 (0.53, 0.65) | 0.72 (0.68, 0.74) | **0.85 (0.74, 0.92)** |
| 3dpoly | 0.84 (0.37, 0.92) | 0.67 (-0.08, 0.86) | 0.67 (0.52, 0.76) | -9.21 (-14.50, -6.17) | **1.00 (0.99, 1.00)** | 0.94 (0.41, 1.00) |
| sin | 0.84 (0.67, 0.93) | 0.88 (0.75, 0.97) | -0.01 (-0.05, 0.02) | 0.49 (0.34, 0.60) | **0.94 (0.87, 0.96)** | 0.92 (0.80, 0.98) |
| linear | 0.99 (0.96, 1.00) | 0.98 (0.93, 0.99) | 0.33 (0.28, 0.41) | **1.00 (0.99, 1.00)** | 0.98 (0.96, 1.00) | 0.97 (0.93, 0.99) |
| rand_pw | 0.91 (0.12, 0.98) | 0.88 (0.45, 0.98) | 0.41 (-0.20, 0.74) | 0.42 (-1.03, 0.98) | 0.81 (-0.03, 0.99) | **0.96 (0.69, 0.99)** |

*Figure 57.* Instrument strength: 0.9, Number of instruments: 3, Number of samples: 1000, Number of experiments: 100, Training steps: 400, Number of critics: 50, Kernel radius: 50 data points, Critic jitter: 1, Test points: Distribution

## D.4. Number of Instruments:5

### D.4.1. INSTRUMENT STRENGTH: 0.5

| function | AGMM Avg | DeepIV | 2SLSpoly | 2SLS | DirectPoly | DirectNN |
|---|---|---|---|---|---|---|
| abs | **0.65 (0.32, 0.92)** | -0.06 (-0.39, 0.19) | -0.45 (-0.76, -0.08) | -0.48 (-0.82, -0.25) | -2.78 (-3.40, -2.18) | -2.50 (-3.29, -1.91) |
| 2dpoly | **0.91 (0.84, 0.96)** | 0.37 (0.20, 0.49) | 0.40 (0.18, 0.50) | 0.34 (0.14, 0.47) | 0.67 (0.59, 0.72) | 0.67 (0.55, 0.74) |
| sigmoid | 0.83 (0.58, 0.96) | 0.57 (0.39, 0.75) | 0.21 (0.14, 0.29) | **0.89 (0.76, 0.95)** | -0.87 (-1.18, -0.59) | -1.12 (-1.74, -0.59) |
| step | 0.49 (0.18, 0.68) | 0.39 (0.27, 0.53) | 0.11 (0.07, 0.16) | **0.62 (0.49, 0.67)** | -0.53 (-0.75, -0.29) | -0.62 (-1.13, -0.16) |
| 3dpoly | -2.37 (-5.82, -0.46) | -1.22 (-2.17, -0.34) | -0.34 (-1.25, 0.19) | -12.02 (-18.81, -8.05) | **0.40 (0.10, 0.58)** | 0.29 (-0.24, 0.59) |
| sin | **0.89 (0.71, 0.97)** | 0.48 (0.32, 0.65) | 0.11 (0.05, 0.18) | 0.74 (0.58, 0.83) | -0.48 (-0.74, -0.21) | -0.86 (-1.48, -0.37) |
| linear | 0.97 (0.89, 0.99) | 0.72 (0.51, 0.87) | 0.40 (0.32, 0.46) | **0.99 (0.97, 1.00)** | 0.00 (-0.16, 0.13) | 0.01 (-0.25, 0.23) |
| rand_pw | 0.56 (-0.91, 0.94) | 0.36 (-0.42, 0.72) | 0.19 (-0.43, 0.45) | **0.58 (-0.91, 0.99)** | 0.25 (-3.70, 0.89) | 0.45 (-2.79, 0.89) |

*Figure 58.* Instrument strength: 0.5, Number of instruments: 5, Number of samples: 1000, Number of experiments: 100, Training steps: 400, Number of critics: 50, Kernel radius: 50 data points, Critic jitter: 1, Test points: Distribution

### D.4.2. INSTRUMENT STRENGTH: 0.7

| function | AGMM Avg | DeepIV | 2SLSpoly | 2SLS | DirectPoly | DirectNN |
|---|---|---|---|---|---|---|
| abs | **0.79 (0.54, 0.94)** | 0.50 (0.29, 0.68) | 0.56 (0.36, 0.70) | -0.45 (-0.78, -0.25) | -0.22 (-0.55, 0.08) | 0.04 (-0.39, 0.42) |
| 2dpoly | **0.95 (0.90, 0.98)** | 0.83 (0.74, 0.89) | 0.51 (0.41, 0.60) | 0.26 (0.11, 0.41) | 0.91 (0.88, 0.94) | 0.91 (0.85, 0.95) |
| sigmoid | 0.87 (0.62, 0.95) | **0.88 (0.73, 0.94)** | 0.16 (0.11, 0.23) | 0.88 (0.78, 0.92) | 0.51 (0.34, 0.67) | 0.37 (-0.10, 0.65) |
| step | 0.58 (0.26, 0.72) | **0.66 (0.55, 0.77)** | 0.08 (0.05, 0.12) | 0.61 (0.52, 0.66) | 0.42 (0.30, 0.53) | 0.42 (0.12, 0.58) |
| 3dpoly | -0.51 (-2.26, 0.46) | -0.19 (-1.14, 0.35) | 0.23 (-0.17, 0.51) | -11.33 (-17.56, -7.67) | **0.88 (0.81, 0.93)** | 0.78 (0.44, 0.92) |
| sin | **0.83 (0.61, 0.95)** | 0.79 (0.61, 0.89) | 0.04 (-0.00, 0.09) | 0.66 (0.52, 0.75) | 0.67 (0.52, 0.80) | 0.47 (0.15, 0.74) |
| linear | 0.98 (0.91, 1.00) | 0.95 (0.88, 0.99) | 0.34 (0.29, 0.41) | **1.00 (0.98, 1.00)** | 0.73 (0.64, 0.82) | 0.72 (0.52, 0.85) |
| rand_pw | 0.78 (-0.25, 0.97) | 0.73 (0.23, 0.95) | 0.32 (-0.09, 0.62) | 0.51 (-0.96, 0.98) | 0.65 (-0.98, 0.96) | **0.82 (0.15, 0.97)** |

*Figure 59.* Instrument strength: 0.7, Number of instruments: 5, Number of samples: 1000, Number of experiments: 100, Training steps: 400, Number of critics: 50, Kernel radius: 50 data points, Critic jitter: 1, Test points: Distribution

### D.4.3. INSTRUMENT STRENGTH: 0.9

| function | AGMM Avg | DeepIV | 2SLSpoly | 2SLS | DirectPoly | DirectNN |
|---|---|---|---|---|---|---|
| abs | **0.88 (0.71, 0.97)** | 0.77 (0.64, 0.89) | 0.74 (0.67, 0.81) | -0.46 (-0.78, -0.25) | 0.64 (0.55, 0.73) | 0.87 (0.74, 0.97) |
| 2dpoly | 0.96 (0.90, 0.98) | 0.93 (0.83, 0.98) | 0.59 (0.44, 0.69) | 0.08 (-0.11, 0.28) | **0.99 (0.99, 1.00)** | 0.99 (0.97, 1.00) |
| sigmoid | 0.85 (0.66, 0.94) | 0.93 (0.83, 0.98) | 0.14 (0.10, 0.19) | 0.85 (0.76, 0.89) | **0.95 (0.91, 0.96)** | 0.91 (0.80, 0.98) |
| step | 0.56 (0.26, 0.77) | 0.77 (0.65, 0.83) | 0.07 (0.04, 0.11) | 0.61 (0.52, 0.66) | 0.72 (0.68, 0.74) | **0.85 (0.71, 0.91)** |
| 3dpoly | 0.57 (-0.28, 0.86) | 0.64 (-0.26, 0.88) | 0.70 (0.54, 0.81) | -9.28 (-13.28, -6.29) | **1.00 (0.99, 1.00)** | 0.94 (0.37, 0.99) |
| sin | 0.73 (0.52, 0.87) | 0.89 (0.75, 0.96) | -0.01 (-0.05, 0.03) | 0.49 (0.32, 0.62) | **0.94 (0.87, 0.96)** | 0.92 (0.80, 0.97) |
| linear | 0.99 (0.95, 1.00) | 0.98 (0.93, 0.99) | 0.33 (0.29, 0.41) | **1.00 (0.99, 1.00)** | 0.98 (0.96, 1.00) | 0.97 (0.91, 0.99) |
| rand_pw | 0.78 (0.01, 0.98) | 0.87 (0.47, 0.96) | 0.29 (-0.28, 0.67) | 0.39 (-1.02, 0.96) | 0.79 (-0.06, 0.98) | **0.96 (0.62, 0.99)** |

*Figure 60.* Instrument strength: 0.9, Number of instruments: 5, Number of samples: 1000, Number of experiments: 32, Training steps: 400, Number of critics: 50, Kernel radius: 50 data points, Critic jitter: 1, Test points: Distribution

## D.5. Number of Instruments:10

### D.5.1. INSTRUMENT STRENGTH: 0.5

| function | AGMM Avg | DeepIV | 2SLSpoly | 2SLS | DirectPoly | DirectNN |
|---|---|---|---|---|---|---|
| abs | -0.17 (-1.02, 0.56) | **-0.01 (-0.29, 0.18)** | -0.40 (-0.76, -0.19) | -0.46 (-0.78, -0.23) | -2.77 (-3.29, -2.24) | -2.50 (-3.30, -1.94) |
| 2dpoly | **0.83 (0.69, 0.94)** | 0.43 (0.25, 0.55) | 0.41 (0.22, 0.52) | 0.35 (0.15, 0.50) | 0.67 (0.59, 0.73) | 0.66 (0.54, 0.74) |
| sigmoid | 0.60 (0.18, 0.87) | 0.81 (0.58, 0.94) | 0.25 (0.16, 0.32) | **0.89 (0.80, 0.94)** | -0.84 (-1.18, -0.60) | -1.15 (-1.71, -0.53) |
| step | 0.46 (0.15, 0.67) | 0.59 (0.43, 0.70) | 0.13 (0.08, 0.18) | **0.61 (0.53, 0.67)** | -0.51 (-0.79, -0.30) | -0.70 (-1.09, -0.28) |
| 3dpoly | -0.49 (-2.07, 0.41) | -1.30 (-2.26, -0.41) | -0.38 (-1.43, 0.18) | -11.49 (-18.81, -7.93) | **0.37 (0.14, 0.58)** | 0.30 (-0.20, 0.58) |
| sin | **0.83 (0.58, 0.97)** | 0.72 (0.49, 0.88) | 0.15 (0.07, 0.21) | 0.74 (0.59, 0.83) | -0.46 (-0.76, -0.18) | -0.87 (-1.40, -0.38) |
| linear | 0.71 (0.45, 0.92) | 0.87 (0.74, 0.96) | 0.43 (0.34, 0.49) | **0.99 (0.98, 1.00)** | 0.02 (-0.16, 0.13) | -0.02 (-0.35, 0.28) |
| rand_pw | **0.66 (-1.10, 0.94)** | 0.48 (-0.31, 0.80) | 0.22 (-0.47, 0.45) | 0.55 (-0.94, 0.98) | 0.25 (-3.75, 0.89) | 0.49 (-2.88, 0.89) |

*Figure 61.* Instrument strength: 0.5, Number of instruments: 10, Number of samples: 1000, Number of experiments: 100, Training steps: 400, Number of critics: 50, Kernel radius: 50 data points, Critic jitter: 1, Test points: Distribution

### D.5.2. INSTRUMENT STRENGTH: 0.7

| function | AGMM Avg | DeepIV | 2SLSpoly | 2SLS | DirectPoly | DirectNN |
|---|---|---|---|---|---|---|
| abs | **0.63 (0.21, 0.89)** | 0.54 (0.23, 0.66) | 0.59 (0.39, 0.71) | -0.46 (-0.74, -0.25) | -0.20 (-0.52, 0.04) | 0.04 (-0.45, 0.30) |
| 2dpoly | **0.92 (0.83, 0.97)** | 0.83 (0.72, 0.89) | 0.52 (0.42, 0.62) | 0.28 (0.04, 0.42) | 0.92 (0.88, 0.94) | 0.91 (0.84, 0.95) |
| sigmoid | 0.87 (0.69, 0.96) | **0.91 (0.79, 0.98)** | 0.18 (0.11, 0.24) | 0.87 (0.80, 0.92) | 0.50 (0.32, 0.67) | 0.37 (0.03, 0.59) |
| step | 0.63 (0.39, 0.77) | **0.69 (0.59, 0.77)** | 0.09 (0.05, 0.13) | 0.60 (0.54, 0.66) | 0.40 (0.26, 0.54) | 0.42 (0.14, 0.66) |
| 3dpoly | 0.26 (-0.89, 0.66) | -0.32 (-1.25, 0.29) | 0.22 (-0.13, 0.51) | -11.29 (-17.23, -7.29) | **0.88 (0.80, 0.93)** | 0.82 (0.25, 0.94) |
| sin | **0.91 (0.78, 0.98)** | 0.84 (0.70, 0.92) | 0.05 (-0.00, 0.10) | 0.66 (0.55, 0.74) | 0.67 (0.50, 0.79) | 0.49 (0.06, 0.72) |
| linear | 0.94 (0.83, 0.99) | 0.97 (0.90, 0.99) | 0.36 (0.28, 0.43) | **1.00 (0.98, 1.00)** | 0.74 (0.66, 0.81) | 0.71 (0.54, 0.85) |
| rand_pw | 0.78 (-0.32, 0.97) | 0.72 (0.11, 0.93) | 0.33 (-0.11, 0.66) | 0.50 (-0.84, 0.98) | 0.66 (-0.78, 0.96) | **0.80 (-0.16, 0.96)** |

*Figure 62.* Instrument strength: 0.7, Number of instruments: 10, Number of samples: 1000, Number of experiments: 100, Training steps: 400, Number of critics: 50, Kernel radius: 50 data points, Critic jitter: 1, Test points: Distribution

### D.5.3. INSTRUMENT STRENGTH: 0.9

| function | AGMM Avg | DeepIV | 2SLSpoly | 2SLS | DirectPoly | DirectNN |
|---|---|---|---|---|---|---|
| abs | 0.84 (0.66, 0.94) | 0.77 (0.56, 0.88) | 0.75 (0.68, 0.81) | -0.48 (-0.72, -0.28) | 0.66 (0.57, 0.73) | **0.91 (0.78, 0.97)** |
| 2dpoly | 0.95 (0.87, 0.99) | 0.94 (0.83, 0.97) | 0.63 (0.48, 0.75) | 0.08 (-0.15, 0.26) | **1.00 (0.99, 1.00)** | 0.99 (0.97, 1.00) |
| sigmoid | 0.85 (0.58, 0.96) | 0.95 (0.87, 0.99) | 0.14 (0.09, 0.19) | 0.84 (0.78, 0.89) | **0.95 (0.90, 0.96)** | 0.93 (0.82, 0.98) |
| step | 0.59 (0.24, 0.80) | 0.77 (0.64, 0.84) | 0.07 (0.03, 0.11) | 0.60 (0.54, 0.65) | 0.72 (0.69, 0.74) | **0.85 (0.75, 0.92)** |
| 3dpoly | 0.76 (0.17, 0.90) | 0.55 (-0.24, 0.83) | 0.75 (0.60, 0.86) | -8.82 (-14.23, -5.54) | **1.00 (0.99, 1.00)** | 0.96 (0.68, 0.99) |
| sin | 0.76 (0.33, 0.89) | 0.88 (0.75, 0.95) | -0.01 (-0.05, 0.02) | 0.50 (0.37, 0.58) | **0.94 (0.88, 0.96)** | 0.93 (0.81, 0.98) |
| linear | 0.98 (0.91, 1.00) | 0.98 (0.94, 1.00) | 0.36 (0.29, 0.47) | **1.00 (0.99, 1.00)** | 0.98 (0.96, 1.00) | 0.97 (0.93, 0.99) |
| rand_pw | 0.85 (-0.64, 0.98) | 0.84 (0.20, 0.98) | 0.50 (-0.31, 0.71) | 0.20 (-0.93, 0.96) | 0.79 (-0.09, 0.98) | **0.96 (0.64, 0.99)** |

*Figure 63.* Instrument strength: 0.9, Number of instruments: 10, Number of samples: 1000, Number of experiments: 24, Training steps: 400, Number of critics: 50, Kernel radius: 50 data points, Critic jitter: 1, Test points: Distribution

# E. Grid Test Points

## E.1. Number of Instruments:1

### E.1.1. INSTRUMENT STRENGTH: 0.5

| function | AGMM Avg | DeepIV | 2SLSpoly | 2SLS | DirectPoly | DirectNN |
|---|---|---|---|---|---|---|
| abs | **0.79 (0.49, 0.93)** | 0.11 (-0.09, 0.24) | -0.16 (-0.34, -0.02) | -0.18 (-0.38, -0.06) | -3.25 (-3.77, -2.65) | -2.90 (-3.78, -2.15) |
| 2dpoly | **0.97 (0.91, 0.99)** | 0.49 (0.32, 0.60) | 0.47 (0.28, 0.59) | 0.56 (0.42, 0.65) | 0.67 (0.57, 0.72) | 0.66 (0.55, 0.74) |
| sigmoid | **0.90 (0.60, 0.98)** | 0.53 (0.30, 0.71) | 0.23 (0.16, 0.33) | 0.89 (0.79, 0.96) | -1.16 (-1.43, -0.82) | -1.26 (-1.91, -0.75) |
| step | **0.67 (0.41, 0.79)** | 0.40 (0.22, 0.55) | 0.14 (0.08, 0.21) | 0.66 (0.58, 0.73) | -1.05 (-1.35, -0.74) | -1.11 (-1.63, -0.59) |
| 3dpoly | 0.21 (-1.30, 0.78) | -0.34 (-1.31, -0.00) | 0.03 (-0.59, 0.40) | -9.57 (-15.19, -5.14) | **0.44 (0.20, 0.65)** | 0.35 (-0.16, 0.65) |
| sin | **0.89 (0.58, 0.98)** | 0.41 (0.25, 0.58) | 0.12 (0.04, 0.19) | 0.74 (0.59, 0.85) | -0.67 (-0.96, -0.40) | -0.94 (-1.39, -0.55) |
| linear | 0.97 (0.90, 0.99) | 0.67 (0.50, 0.79) | 0.43 (0.35, 0.52) | **1.00 (0.98, 1.00)** | 0.00 (-0.14, 0.14) | 0.00 (-0.27, 0.25) |
| rand_pw | **0.86 (0.04, 0.98)** | 0.41 (-0.14, 0.71) | 0.26 (-0.15, 0.48) | 0.61 (-0.69, 0.99) | 0.26 (-3.73, 0.88) | 0.38 (-3.61, 0.88) |

*Figure 64.* Instrument strength: 0.5, Number of instruments: 1, Number of samples: 1000, Number of experiments: 100, Training steps: 400, Number of critics: 50, Kernel radius: 50 data points, Critic jitter: 1, Test points: Grid

### E.1.2. INSTRUMENT STRENGTH: 0.7

| function | AGMM Avg | DeepIV | 2SLSpoly | 2SLS | DirectPoly | DirectNN |
|---|---|---|---|---|---|---|
| abs | **0.90 (0.75, 0.97)** | 0.61 (0.43, 0.71) | 0.63 (0.49, 0.72) | -0.17 (-0.36, -0.05) | -0.28 (-0.67, 0.02) | -0.08 (-0.73, 0.25) |
| 2dpoly | **0.99 (0.96, 1.00)** | 0.89 (0.81, 0.93) | 0.56 (0.42, 0.64) | 0.51 (0.33, 0.61) | 0.91 (0.89, 0.94) | 0.91 (0.85, 0.95) |
| sigmoid | **0.95 (0.85, 0.98)** | 0.89 (0.75, 0.97) | 0.20 (0.14, 0.26) | 0.89 (0.82, 0.94) | 0.43 (0.21, 0.60) | 0.33 (-0.09, 0.64) |
| step | **0.81 (0.70, 0.88)** | 0.71 (0.59, 0.82) | 0.12 (0.07, 0.16) | 0.66 (0.59, 0.72) | 0.26 (0.05, 0.43) | 0.32 (-0.05, 0.54) |
| 3dpoly | 0.89 (0.78, 0.96) | 0.15 (-0.73, 0.59) | 0.37 (0.09, 0.59) | -8.53 (-14.56, -4.68) | **0.90 (0.83, 0.94)** | 0.83 (0.46, 0.94) |
| sin | **0.95 (0.85, 0.98)** | 0.78 (0.64, 0.92) | 0.05 (-0.01, 0.09) | 0.68 (0.57, 0.77) | 0.61 (0.42, 0.73) | 0.45 (0.05, 0.77) |
| linear | 0.98 (0.94, 0.99) | 0.96 (0.87, 0.98) | 0.39 (0.32, 0.45) | **1.00 (0.99, 1.00)** | 0.73 (0.65, 0.81) | 0.72 (0.59, 0.85) |
| rand_pw | **0.94 (0.62, 0.99)** | 0.80 (0.20, 0.96) | 0.35 (-0.13, 0.68) | 0.61 (-0.79, 0.98) | 0.68 (-0.79, 0.95) | 0.82 (-0.09, 0.96) |

*Figure 65.* Instrument strength: 0.7, Number of instruments: 1, Number of samples: 1000, Number of experiments: 100, Training steps: 400, Number of critics: 50, Kernel radius: 50 data points, Critic jitter: 1, Test points: Grid

### E.1.3. INSTRUMENT STRENGTH: 0.9

| function | AGMM Avg | DeepIV | 2SLSpoly | 2SLS | DirectPoly | DirectNN |
|---|---|---|---|---|---|---|
| abs | **0.96 (0.89, 0.98)** | 0.86 (0.74, 0.95) | 0.78 (0.71, 0.83) | -0.17 (-0.35, -0.07) | 0.69 (0.57, 0.75) | 0.91 (0.77, 0.97) |
| 2dpoly | 1.00 (0.99, 1.00) | 0.96 (0.90, 0.99) | 0.57 (0.43, 0.65) | 0.38 (0.21, 0.50) | **1.00 (0.99, 1.00)** | 0.99 (0.97, 1.00) |
| sigmoid | **0.97 (0.91, 0.99)** | 0.96 (0.87, 0.99) | 0.18 (0.11, 0.22) | 0.87 (0.81, 0.92) | 0.94 (0.90, 0.96) | 0.93 (0.85, 0.98) |
| step | **0.90 (0.84, 0.94)** | 0.85 (0.75, 0.90) | 0.11 (0.05, 0.15) | 0.67 (0.60, 0.72) | 0.75 (0.70, 0.77) | 0.86 (0.73, 0.94) |
| 3dpoly | 0.99 (0.98, 1.00) | 0.72 (0.06, 0.92) | 0.71 (0.60, 0.79) | -6.05 (-9.81, -3.65) | **1.00 (0.99, 1.00)** | 0.96 (0.73, 0.99) |
| sin | **0.96 (0.90, 0.98)** | 0.91 (0.78, 0.96) | -0.01 (-0.06, 0.03) | 0.52 (0.39, 0.64) | 0.94 (0.90, 0.96) | 0.93 (0.80, 0.98) |
| linear | 0.99 (0.97, 1.00) | 0.98 (0.95, 0.99) | 0.37 (0.32, 0.43) | **1.00 (0.99, 1.00)** | 0.98 (0.96, 1.00) | 0.98 (0.94, 1.00) |
| rand_pw | **0.98 (0.87, 1.00)** | 0.93 (0.59, 0.98) | 0.44 (-0.25, 0.74) | 0.54 (-0.83, 0.98) | 0.83 (0.03, 0.99) | 0.96 (0.81, 1.00) |

*Figure 66.* Instrument strength: 0.9, Number of instruments: 1, Number of samples: 1000, Number of experiments: 100, Training steps: 400, Number of critics: 50, Kernel radius: 50 data points, Critic jitter: 1, Test points: Grid

## E.2. Number of Instruments:2

### E.2.1. INSTRUMENT STRENGTH: 0.5

| function | AGMM Avg | DeepIV | 2SLSpoly | 2SLS | DirectPoly | DirectNN |
|---|---|---|---|---|---|---|
| abs | **0.70 (0.39, 0.88)** | 0.12 (-0.11, 0.24) | -0.16 (-0.38, 0.00) | -0.18 (-0.39, -0.06) | -3.13 (-4.15, -2.66) | -2.89 (-3.71, -1.97) |
| 2dpoly | **0.96 (0.89, 0.98)** | 0.47 (0.36, 0.59) | 0.45 (0.30, 0.59) | 0.55 (0.42, 0.63) | 0.67 (0.57, 0.74) | 0.67 (0.55, 0.76) |
| sigmoid | 0.82 (0.57, 0.94) | 0.54 (0.33, 0.75) | 0.25 (0.14, 0.33) | **0.91 (0.80, 0.96)** | -1.13 (-1.47, -0.85) | -1.28 (-1.92, -0.78) |
| step | 0.52 (0.18, 0.70) | 0.40 (0.26, 0.55) | 0.15 (0.08, 0.21) | **0.67 (0.56, 0.73)** | -1.04 (-1.43, -0.76) | -1.04 (-1.74, -0.58) |
| 3dpoly | -1.46 (-3.57, -0.09) | -0.43 (-1.13, 0.08) | 0.03 (-0.72, 0.47) | -9.20 (-15.45, -5.27) | **0.46 (0.07, 0.67)** | 0.36 (-0.16, 0.66) |
| sin | **0.84 (0.47, 0.96)** | 0.42 (0.24, 0.61) | 0.13 (0.04, 0.20) | 0.75 (0.59, 0.85) | -0.68 (-0.97, -0.45) | -0.95 (-1.56, -0.37) |
| linear | 0.97 (0.89, 1.00) | 0.67 (0.52, 0.81) | 0.45 (0.35, 0.53) | **0.99 (0.98, 1.00)** | 0.03 (-0.16, 0.14) | 0.01 (-0.25, 0.23) |
| rand_pw | **0.68 (-0.62, 0.95)** | 0.41 (-0.09, 0.70) | 0.26 (-0.17, 0.48) | 0.64 (-0.71, 0.98) | 0.22 (-3.51, 0.88) | 0.41 (-3.48, 0.89) |

*Figure 67.* Instrument strength: 0.5, Number of instruments: 2, Number of samples: 1000, Number of experiments: 100, Training steps: 400, Number of critics: 50, Kernel radius: 50 data points, Critic jitter: 1, Test points: Grid

### E.2.2. INSTRUMENT STRENGTH: 0.7

| function | AGMM Avg | DeepIV | 2SLSpoly | 2SLS | DirectPoly | DirectNN |
|---|---|---|---|---|---|---|
| abs | **0.85 (0.64, 0.94)** | 0.65 (0.47, 0.73) | 0.65 (0.46, 0.75) | -0.16 (-0.32, -0.07) | -0.34 (-0.68, 0.06) | -0.15 (-0.61, 0.30) |
| 2dpoly | **0.97 (0.94, 0.99)** | 0.87 (0.77, 0.93) | 0.55 (0.44, 0.65) | 0.50 (0.37, 0.61) | 0.91 (0.87, 0.94) | 0.91 (0.85, 0.95) |
| sigmoid | **0.90 (0.75, 0.96)** | 0.88 (0.74, 0.97) | 0.21 (0.13, 0.27) | 0.90 (0.82, 0.94) | 0.40 (0.18, 0.61) | 0.28 (-0.14, 0.59) |
| step | 0.70 (0.48, 0.80) | **0.72 (0.51, 0.80)** | 0.12 (0.07, 0.16) | 0.68 (0.60, 0.72) | 0.24 (0.05, 0.43) | 0.28 (-0.07, 0.58) |
| 3dpoly | 0.61 (-0.04, 0.87) | 0.25 (-0.47, 0.57) | 0.43 (0.06, 0.63) | -7.70 (-12.40, -4.31) | **0.90 (0.83, 0.94)** | 0.86 (0.48, 0.95) |
| sin | **0.90 (0.70, 0.97)** | 0.79 (0.63, 0.90) | 0.05 (-0.00, 0.10) | 0.69 (0.58, 0.78) | 0.58 (0.40, 0.75) | 0.44 (0.05, 0.70) |
| linear | 0.98 (0.93, 1.00) | 0.95 (0.88, 0.98) | 0.40 (0.33, 0.46) | **1.00 (0.99, 1.00)** | 0.73 (0.64, 0.81) | 0.72 (0.54, 0.83) |
| rand_pw | **0.89 (0.20, 0.98)** | 0.76 (0.28, 0.96) | 0.33 (-0.11, 0.69) | 0.60 (-0.62, 0.98) | 0.67 (-0.71, 0.95) | 0.78 (0.15, 0.96) |

*Figure 68.* Instrument strength: 0.7, Number of instruments: 2, Number of samples: 1000, Number of experiments: 100, Training steps: 400, Number of critics: 50, Kernel radius: 50 data points, Critic jitter: 1, Test points: Grid

### E.2.3. INSTRUMENT STRENGTH: 0.9

| function | AGMM Avg | DeepIV | 2SLSpoly | 2SLS | DirectPoly | DirectNN |
|---|---|---|---|---|---|---|
| abs | **0.92 (0.80, 0.97)** | 0.86 (0.71, 0.93) | 0.78 (0.71, 0.84) | -0.17 (-0.35, -0.07) | 0.68 (0.57, 0.76) | 0.89 (0.76, 0.96) |
| 2dpoly | 0.98 (0.96, 0.99) | 0.96 (0.87, 0.99) | 0.58 (0.45, 0.67) | 0.38 (0.21, 0.51) | **0.99 (0.99, 1.00)** | 0.99 (0.97, 1.00) |
| sigmoid | 0.93 (0.82, 0.98) | **0.96 (0.90, 0.99)** | 0.18 (0.11, 0.24) | 0.88 (0.81, 0.92) | 0.94 (0.88, 0.96) | 0.91 (0.79, 0.98) |
| step | 0.78 (0.63, 0.86) | **0.83 (0.73, 0.88)** | 0.11 (0.06, 0.16) | 0.68 (0.60, 0.72) | 0.74 (0.69, 0.77) | 0.83 (0.69, 0.93) |
| 3dpoly | 0.87 (0.60, 0.94) | 0.78 (0.16, 0.92) | 0.73 (0.59, 0.82) | -5.69 (-9.21, -3.38) | **1.00 (0.99, 1.00)** | 0.97 (0.79, 1.00) |
| sin | 0.91 (0.74, 0.96) | 0.90 (0.76, 0.96) | -0.00 (-0.05, 0.02) | 0.53 (0.38, 0.62) | **0.93 (0.90, 0.96)** | 0.92 (0.80, 0.98) |
| linear | 0.99 (0.96, 1.00) | 0.98 (0.92, 0.99) | 0.39 (0.32, 0.44) | **1.00 (0.99, 1.00)** | 0.98 (0.96, 1.00) | 0.97 (0.93, 0.99) |
| rand_pw | 0.93 (0.43, 0.99) | 0.92 (0.53, 0.99) | 0.42 (-0.20, 0.76) | 0.50 (-0.88, 0.98) | 0.85 (-0.05, 0.99) | **0.96 (0.86, 0.99)** |

*Figure 69.* Instrument strength: 0.9, Number of instruments: 2, Number of samples: 1000, Number of experiments: 100, Training steps: 400, Number of critics: 50, Kernel radius: 50 data points, Critic jitter: 1, Test points: Grid

### E.3. Number of Instruments:3

E.3.1. INSTRUMENT STRENGTH: 0.5

| function | AGMM Avg | DeepIV | 2SLSpoly | 2SLS | DirectPoly | DirectNN |
|---|---|---|---|---|---|---|
| abs | **0.71 (0.47, 0.87)** | 0.11 (-0.10, 0.26) | -0.15 (-0.37, -0.00) | -0.18 (-0.38, -0.07) | -3.19 (-3.91, -2.65) | -2.89 (-3.94, -2.24) |
| 2dpoly | **0.95 (0.89, 0.98)** | 0.45 (0.31, 0.56) | 0.47 (0.27, 0.57) | 0.55 (0.43, 0.64) | 0.66 (0.59, 0.73) | 0.65 (0.56, 0.75) |
| sigmoid | 0.85 (0.61, 0.94) | 0.56 (0.36, 0.73) | 0.25 (0.16, 0.33) | **0.91 (0.78, 0.95)** | -1.11 (-1.42, -0.76) | -1.24 (-2.10, -0.79) |
| step | 0.54 (0.22, 0.71) | 0.41 (0.26, 0.56) | 0.15 (0.08, 0.20) | **0.68 (0.55, 0.73)** | -1.01 (-1.37, -0.70) | -1.02 (-1.69, -0.59) |
| 3dpoly | -1.95 (-3.96, -0.15) | -0.31 (-1.31, 0.06) | 0.06 (-0.57, 0.34) | -9.14 (-14.72, -5.14) | **0.47 (0.12, 0.63)** | 0.35 (-0.23, 0.68) |
| sin | **0.84 (0.60, 0.96)** | 0.43 (0.27, 0.56) | 0.13 (0.05, 0.20) | 0.75 (0.59, 0.84) | -0.63 (-0.93, -0.35) | -0.96 (-1.56, -0.43) |
| linear | 0.98 (0.93, 0.99) | 0.66 (0.48, 0.80) | 0.43 (0.36, 0.52) | **1.00 (0.98, 1.00)** | 0.01 (-0.14, 0.14) | 0.01 (-0.22, 0.27) |
| rand_pw | **0.70 (-0.63, 0.94)** | 0.37 (-0.09, 0.70) | 0.26 (-0.10, 0.50) | 0.62 (-0.48, 0.99) | 0.26 (-3.44, 0.88) | 0.42 (-3.24, 0.89) |

*Figure 70.* Instrument strength: 0.5, Number of instruments: 3, Number of samples: 1000, Number of experiments: 100, Training steps: 400, Number of critics: 50, Kernel radius: 50 data points, Critic jitter: 1, Test points: Grid

E.3.2. INSTRUMENT STRENGTH: 0.7

| function | AGMM Avg | DeepIV | 2SLSpoly | 2SLS | DirectPoly | DirectNN |
|---|---|---|---|---|---|---|
| abs | **0.86 (0.65, 0.94)** | 0.61 (0.46, 0.73) | 0.64 (0.51, 0.74) | -0.17 (-0.38, -0.05) | -0.26 (-0.68, 0.01) | -0.12 (-0.56, 0.30) |
| 2dpoly | **0.96 (0.92, 0.99)** | 0.88 (0.78, 0.93) | 0.54 (0.42, 0.62) | 0.49 (0.38, 0.59) | 0.91 (0.87, 0.93) | 0.91 (0.86, 0.95) |
| sigmoid | **0.91 (0.77, 0.96)** | 0.88 (0.74, 0.98) | 0.20 (0.13, 0.27) | 0.90 (0.81, 0.94) | 0.43 (0.21, 0.60) | 0.33 (-0.04, 0.65) |
| step | 0.67 (0.51, 0.79) | **0.71 (0.58, 0.81)** | 0.12 (0.07, 0.16) | 0.67 (0.59, 0.72) | 0.28 (0.04, 0.42) | 0.26 (-0.05, 0.56) |
| 3dpoly | 0.47 (-0.67, 0.85) | 0.17 (-0.58, 0.55) | 0.42 (0.06, 0.61) | -7.90 (-14.04, -4.80) | **0.90 (0.84, 0.94)** | 0.85 (0.54, 0.94) |
| sin | **0.88 (0.70, 0.96)** | 0.80 (0.66, 0.91) | 0.05 (-0.01, 0.10) | 0.69 (0.55, 0.78) | 0.61 (0.41, 0.76) | 0.45 (0.09, 0.77) |
| linear | 0.98 (0.95, 1.00) | 0.94 (0.85, 0.98) | 0.39 (0.30, 0.46) | **1.00 (0.99, 1.00)** | 0.74 (0.66, 0.82) | 0.72 (0.60, 0.84) |
| rand_pw | **0.86 (0.09, 0.99)** | 0.77 (0.19, 0.95) | 0.36 (-0.14, 0.67) | 0.59 (-0.61, 0.99) | 0.70 (-0.69, 0.95) | 0.81 (-0.22, 0.97) |

*Figure 71.* Instrument strength: 0.7, Number of instruments: 3, Number of samples: 1000, Number of experiments: 100, Training steps: 400, Number of critics: 50, Kernel radius: 50 data points, Critic jitter: 1, Test points: Grid

E.3.3. INSTRUMENT STRENGTH: 0.9

| function | AGMM Avg | DeepIV | 2SLSpoly | 2SLS | DirectPoly | DirectNN |
|---|---|---|---|---|---|---|
| abs | **0.91 (0.78, 0.97)** | 0.82 (0.67, 0.90) | 0.78 (0.72, 0.83) | -0.17 (-0.35, -0.08) | 0.69 (0.58, 0.76) | 0.91 (0.77, 0.98) |
| 2dpoly | 0.97 (0.93, 0.99) | 0.95 (0.85, 0.98) | 0.59 (0.46, 0.72) | 0.36 (0.24, 0.47) | **0.99 (0.99, 1.00)** | 0.99 (0.97, 1.00) |
| sigmoid | 0.92 (0.80, 0.97) | **0.95 (0.89, 0.99)** | 0.17 (0.11, 0.22) | 0.87 (0.80, 0.92) | 0.94 (0.89, 0.96) | 0.94 (0.83, 0.98) |
| step | 0.74 (0.57, 0.84) | 0.82 (0.71, 0.88) | 0.10 (0.05, 0.14) | 0.67 (0.60, 0.71) | 0.75 (0.68, 0.77) | **0.86 (0.74, 0.93)** |
| 3dpoly | 0.88 (0.47, 0.95) | 0.76 (0.14, 0.91) | 0.75 (0.63, 0.83) | -6.29 (-10.45, -3.75) | **1.00 (0.99, 1.00)** | 0.96 (0.60, 1.00) |
| sin | 0.85 (0.68, 0.94) | 0.89 (0.76, 0.97) | -0.01 (-0.05, 0.03) | 0.52 (0.36, 0.63) | **0.94 (0.89, 0.96)** | 0.92 (0.78, 0.98) |
| linear | 0.99 (0.97, 1.00) | 0.98 (0.93, 1.00) | 0.37 (0.31, 0.46) | **1.00 (0.99, 1.00)** | 0.98 (0.96, 1.00) | 0.97 (0.93, 0.99) |
| rand_pw | 0.92 (0.17, 0.99) | 0.90 (0.45, 0.98) | 0.45 (-0.21, 0.75) | 0.52 (-0.74, 0.98) | 0.83 (-0.02, 0.99) | **0.96 (0.69, 0.99)** |

*Figure 72.* Instrument strength: 0.9, Number of instruments: 3, Number of samples: 1000, Number of experiments: 100, Training steps: 400, Number of critics: 50, Kernel radius: 50 data points, Critic jitter: 1, Test points: Grid

### E.4. Number of Instruments:5

#### E.4.1. INSTRUMENT STRENGTH: 0.5

| function | AGMM Avg | DeepIV | 2SLSpoly | 2SLS | DirectPoly | DirectNN |
|---|---|---|---|---|---|---|
| abs | **0.68 (0.43, 0.91)** | 0.15 (-0.07, 0.30) | -0.15 (-0.34, 0.01) | -0.17 (-0.39, -0.06) | -3.27 (-3.97, -2.60) | -2.97 (-3.83, -2.37) |
| 2dpoly | **0.93 (0.87, 0.97)** | 0.49 (0.35, 0.59) | 0.47 (0.30, 0.55) | 0.55 (0.43, 0.65) | 0.67 (0.58, 0.72) | 0.67 (0.54, 0.74) |
| sigmoid | 0.86 (0.64, 0.96) | 0.59 (0.43, 0.77) | 0.25 (0.18, 0.34) | **0.91 (0.78, 0.95)** | -1.12 (-1.51, -0.79) | -1.33 (-1.99, -0.78) |
| step | 0.55 (0.28, 0.71) | 0.45 (0.31, 0.58) | 0.15 (0.09, 0.21) | **0.68 (0.56, 0.73)** | -1.04 (-1.38, -0.68) | -1.03 (-1.66, -0.45) |
| 3dpoly | -1.79 (-4.77, -0.12) | -0.54 (-1.34, 0.11) | -0.03 (-0.80, 0.39) | -9.36 (-15.83, -5.84) | **0.46 (0.13, 0.65)** | 0.38 (-0.08, 0.66) |
| sin | **0.91 (0.73, 0.98)** | 0.50 (0.33, 0.67) | 0.14 (0.07, 0.21) | 0.76 (0.60, 0.85) | -0.68 (-0.99, -0.38) | -0.98 (-1.58, -0.48) |
| linear | 0.97 (0.92, 0.99) | 0.72 (0.54, 0.87) | 0.44 (0.36, 0.51) | **0.99 (0.98, 1.00)** | 0.00 (-0.15, 0.13) | 0.01 (-0.22, 0.23) |
| rand_pw | **0.62 (-0.63, 0.95)** | 0.43 (-0.19, 0.74) | 0.23 (-0.22, 0.52) | 0.60 (-0.63, 0.99) | 0.24 (-3.87, 0.88) | 0.42 (-3.23, 0.88) |

*Figure 73.* Instrument strength: 0.5, Number of instruments: 5, Number of samples: 1000, Number of experiments: 100, Training steps: 400, Number of critics: 50, Kernel radius: 50 data points, Critic jitter: 1, Test points: Grid

#### E.4.2. INSTRUMENT STRENGTH: 0.7

| function | AGMM Avg | DeepIV | 2SLSpoly | 2SLS | DirectPoly | DirectNN |
|---|---|---|---|---|---|---|
| abs | **0.81 (0.60, 0.94)** | 0.60 (0.43, 0.74) | 0.62 (0.49, 0.74) | -0.17 (-0.35, -0.07) | -0.31 (-0.68, 0.03) | -0.08 (-0.54, 0.31) |
| 2dpoly | **0.96 (0.92, 0.99)** | 0.88 (0.81, 0.92) | 0.55 (0.42, 0.64) | 0.51 (0.37, 0.61) | 0.91 (0.88, 0.94) | 0.91 (0.85, 0.95) |
| sigmoid | 0.89 (0.66, 0.95) | 0.89 (0.75, 0.95) | 0.19 (0.14, 0.27) | **0.90 (0.81, 0.94)** | 0.41 (0.24, 0.60) | 0.32 (-0.16, 0.63) |
| step | 0.63 (0.33, 0.76) | **0.72 (0.60, 0.81)** | 0.11 (0.07, 0.17) | 0.67 (0.59, 0.72) | 0.26 (0.08, 0.43) | 0.29 (-0.07, 0.50) |
| 3dpoly | -0.14 (-1.54, 0.62) | 0.16 (-0.55, 0.58) | 0.38 (0.03, 0.62) | -8.43 (-14.29, -5.05) | **0.90 (0.83, 0.94)** | 0.82 (0.56, 0.94) |
| sin | **0.85 (0.65, 0.95)** | 0.80 (0.62, 0.89) | 0.06 (-0.00, 0.11) | 0.69 (0.55, 0.78) | 0.59 (0.42, 0.75) | 0.42 (0.08, 0.72) |
| linear | 0.98 (0.94, 1.00) | 0.95 (0.88, 0.99) | 0.38 (0.32, 0.45) | **1.00 (0.99, 1.00)** | 0.73 (0.64, 0.81) | 0.73 (0.52, 0.86) |
| rand_pw | **0.82 (-0.11, 0.97)** | 0.78 (0.26, 0.95) | 0.36 (-0.15, 0.65) | 0.59 (-0.84, 0.98) | 0.66 (-1.11, 0.96) | 0.80 (0.13, 0.96) |

*Figure 74.* Instrument strength: 0.7, Number of instruments: 5, Number of samples: 1000, Number of experiments: 100, Training steps: 400, Number of critics: 50, Kernel radius: 50 data points, Critic jitter: 1, Test points: Grid

#### E.4.3. INSTRUMENT STRENGTH: 0.9

| function | AGMM Avg | DeepIV | 2SLSpoly | 2SLS | DirectPoly | DirectNN |
|---|---|---|---|---|---|---|
| abs | **0.89 (0.76, 0.97)** | 0.81 (0.69, 0.90) | 0.77 (0.70, 0.83) | -0.17 (-0.35, -0.06) | 0.67 (0.58, 0.76) | 0.87 (0.73, 0.98) |
| 2dpoly | 0.96 (0.92, 0.98) | 0.95 (0.87, 0.98) | 0.63 (0.47, 0.73) | 0.37 (0.22, 0.49) | **0.99 (0.99, 1.00)** | 0.99 (0.97, 1.00) |
| sigmoid | 0.87 (0.68, 0.95) | **0.94 (0.85, 0.98)** | 0.17 (0.11, 0.23) | 0.87 (0.79, 0.92) | 0.94 (0.90, 0.96) | 0.91 (0.81, 0.98) |
| step | 0.62 (0.33, 0.80) | 0.81 (0.69, 0.86) | 0.10 (0.07, 0.16) | 0.67 (0.59, 0.72) | 0.74 (0.69, 0.77) | **0.85 (0.70, 0.92)** |
| 3dpoly | 0.67 (-0.02, 0.90) | 0.74 (0.03, 0.93) | 0.77 (0.64, 0.87) | -6.45 (-9.62, -4.13) | **1.00 (0.99, 1.00)** | 0.95 (0.57, 1.00) |
| sin | 0.75 (0.55, 0.88) | 0.89 (0.76, 0.96) | -0.01 (-0.06, 0.03) | 0.51 (0.34, 0.64) | **0.94 (0.89, 0.96)** | 0.92 (0.80, 0.98) |
| linear | 0.99 (0.96, 1.00) | 0.98 (0.93, 0.99) | 0.37 (0.32, 0.45) | **1.00 (0.99, 1.00)** | 0.98 (0.97, 1.00) | 0.97 (0.93, 1.00) |
| rand_pw | 0.81 (0.09, 0.98) | 0.89 (0.54, 0.97) | 0.33 (-0.29, 0.69) | 0.51 (-0.89, 0.97) | 0.82 (0.00, 0.98) | **0.96 (0.63, 0.99)** |

*Figure 75.* Instrument strength: 0.9, Number of instruments: 5, Number of samples: 1000, Number of experiments: 32, Training steps: 400, Number of critics: 50, Kernel radius: 50 data points, Critic jitter: 1, Test points: Grid

## E.5. Number of Instruments:10

### E.5.1. INSTRUMENT STRENGTH: 0.5

| function | AGMM Avg | DeepIV | 2SLSpoly | 2SLS | DirectPoly | DirectNN |
|---|---|---|---|---|---|---|
| abs | -0.26 (-1.20, 0.53) | **0.16 (-0.02, 0.32)** | -0.13 (-0.36, 0.01) | -0.17 (-0.38, -0.04) | -3.24 (-3.85, -2.60) | -3.00 (-3.97, -2.38) |
| 2dpoly | **0.84 (0.70, 0.95)** | 0.52 (0.36, 0.64) | 0.47 (0.34, 0.58) | 0.56 (0.40, 0.66) | 0.67 (0.57, 0.74) | 0.66 (0.54, 0.74) |
| sigmoid | 0.58 (0.09, 0.86) | 0.84 (0.59, 0.95) | 0.30 (0.20, 0.38) | **0.91 (0.82, 0.95)** | -1.10 (-1.48, -0.83) | -1.31 (-1.90, -0.68) |
| step | 0.38 (-0.08, 0.67) | 0.65 (0.48, 0.75) | 0.19 (0.11, 0.24) | **0.68 (0.59, 0.73)** | -1.01 (-1.38, -0.72) | -1.12 (-1.60, -0.60) |
| 3dpoly | -0.23 (-1.66, 0.55) | -0.67 (-1.63, 0.07) | -0.15 (-0.90, 0.33) | -9.12 (-15.43, -5.94) | **0.43 (0.18, 0.65)** | 0.37 (-0.09, 0.65) |
| sin | **0.81 (0.53, 0.97)** | 0.73 (0.50, 0.89) | 0.18 (0.09, 0.26) | 0.76 (0.61, 0.85) | -0.64 (-0.97, -0.35) | -1.01 (-1.52, -0.44) |
| linear | 0.73 (0.50, 0.92) | 0.87 (0.74, 0.96) | 0.48 (0.38, 0.55) | **0.99 (0.98, 1.00)** | 0.02 (-0.15, 0.12) | -0.01 (-0.32, 0.27) |
| rand_pw | **0.69 (-1.27, 0.95)** | 0.51 (-0.21, 0.85) | 0.30 (-0.25, 0.51) | 0.63 (-0.66, 0.98) | 0.24 (-3.72, 0.89) | 0.43 (-3.25, 0.89) |

*Figure 76.* Instrument strength: 0.5, Number of instruments: 10, Number of samples: 1000, Number of experiments: 100, Training steps: 400, Number of critics: 50, Kernel radius: 50 data points, Critic jitter: 1, Test points: Grid

### E.5.2. INSTRUMENT STRENGTH: 0.7

| function | AGMM Avg | DeepIV | 2SLSpoly | 2SLS | DirectPoly | DirectNN |
|---|---|---|---|---|---|---|
| abs | **0.65 (0.21, 0.91)** | 0.61 (0.39, 0.71) | 0.64 (0.50, 0.73) | -0.17 (-0.36, -0.05) | -0.30 (-0.65, 0.00) | -0.08 (-0.68, 0.21) |
| 2dpoly | **0.93 (0.85, 0.98)** | 0.87 (0.77, 0.92) | 0.55 (0.44, 0.66) | 0.51 (0.32, 0.61) | 0.92 (0.88, 0.94) | 0.91 (0.84, 0.95) |
| sigmoid | 0.87 (0.70, 0.95) | **0.92 (0.80, 0.98)** | 0.22 (0.13, 0.29) | 0.89 (0.82, 0.94) | 0.41 (0.21, 0.61) | 0.33 (-0.01, 0.54) |
| step | 0.65 (0.43, 0.79) | **0.75 (0.64, 0.82)** | 0.13 (0.07, 0.18) | 0.67 (0.59, 0.72) | 0.26 (0.05, 0.44) | 0.31 (0.01, 0.58) |
| 3dpoly | 0.42 (-0.50, 0.78) | 0.04 (-0.66, 0.53) | 0.37 (0.06, 0.62) | -8.35 (-13.05, -4.90) | **0.90 (0.82, 0.94)** | 0.85 (0.40, 0.95) |
| sin | **0.92 (0.79, 0.97)** | 0.85 (0.70, 0.92) | 0.06 (0.00, 0.12) | 0.68 (0.58, 0.77) | 0.59 (0.41, 0.74) | 0.44 (0.03, 0.70) |
| linear | 0.95 (0.84, 0.99) | 0.97 (0.91, 0.99) | 0.40 (0.32, 0.48) | **1.00 (0.99, 1.00)** | 0.74 (0.66, 0.81) | 0.72 (0.54, 0.85) |
| rand_pw | **0.82 (-0.26, 0.97)** | 0.77 (0.14, 0.94) | 0.37 (-0.13, 0.69) | 0.60 (-0.67, 0.99) | 0.66 (-0.78, 0.95) | 0.80 (-0.06, 0.96) |

*Figure 77.* Instrument strength: 0.7, Number of instruments: 10, Number of samples: 1000, Number of experiments: 100, Training steps: 400, Number of critics: 50, Kernel radius: 50 data points, Critic jitter: 1, Test points: Grid

### E.5.3. INSTRUMENT STRENGTH: 0.9

| function | AGMM Avg | DeepIV | 2SLSpoly | 2SLS | DirectPoly | DirectNN |
|---|---|---|---|---|---|---|
| abs | 0.86 (0.66, 0.95) | 0.82 (0.65, 0.90) | 0.77 (0.70, 0.84) | -0.17 (-0.35, -0.06) | 0.68 (0.58, 0.76) | **0.91 (0.77, 0.98)** |
| 2dpoly | 0.96 (0.89, 0.99) | 0.95 (0.87, 0.98) | 0.67 (0.51, 0.78) | 0.38 (0.17, 0.50) | **1.00 (0.99, 1.00)** | 0.99 (0.97, 1.00) |
| sigmoid | 0.87 (0.65, 0.96) | **0.95 (0.89, 0.99)** | 0.17 (0.11, 0.23) | 0.87 (0.81, 0.91) | 0.94 (0.90, 0.96) | 0.93 (0.84, 0.98) |
| step | 0.64 (0.34, 0.81) | 0.81 (0.69, 0.87) | 0.10 (0.06, 0.15) | 0.66 (0.60, 0.72) | 0.75 (0.70, 0.77) | **0.85 (0.75, 0.94)** |
| 3dpoly | 0.83 (0.37, 0.93) | 0.65 (-0.21, 0.88) | 0.81 (0.69, 0.90) | -6.07 (-10.31, -3.50) | **1.00 (0.99, 1.00)** | 0.97 (0.78, 1.00) |
| sin | 0.77 (0.35, 0.90) | 0.88 (0.77, 0.95) | -0.01 (-0.06, 0.02) | 0.52 (0.39, 0.61) | **0.94 (0.90, 0.96)** | 0.93 (0.79, 0.98) |
| linear | 0.98 (0.93, 1.00) | 0.98 (0.94, 1.00) | 0.40 (0.32, 0.50) | **1.00 (0.99, 1.00)** | 0.98 (0.96, 1.00) | 0.98 (0.94, 1.00) |
| rand_pw | 0.88 (-0.61, 0.98) | 0.86 (0.24, 0.98) | 0.52 (-0.33, 0.73) | 0.44 (-0.90, 0.97) | 0.81 (-0.06, 0.98) | **0.96 (0.60, 0.99)** |

*Figure 78.* Instrument strength: 0.9, Number of instruments: 10, Number of samples: 1000, Number of experiments: 25, Training steps: 400, Number of critics: 50, Kernel radius: 50 data points, Critic jitter: 1, Test points: Grid

# Second DGP

## F. Marginal Distribution Test Points

### F.1. Number of Instruments:2

#### F.1.1. INSTRUMENT STRENGTH: 0.5

| function | AGMM Avg | DeepIV | 2SLSpoly | 2SLS | DirectPoly | DirectNN |
|---|---|---|---|---|---|---|
| abs | **0.67** **(0.16, 0.85)** | -0.16 (-0.49, 0.08) | -0.47 (-0.91, -0.27) | -0.50 (-0.92, -0.28) | -4.68 (-5.62, -4.03) | -4.63 (-6.03, -3.59) |
| 2dpoly | **0.93** **(0.79, 0.98)** | 0.20 (0.02, 0.33) | 0.18 (-0.15, 0.36) | 0.39 (0.24, 0.52) | 0.45 (0.35, 0.55) | 0.41 (0.27, 0.55) |
| sigmoid | 0.84 (0.35, 0.96) | 0.37 (0.20, 0.52) | 0.17 (0.08, 0.24) | **0.88** **(0.69, 0.96)** | -1.70 (-2.13, -1.29) | -2.21 (-3.09, -1.44) |
| step | 0.53 (0.07, 0.70) | 0.25 (0.14, 0.38) | 0.08 (0.01, 0.12) | **0.59** **(0.40, 0.67)** | -0.87 (-1.15, -0.55) | -1.17 (-1.65, -0.74) |
| 3dpoly | -2.29 (-4.85, -0.39) | -1.29 (-2.61, -0.53) | -1.40 (-2.16, -0.59) | -15.82 (-26.59, -10.85) | **-0.38** **(-0.81, 0.07)** | -0.57 (-1.47, 0.08) |
| sin | **0.84** **(0.40, 0.96)** | 0.29 (0.15, 0.47) | 0.08 (-0.00, 0.15) | 0.74 (0.50, 0.87) | -1.19 (-1.58, -0.82) | -1.87 (-2.50, -1.31) |
| linear | 0.95 (0.79, 0.99) | 0.49 (0.31, 0.65) | 0.35 (0.27, 0.43) | **0.99** **(0.96, 1.00)** | -0.58 (-0.76, -0.37) | -0.70 (-1.04, -0.35) |
| rand_pw | **0.67** **(-0.63, 0.97)** | 0.25 (-0.42, 0.56) | 0.11 (-0.62, 0.38) | 0.60 (-0.83, 0.98) | -0.09 (-6.81, 0.83) | 0.05 (-5.34, 0.84) |

*Figure 79.* Instrument strength: 0.5, Number of instruments: 2, Number of samples: 1000, Number of experiments: 100, Training steps: 400, Number of critics: 50, Kernel radius: 50 data points, Critic jitter: 1, Test points: Distribution

#### F.1.2. INSTRUMENT STRENGTH: 0.7

| function | AGMM Avg | DeepIV | 2SLSpoly | 2SLS | DirectPoly | DirectNN |
|---|---|---|---|---|---|---|
| abs | **0.80** **(0.49, 0.92)** | 0.43 (0.17, 0.64) | 0.23 (-0.23, 0.48) | -0.62 (-1.02, -0.34) | -1.67 (-2.53, -1.06) | -2.78 (-4.13, -1.75) |
| 2dpoly | **0.93** **(0.86, 0.98)** | 0.75 (0.56, 0.83) | 0.43 (0.34, 0.53) | 0.20 (0.02, 0.39) | 0.74 (0.66, 0.82) | 0.57 (0.39, 0.70) |
| sigmoid | **0.90** **(0.61, 0.98)** | 0.74 (0.55, 0.85) | 0.13 (0.08, 0.18) | 0.86 (0.73, 0.93) | -0.09 (-0.40, 0.19) | -1.26 (-2.18, -0.58) |
| step | **0.62** **(0.34, 0.74)** | 0.46 (0.32, 0.57) | 0.05 (0.02, 0.08) | 0.54 (0.42, 0.61) | 0.20 (0.05, 0.39) | -0.41 (-0.78, -0.03) |
| 3dpoly | 0.25 (-1.14, 0.68) | -0.99 (-2.08, -0.24) | -0.69 (-1.46, -0.15) | -20.59 (-31.26, -14.27) | **0.39** **(0.02, 0.61)** | -0.17 (-1.62, 0.35) |
| sin | **0.89** **(0.61, 0.97)** | 0.65 (0.43, 0.80) | 0.05 (-0.00, 0.09) | 0.69 (0.53, 0.81) | 0.17 (-0.11, 0.46) | -0.92 (-1.71, -0.20) |
| linear | 0.96 (0.86, 0.99) | 0.85 (0.69, 0.94) | 0.28 (0.22, 0.35) | **0.99** **(0.97, 1.00)** | 0.27 (0.05, 0.44) | -0.16 (-0.64, 0.19) |
| rand_pw | **0.86** **(0.06, 0.99)** | 0.63 (-0.09, 0.86) | 0.23 (-0.25, 0.46) | 0.53 (-1.39, 0.98) | 0.37 (-3.40, 0.92) | 0.40 (-4.61, 0.87) |

*Figure 80.* Instrument strength: 0.7, Number of instruments: 2, Number of samples: 1000, Number of experiments: 100, Training steps: 400, Number of critics: 50, Kernel radius: 50 data points, Critic jitter: 1, Test points: Distribution

#### F.1.3. INSTRUMENT STRENGTH: 0.9

| function | AGMM Avg | DeepIV | 2SLSpoly | 2SLS | DirectPoly | DirectNN |
|---|---|---|---|---|---|---|
| abs | **0.90** **(0.77, 0.97)** | 0.80 (0.64, 0.89) | 0.58 (0.46, 0.67) | -0.66 (-1.00, -0.40) | 0.49 (0.35, 0.62) | -0.32 (-1.08, 0.16) |
| 2dpoly | 0.97 (0.93, 0.99) | 0.91 (0.76, 0.96) | 0.49 (0.35, 0.62) | -0.20 (-0.53, 0.05) | **0.98** **(0.96, 0.99)** | 0.91 (0.77, 0.94) |
| sigmoid | 0.92 (0.74, 0.98) | 0.85 (0.70, 0.94) | 0.11 (0.06, 0.17) | 0.84 (0.73, 0.91) | **0.93** **(0.86, 0.97)** | -0.29 (-0.91, 0.20) |
| step | **0.60** **(0.44, 0.71)** | 0.51 (0.40, 0.59) | 0.04 (0.01, 0.06) | 0.46 (0.38, 0.51) | 0.58 (0.55, 0.60) | 0.18 (-0.48, 0.48) |
| 3dpoly | 0.65 (-0.29, 0.91) | 0.34 (-1.21, 0.68) | -0.01 (-0.71, 0.39) | -20.02 (-34.00, -10.20) | **0.97** **(0.93, 0.99)** | 0.76 (-0.50, 0.93) |
| sin | 0.90 (0.72, 0.97) | 0.79 (0.59, 0.90) | 0.01 (-0.04, 0.04) | 0.57 (0.41, 0.70) | **0.96** **(0.92, 0.98)** | -0.08 (-0.69, 0.35) |
| linear | 0.97 (0.90, 1.00) | 0.93 (0.84, 0.98) | 0.28 (0.23, 0.36) | **0.99** **(0.98, 1.00)** | 0.94 (0.88, 0.97) | 0.45 (0.16, 0.67) |
| rand_pw | **0.91** **(0.52, 0.99)** | 0.81 (0.36, 0.95) | 0.25 (-0.26, 0.55) | 0.37 (-1.46, 0.98) | 0.72 (-0.27, 0.98) | 0.73 (-1.16, 0.95) |

*Figure 81.* Instrument strength: 0.9, Number of instruments: 2, Number of samples: 1000, Number of experiments: 100, Training steps: 400, Number of critics: 50, Kernel radius: 50 data points, Critic jitter: 1, Test points: Distribution

## F.2. Number of Instruments:3

### F.2.1. INSTRUMENT STRENGTH: 0.5

| function | AGMM Avg | DeepIV | 2SLSpoly | 2SLS | DirectPoly | DirectNN |
|---|---|---|---|---|---|---|
| abs | **0.66** (0.27, 0.84) | -0.16 (-0.59, 0.09) | -0.47 (-0.89, -0.24) | -0.50 (-0.92, -0.25) | -4.70 (-5.32, -3.96) | -4.63 (-5.70, -3.78) |
| 2dpoly | **0.90** (0.81, 0.97) | 0.17 (-0.04, 0.32) | 0.17 (-0.17, 0.36) | 0.39 (0.24, 0.51) | 0.44 (0.37, 0.52) | 0.41 (0.28, 0.53) |
| sigmoid | 0.84 (0.62, 0.96) | 0.37 (0.22, 0.53) | 0.16 (0.09, 0.24) | **0.87** (0.74, 0.95) | -1.70 (-2.00, -1.38) | -2.10 (-2.81, -1.54) |
| step | 0.52 (0.33, 0.69) | 0.26 (0.10, 0.35) | 0.07 (0.03, 0.12) | **0.58** (0.45, 0.66) | -0.86 (-1.09, -0.62) | -1.15 (-1.62, -0.79) |
| 3dpoly | -3.88 (-7.16, -1.17) | -1.17 (-2.35, -0.36) | -1.37 (-2.39, -0.72) | -14.91 (-24.45, -10.07) | **-0.34 (-0.75, -0.01)** | -0.58 (-1.66, -0.02) |
| sin | **0.84** (0.60, 0.97) | 0.30 (0.16, 0.46) | 0.07 (0.01, 0.15) | 0.74 (0.54, 0.85) | -1.20 (-1.48, -0.83) | -1.78 (-2.54, -1.01) |
| linear | 0.96 (0.88, 0.99) | 0.49 (0.32, 0.69) | 0.36 (0.28, 0.42) | **0.99 (0.97, 1.00)** | -0.57 (-0.71, -0.42) | -0.59 (-1.01, -0.32) |
| rand_pw | **0.60 (-0.75, 0.95)** | 0.22 (-0.41, 0.51) | 0.11 (-0.48, 0.37) | 0.60 (-1.14, 0.98) | -0.08 (-6.56, 0.83) | 0.11 (-5.77, 0.83) |

*Figure 82.* Instrument strength: 0.5, Number of instruments: 3, Number of samples: 1000, Number of experiments: 100, Training steps: 400, Number of critics: 50, Kernel radius: 50 data points, Critic jitter: 1, Test points: Distribution

### F.2.2. INSTRUMENT STRENGTH: 0.7

| function | AGMM Avg | DeepIV | 2SLSpoly | 2SLS | DirectPoly | DirectNN |
|---|---|---|---|---|---|---|
| abs | **0.78 (0.52, 0.92)** | 0.41 (0.07, 0.62) | 0.26 (-0.31, 0.48) | -0.61 (-1.16, -0.35) | -1.69 (-2.44, -1.06) | -2.85 (-4.04, -1.73) |
| 2dpoly | **0.91 (0.78, 0.97)** | 0.74 (0.63, 0.82) | 0.43 (0.35, 0.52) | 0.22 (-0.00, 0.36) | 0.74 (0.65, 0.81) | 0.56 (0.39, 0.67) |
| sigmoid | **0.89 (0.71, 0.97)** | 0.71 (0.53, 0.86) | 0.12 (0.07, 0.18) | 0.84 (0.75, 0.92) | -0.09 (-0.42, 0.23) | -1.27 (-1.97, -0.62) |
| step | **0.60 (0.43, 0.74)** | 0.44 (0.31, 0.57) | 0.05 (0.02, 0.08) | 0.53 (0.44, 0.59) | 0.19 (0.02, 0.37) | -0.37 (-0.89, -0.04) |
| 3dpoly | -0.37 (-2.54, 0.42) | -0.79 (-2.00, -0.10) | -0.76 (-1.26, -0.35) | -21.01 (-30.47, -13.83) | **0.38 (0.08, 0.57)** | -0.17 (-1.04, 0.28) |
| sin | **0.86 (0.67, 0.96)** | 0.65 (0.49, 0.79) | 0.04 (0.00, 0.09) | 0.68 (0.55, 0.78) | 0.14 (-0.12, 0.44) | -0.87 (-1.76, -0.28) |
| linear | 0.97 (0.89, 0.99) | 0.84 (0.73, 0.93) | 0.29 (0.22, 0.35) | **0.99 (0.98, 1.00)** | 0.26 (0.05, 0.43) | -0.23 (-0.55, 0.13) |
| rand_pw | **0.83 (-0.09, 0.97)** | 0.59 (-0.01, 0.87) | 0.18 (-0.24, 0.41) | 0.52 (-1.31, 0.98) | 0.37 (-3.44, 0.92) | 0.37 (-4.54, 0.86) |

*Figure 83.* Instrument strength: 0.7, Number of instruments: 3, Number of samples: 1000, Number of experiments: 100, Training steps: 400, Number of critics: 50, Kernel radius: 50 data points, Critic jitter: 1, Test points: Distribution

### F.2.3. INSTRUMENT STRENGTH: 0.9

| function | AGMM Avg | DeepIV | 2SLSpoly | 2SLS | DirectPoly | DirectNN |
|---|---|---|---|---|---|---|
| abs | **0.91 (0.77, 0.97)** | 0.77 (0.62, 0.86) | 0.57 (0.42, 0.68) | -0.63 (-1.07, -0.40) | 0.49 (0.29, 0.59) | -0.39 (-0.96, 0.25) |
| 2dpoly | 0.94 (0.86, 0.98) | 0.91 (0.80, 0.96) | 0.48 (0.36, 0.60) | -0.21 (-0.51, 0.06) | **0.98 (0.96, 0.99)** | 0.91 (0.82, 0.95) |
| sigmoid | 0.90 (0.76, 0.97) | 0.82 (0.64, 0.90) | 0.11 (0.07, 0.17) | 0.83 (0.75, 0.90) | **0.93 (0.85, 0.97)** | -0.31 (-0.93, 0.24) |
| step | 0.57 (0.39, 0.65) | 0.48 (0.36, 0.57) | 0.04 (0.01, 0.06) | 0.45 (0.39, 0.51) | **0.58 (0.55, 0.60)** | 0.14 (-0.30, 0.43) |
| 3dpoly | 0.61 (0.12, 0.85) | 0.45 (-0.71, 0.75) | 0.02 (-0.52, 0.34) | -20.05 (-31.22, -11.81) | **0.97 (0.94, 0.99)** | 0.72 (-0.78, 0.93) |
| sin | 0.84 (0.63, 0.96) | 0.75 (0.58, 0.89) | 0.00 (-0.03, 0.05) | 0.55 (0.43, 0.68) | **0.96 (0.89, 0.98)** | -0.15 (-0.78, 0.26) |
| linear | 0.97 (0.92, 0.99) | 0.92 (0.83, 0.97) | 0.29 (0.22, 0.35) | **1.00 (0.98, 1.00)** | 0.94 (0.88, 0.97) | 0.43 (0.12, 0.67) |
| rand_pw | **0.88 (0.31, 0.99)** | 0.77 (0.22, 0.93) | 0.23 (-0.13, 0.60) | 0.38 (-1.24, 0.97) | 0.74 (-0.40, 0.98) | 0.72 (-0.66, 0.95) |

*Figure 84.* Instrument strength: 0.9, Number of instruments: 3, Number of samples: 1000, Number of experiments: 100, Training steps: 400, Number of critics: 50, Kernel radius: 50 data points, Critic jitter: 1, Test points: Distribution

## F.3. Number of Instruments:5

### F.3.1. INSTRUMENT STRENGTH: 0.5

| function | AGMM Avg | DeepIV | 2SLSpoly | 2SLS | DirectPoly | DirectNN |
|---|---|---|---|---|---|---|
| abs | **0.48** (-0.14, 0.84) | -0.17 (-0.56, 0.08) | -0.46 (-0.83, -0.24) | -0.49 (-0.94, -0.27) | -4.75 (-5.57, -4.04) | -4.67 (-5.73, -3.56) |
| 2dpoly | **0.88** (0.75, 0.96) | 0.20 (0.04, 0.32) | 0.17 (-0.05, 0.37) | 0.40 (0.21, 0.52) | 0.44 (0.34, 0.52) | 0.43 (0.27, 0.54) |
| sigmoid | 0.85 (0.62, 0.95) | 0.46 (0.24, 0.65) | 0.19 (0.10, 0.28) | **0.89** (0.72, 0.96) | -1.75 (-2.12, -1.36) | -2.18 (-2.93, -1.52) |
| step | 0.56 (0.32, 0.69) | 0.29 (0.16, 0.42) | 0.09 (0.03, 0.14) | **0.60** (0.45, 0.67) | -0.90 (-1.21, -0.64) | -1.18 (-1.67, -0.81) |
| 3dpoly | -3.62 (-7.85, -0.68) | -1.26 (-2.58, -0.58) | -1.43 (-2.59, -0.65) | -15.63 (-23.19, -9.28) | **-0.38** (-0.84, 0.05) | -0.52 (-1.41, -0.00) |
| sin | **0.93** (0.66, 0.98) | 0.37 (0.18, 0.55) | 0.10 (0.01, 0.18) | 0.76 (0.54, 0.88) | -1.23 (-1.59, -0.86) | -1.77 (-2.36, -1.18) |
| linear | 0.93 (0.75, 0.99) | 0.55 (0.41, 0.73) | 0.37 (0.29, 0.46) | **0.99** (0.95, 1.00) | -0.58 (-0.78, -0.41) | -0.68 (-1.11, -0.35) |
| rand_pw | **0.57** (-1.18, 0.95) | 0.21 (-0.30, 0.57) | 0.14 (-0.51, 0.39) | 0.54 (-0.94, 0.98) | -0.05 (-7.16, 0.84) | 0.12 (-5.11, 0.82) |

*Figure 85.* Instrument strength: 0.5, Number of instruments: 5, Number of samples: 1000, Number of experiments: 100, Training steps: 400, Number of critics: 50, Kernel radius: 50 data points, Critic jitter: 1, Test points: Distribution

### F.3.2. INSTRUMENT STRENGTH: 0.7

| function | AGMM Avg | DeepIV | 2SLSpoly | 2SLS | DirectPoly | DirectNN |
|---|---|---|---|---|---|---|
| abs | **0.73** (0.39, 0.90) | 0.39 (0.15, 0.56) | 0.19 (-0.21, 0.47) | -0.62 (-1.09, -0.35) | -1.80 (-2.41, -1.15) | -2.79 (-3.93, -1.71) |
| 2dpoly | **0.89** (0.77, 0.96) | 0.72 (0.56, 0.81) | 0.43 (0.32, 0.53) | 0.25 (-0.07, 0.40) | 0.74 (0.63, 0.81) | 0.57 (0.37, 0.73) |
| sigmoid | **0.88** (0.62, 0.97) | 0.73 (0.52, 0.88) | 0.13 (0.07, 0.18) | 0.86 (0.75, 0.93) | -0.11 (-0.48, 0.24) | -1.36 (-2.08, -0.52) |
| step | **0.60** (0.36, 0.73) | 0.48 (0.34, 0.60) | 0.05 (0.02, 0.08) | 0.54 (0.44, 0.61) | 0.19 (-0.02, 0.39) | -0.40 (-1.00, 0.03) |
| 3dpoly | -1.28 (-4.59, 0.10) | -0.94 (-2.22, -0.08) | -0.87 (-1.46, -0.14) | -21.28 (-34.52, -14.29) | **0.34** (0.04, 0.61) | -0.26 (-1.39, 0.30) |
| sin | **0.87** (0.62, 0.98) | 0.66 (0.49, 0.81) | 0.05 (-0.00, 0.09) | 0.71 (0.52, 0.82) | 0.15 (-0.15, 0.47) | -0.91 (-1.77, -0.30) |
| linear | 0.96 (0.90, 0.99) | 0.86 (0.72, 0.94) | 0.29 (0.22, 0.36) | **0.99** (0.97, 1.00) | 0.24 (0.05, 0.43) | -0.23 (-0.68, 0.17) |
| rand_pw | **0.79** (-0.15, 0.98) | 0.63 (-0.04, 0.84) | 0.21 (-0.20, 0.40) | 0.46 (-1.08, 0.98) | 0.35 (-2.68, 0.92) | 0.35 (-3.77, 0.85) |

*Figure 86.* Instrument strength: 0.7, Number of instruments: 5, Number of samples: 1000, Number of experiments: 100, Training steps: 400, Number of critics: 50, Kernel radius: 50 data points, Critic jitter: 1, Test points: Distribution

### F.3.3. INSTRUMENT STRENGTH: 0.9

| function | AGMM Avg | DeepIV | 2SLSpoly | 2SLS | DirectPoly | DirectNN |
|---|---|---|---|---|---|---|
| abs | **0.88** (0.71, 0.96) | 0.75 (0.60, 0.87) | 0.56 (0.44, 0.66) | -0.65 (-0.99, -0.38) | 0.47 (0.31, 0.60) | -0.30 (-1.16, 0.22) |
| 2dpoly | 0.90 (0.83, 0.95) | 0.89 (0.76, 0.94) | 0.49 (0.37, 0.60) | -0.18 (-0.55, 0.06) | **0.98** (0.96, 0.99) | 0.91 (0.79, 0.95) |
| sigmoid | 0.85 (0.61, 0.95) | 0.81 (0.67, 0.91) | 0.12 (0.06, 0.17) | 0.84 (0.74, 0.91) | **0.93** (0.84, 0.97) | -0.36 (-1.14, 0.22) |
| step | 0.49 (0.32, 0.60) | 0.49 (0.35, 0.57) | 0.04 (0.01, 0.06) | 0.46 (0.38, 0.51) | **0.58** (0.55, 0.60) | 0.19 (-0.23, 0.49) |
| 3dpoly | 0.26 (-1.33, 0.65) | 0.38 (-0.81, 0.73) | -0.08 (-0.63, 0.44) | -21.29 (-34.54, -11.35) | **0.97** (0.93, 0.99) | 0.67 (-0.26, 0.93) |
| sin | 0.77 (0.54, 0.93) | 0.76 (0.55, 0.90) | 0.01 (-0.04, 0.05) | 0.58 (0.42, 0.70) | **0.96** (0.90, 0.98) | -0.18 (-0.78, 0.35) |
| linear | 0.97 (0.90, 0.99) | 0.92 (0.82, 0.98) | 0.29 (0.21, 0.35) | **1.00** (0.97, 1.00) | 0.94 (0.88, 0.98) | 0.42 (0.14, 0.65) |
| rand_pw | **0.85** (-0.15, 0.98) | 0.79 (0.33, 0.93) | 0.23 (-0.26, 0.55) | 0.32 (-1.29, 0.98) | 0.71 (-0.36, 0.98) | 0.73 (-1.49, 0.96) |

*Figure 87.* Instrument strength: 0.9, Number of instruments: 5, Number of samples: 1000, Number of experiments: 100, Training steps: 400, Number of critics: 50, Kernel radius: 50 data points, Critic jitter: 1, Test points: Distribution

## F.4. Number of Instruments:10

### F.4.1. INSTRUMENT STRENGTH: 0.5

| function | AGMM Avg | DeepIV | 2SLSpoly | 2SLS | DirectPoly | DirectNN |
|---|---|---|---|---|---|---|
| abs | -1.04 (-2.96, -0.10) | **-0.11 (-0.55, 0.12)** | -0.47 (-0.87, -0.26) | -0.54 (-0.92, -0.27) | -4.81 (-5.53, -3.95) | -4.63 (-5.62, -3.51) |
| 2dpoly | **0.68 (0.41, 0.85)** | 0.24 (0.07, 0.39) | 0.18 (-0.07, 0.37) | 0.38 (0.17, 0.54) | 0.45 (0.34, 0.51) | 0.40 (0.27, 0.52) |
| sigmoid | 0.18 (-0.70, 0.73) | 0.71 (0.41, 0.90) | 0.22 (0.11, 0.28) | **0.89 (0.70, 0.97)** | -1.70 (-2.16, -1.37) | -2.28 (-3.02, -1.50) |
| step | 0.26 (-0.26, 0.58) | 0.50 (0.30, 0.65) | 0.11 (0.04, 0.15) | **0.60 (0.42, 0.67)** | -0.83 (-1.19, -0.65) | -1.17 (-1.73, -0.75) |
| 3dpoly | -1.22 (-3.78, -0.03) | -1.41 (-2.61, -0.60) | -1.50 (-2.55, -0.76) | -15.81 (-23.45, -10.25) | **-0.36 (-0.84, 0.04)** | -0.56 (-1.36, -0.00) |
| sin | 0.51 (-0.25, 0.90) | 0.65 (0.41, 0.82) | 0.13 (0.04, 0.19) | **0.76 (0.53, 0.89)** | -1.18 (-1.56, -0.89) | -1.76 (-2.52, -1.10) |
| linear | 0.42 (-0.22, 0.75) | 0.76 (0.61, 0.91) | 0.39 (0.30, 0.46) | **0.99 (0.96, 1.00)** | -0.59 (-0.78, -0.41) | -0.69 (-1.04, -0.32) |
| rand_pw | 0.43 (-2.30, 0.90) | 0.29 (-0.49, 0.72) | 0.15 (-0.46, 0.42) | **0.56 (-0.99, 0.98)** | -0.03 (-6.27, 0.81) | 0.14 (-4.67, 0.83) |

*Figure 88.* Instrument strength: 0.5, Number of instruments: 10, Number of samples: 1000, Number of experiments: 100, Training steps: 400, Number of critics: 50, Kernel radius: 50 data points, Critic jitter: 1, Test points: Distribution

### F.4.2. INSTRUMENT STRENGTH: 0.7

| function | AGMM Avg | DeepIV | 2SLSpoly | 2SLS | DirectPoly | DirectNN |
|---|---|---|---|---|---|---|
| abs | 0.31 (-0.76, 0.70) | **0.38 (0.10, 0.54)** | 0.26 (-0.19, 0.47) | -0.63 (-1.03, -0.35) | -1.60 (-2.41, -1.15) | -2.60 (-4.06, -1.88) |
| 2dpoly | **0.79 (0.52, 0.91)** | 0.68 (0.52, 0.78) | 0.43 (0.30, 0.51) | 0.21 (-0.02, 0.41) | 0.75 (0.65, 0.81) | 0.57 (0.37, 0.68) |
| sigmoid | 0.81 (0.33, 0.95) | 0.79 (0.65, 0.91) | 0.14 (0.06, 0.19) | **0.85 (0.73, 0.94)** | -0.06 (-0.46, 0.25) | -1.19 (-2.10, -0.64) |
| step | **0.64 (0.43, 0.74)** | 0.54 (0.40, 0.63) | 0.05 (0.02, 0.09) | 0.53 (0.41, 0.61) | 0.23 (0.03, 0.38) | -0.29 (-0.89, 0.03) |
| 3dpoly | -0.44 (-2.35, 0.26) | -1.16 (-2.61, -0.22) | -0.89 (-1.50, -0.26) | -21.53 (-32.38, -13.40) | **0.38 (0.05, 0.60)** | -0.21 (-1.39, 0.32) |
| sin | **0.91 (0.63, 0.98)** | 0.74 (0.61, 0.87) | 0.05 (0.00, 0.10) | 0.69 (0.52, 0.82) | 0.20 (-0.18, 0.45) | -0.87 (-1.59, -0.14) |
| linear | 0.84 (0.47, 0.97) | 0.90 (0.77, 0.96) | 0.29 (0.21, 0.35) | **0.99 (0.97, 1.00)** | 0.28 (0.02, 0.46) | -0.16 (-0.73, 0.20) |
| rand_pw | **0.68 (-1.05, 0.94)** | 0.58 (0.02, 0.88) | 0.19 (-0.26, 0.47) | 0.51 (-1.42, 0.98) | 0.36 (-3.20, 0.93) | 0.35 (-3.92, 0.88) |

*Figure 89.* Instrument strength: 0.7, Number of instruments: 10, Number of samples: 1000, Number of experiments: 100, Training steps: 400, Number of critics: 50, Kernel radius: 50 data points, Critic jitter: 1, Test points: Distribution

### F.4.3. INSTRUMENT STRENGTH: 0.9

| function | AGMM Avg | DeepIV | 2SLSpoly | 2SLS | DirectPoly | DirectNN |
|---|---|---|---|---|---|---|
| abs | **0.80 (0.55, 0.95)** | 0.73 (0.60, 0.86) | 0.58 (0.39, 0.67) | -0.65 (-0.98, -0.32) | 0.50 (0.35, 0.60) | -0.34 (-1.10, 0.17) |
| 2dpoly | 0.88 (0.75, 0.95) | 0.89 (0.69, 0.95) | 0.49 (0.34, 0.58) | -0.21 (-0.47, 0.05) | **0.98 (0.96, 0.99)** | 0.92 (0.81, 0.95) |
| sigmoid | 0.87 (0.60, 0.97) | 0.83 (0.69, 0.92) | 0.12 (0.05, 0.17) | 0.83 (0.72, 0.91) | **0.93 (0.86, 0.97)** | -0.25 (-0.99, 0.21) |
| step | 0.58 (0.32, 0.72) | 0.49 (0.38, 0.57) | 0.04 (0.01, 0.07) | 0.45 (0.37, 0.52) | **0.58 (0.55, 0.60)** | 0.16 (-0.35, 0.53) |
| 3dpoly | 0.35 (-0.86, 0.73) | 0.21 (-0.81, 0.56) | 0.02 (-0.68, 0.43) | -20.47 (-35.69, -10.74) | **0.97 (0.94, 0.99)** | 0.77 (-0.76, 0.94) |
| sin | 0.79 (0.46, 0.96) | 0.74 (0.62, 0.86) | 0.00 (-0.04, 0.04) | 0.56 (0.38, 0.71) | **0.96 (0.91, 0.98)** | -0.05 (-0.58, 0.42) |
| linear | 0.95 (0.83, 0.99) | 0.92 (0.82, 0.98) | 0.28 (0.21, 0.36) | **0.99 (0.98, 1.00)** | 0.94 (0.89, 0.98) | 0.44 (0.17, 0.70) |
| rand_pw | **0.84 (-0.39, 0.97)** | 0.76 (0.28, 0.94) | 0.23 (-0.19, 0.62) | 0.30 (-1.46, 0.98) | 0.69 (-0.21, 0.98) | 0.72 (-1.15, 0.95) |

*Figure 90.* Instrument strength: 0.9, Number of instruments: 10, Number of samples: 1000, Number of experiments: 100, Training steps: 400, Number of critics: 50, Kernel radius: 50 data points, Critic jitter: 1, Test points: Distribution

# G. Grid Test Points

## G.1. Number of Instruments:2

### G.1.1. INSTRUMENT STRENGTH: 0.5

| function | AGMM Avg | DeepIV | 2SLSpoly | 2SLS | DirectPoly | DirectNN |
|---|---|---|---|---|---|---|
| abs | **0.67** (0.17, 0.86) | 0.08 (-0.12, 0.22) | -0.16 (-0.41, -0.05) | -0.20 (-0.45, -0.06) | -5.43 (-6.51, -4.71) | -5.33 (-6.78, -4.27) |
| 2dpoly | **0.95** (0.82, 0.98) | 0.33 (0.17, 0.44) | 0.33 (0.09, 0.45) | 0.60 (0.46, 0.68) | 0.45 (0.33, 0.55) | 0.42 (0.25, 0.56) |
| sigmoid | 0.86 (0.43, 0.96) | 0.39 (0.21, 0.54) | 0.20 (0.09, 0.30) | **0.90** (0.71, 0.96) | -2.02 (-2.46, -1.53) | -2.38 (-3.27, -1.66) |
| step | 0.59 (0.15, 0.74) | 0.29 (0.17, 0.44) | 0.11 (0.03, 0.17) | **0.66** (0.47, 0.73) | -1.52 (-1.93, -1.12) | -1.74 (-2.31, -1.19) |
| 3dpoly | -1.91 (-4.67, -0.15) | -0.68 (-1.90, -0.04) | -1.11 (-2.02, -0.18) | -13.88 (-24.32, -8.76) | **-0.36 (-1.02, 0.19)** | -0.54 (-1.47, 0.13) |
| sin | **0.87** (0.42, 0.96) | 0.30 (0.16, 0.47) | 0.11 (-0.00, 0.18) | 0.76 (0.51, 0.89) | -1.41 (-1.81, -1.00) | -2.00 (-2.60, -1.41) |
| linear | 0.96 (0.79, 1.00) | 0.50 (0.32, 0.66) | 0.40 (0.31, 0.48) | **0.99 (0.96, 1.00)** | -0.57 (-0.75, -0.37) | -0.65 (-0.98, -0.32) |
| rand_pw | **0.73 (-0.53, 0.98)** | 0.28 (-0.27, 0.58) | 0.19 (-0.37, 0.44) | 0.67 (-0.77, 0.99) | -0.14 (-6.99, 0.84) | -0.00 (-6.85, 0.83) |

*Figure 91.* Instrument strength: 0.5, Number of instruments: 2, Number of samples: 1000, Number of experiments: 100, Training steps: 400, Number of critics: 50, Kernel radius: 50 data points, Critic jitter: 1, Test points: Grid

### G.1.2. INSTRUMENT STRENGTH: 0.7

| function | AGMM Avg | DeepIV | 2SLSpoly | 2SLS | DirectPoly | DirectNN |
|---|---|---|---|---|---|---|
| abs | **0.82 (0.51, 0.93)** | 0.59 (0.41, 0.74) | 0.43 (0.15, 0.57) | -0.15 (-0.33, -0.03) | -2.14 (-3.20, -1.43) | -2.90 (-4.43, -1.90) |
| 2dpoly | **0.95 (0.90, 0.98)** | 0.80 (0.66, 0.89) | 0.47 (0.36, 0.58) | 0.57 (0.44, 0.67) | 0.74 (0.65, 0.81) | 0.64 (0.45, 0.76) |
| sigmoid | **0.91 (0.66, 0.98)** | 0.77 (0.59, 0.88) | 0.17 (0.10, 0.23) | 0.88 (0.75, 0.95) | -0.29 (-0.62, 0.05) | -1.06 (-2.00, -0.46) |
| step | **0.71 (0.47, 0.80)** | 0.56 (0.40, 0.66) | 0.09 (0.04, 0.13) | 0.65 (0.52, 0.71) | -0.20 (-0.43, 0.09) | -0.69 (-1.16, -0.21) |
| 3dpoly | 0.36 (-1.24, 0.74) | -0.67 (-1.96, -0.02) | -0.72 (-1.82, -0.01) | -18.36 (-32.25, -10.79) | **0.38 (-0.11, 0.65)** | -0.08 (-1.47, 0.52) |
| sin | **0.91 (0.63, 0.98)** | 0.66 (0.44, 0.80) | 0.06 (0.01, 0.11) | 0.71 (0.55, 0.83) | 0.05 (-0.27, 0.37) | -0.69 (-1.53, -0.05) |
| linear | 0.97 (0.89, 0.99) | 0.87 (0.71, 0.94) | 0.33 (0.26, 0.41) | **0.99 (0.98, 1.00)** | 0.28 (0.07, 0.45) | 0.03 (-0.39, 0.34) |
| rand_pw | **0.89 (0.01, 0.99)** | 0.70 (0.01, 0.88) | 0.27 (-0.15, 0.50) | 0.61 (-1.01, 0.98) | 0.39 (-3.09, 0.92) | 0.43 (-3.39, 0.89) |

*Figure 92.* Instrument strength: 0.7, Number of instruments: 2, Number of samples: 1000, Number of experiments: 100, Training steps: 400, Number of critics: 50, Kernel radius: 50 data points, Critic jitter: 1, Test points: Grid

### G.1.3. INSTRUMENT STRENGTH: 0.9

| function | AGMM Avg | DeepIV | 2SLSpoly | 2SLS | DirectPoly | DirectNN |
|---|---|---|---|---|---|---|
| abs | **0.91 (0.74, 0.97)** | 0.83 (0.64, 0.93) | 0.62 (0.52, 0.70) | -0.04 (-0.16, 0.01) | 0.41 (0.16, 0.59) | 0.05 (-0.49, 0.43) |
| 2dpoly | 0.98 (0.94, 0.99) | 0.94 (0.84, 0.97) | 0.49 (0.35, 0.63) | 0.49 (0.35, 0.60) | **0.98 (0.96, 0.99)** | 0.94 (0.89, 0.96) |
| sigmoid | **0.93 (0.79, 0.99)** | 0.88 (0.76, 0.95) | 0.14 (0.09, 0.20) | 0.87 (0.76, 0.93) | 0.91 (0.84, 0.96) | 0.41 (0.10, 0.68) |
| step | **0.77 (0.62, 0.85)** | 0.70 (0.60, 0.78) | 0.08 (0.04, 0.12) | 0.66 (0.56, 0.71) | 0.71 (0.64, 0.75) | 0.44 (-0.01, 0.65) |
| 3dpoly | 0.68 (-0.29, 0.92) | 0.42 (-1.14, 0.76) | 0.03 (-0.91, 0.42) | -18.29 (-31.99, -7.88) | **0.97 (0.94, 0.99)** | 0.84 (-0.17, 0.95) |
| sin | 0.92 (0.77, 0.98) | 0.80 (0.60, 0.90) | 0.01 (-0.03, 0.05) | 0.60 (0.43, 0.72) | **0.97 (0.92, 0.98)** | 0.50 (0.19, 0.72) |
| linear | 0.98 (0.93, 1.00) | 0.95 (0.87, 0.99) | 0.32 (0.26, 0.40) | **1.00 (0.98, 1.00)** | 0.94 (0.89, 0.97) | 0.76 (0.64, 0.87) |
| rand_pw | **0.92 (0.56, 0.99)** | 0.84 (0.45, 0.97) | 0.28 (-0.20, 0.58) | 0.53 (-1.01, 0.98) | 0.76 (-0.23, 0.98) | 0.85 (-0.18, 0.97) |

*Figure 93.* Instrument strength: 0.9, Number of instruments: 2, Number of samples: 1000, Number of experiments: 100, Training steps: 400, Number of critics: 50, Kernel radius: 50 data points, Critic jitter: 1, Test points: Grid

## G.2. Number of Instruments:3

### G.2.1. INSTRUMENT STRENGTH: 0.5

| function | AGMM Avg | DeepIV | 2SLSpoly | 2SLS | DirectPoly | DirectNN |
|---|---|---|---|---|---|---|
| abs | **0.68** (0.33, 0.84) | 0.06 (-0.18, 0.20) | -0.16 (-0.39, -0.04) | -0.19 (-0.43, -0.05) | -5.50 (-6.18, -4.69) | -5.40 (-6.51, -4.48) |
| 2dpoly | **0.93** (0.83, 0.97) | 0.31 (0.14, 0.43) | 0.32 (0.08, 0.43) | 0.60 (0.48, 0.68) | 0.44 (0.35, 0.52) | 0.41 (0.27, 0.54) |
| sigmoid | 0.87 (0.69, 0.96) | 0.40 (0.24, 0.55) | 0.19 (0.12, 0.29) | **0.89** (0.77, 0.96) | -2.00 (-2.36, -1.68) | -2.28 (-2.98, -1.69) |
| step | 0.60 (0.43, 0.72) | 0.31 (0.12, 0.40) | 0.10 (0.04, 0.17) | **0.65** (0.51, 0.72) | -1.53 (-1.91, -1.23) | -1.72 (-2.30, -1.30) |
| 3dpoly | -3.25 (-6.39, -0.51) | -0.61 (-1.54, 0.01) | -0.96 (-1.99, -0.44) | -12.73 (-22.77, -8.33) | **-0.32** (-0.83, 0.12) | -0.52 (-1.57, 0.07) |
| sin | **0.86** (0.62, 0.97) | 0.31 (0.18, 0.47) | 0.09 (0.02, 0.18) | 0.76 (0.56, 0.87) | -1.41 (-1.73, -1.03) | -1.89 (-2.71, -1.07) |
| linear | 0.97 (0.91, 0.99) | 0.50 (0.33, 0.69) | 0.41 (0.32, 0.48) | **0.99** (0.97, 1.00) | -0.56 (-0.70, -0.41) | -0.56 (-0.96, -0.30) |
| rand_pw | 0.66 (-0.48, 0.96) | 0.26 (-0.14, 0.53) | 0.16 (-0.28, 0.43) | **0.68** (-0.67, 0.98) | -0.10 (-6.55, 0.83) | 0.04 (-5.42, 0.83) |

*Figure 94.* Instrument strength: 0.5, Number of instruments: 3, Number of samples: 1000, Number of experiments: 100, Training steps: 400, Number of critics: 50, Kernel radius: 50 data points, Critic jitter: 1, Test points: Grid

### G.2.2. INSTRUMENT STRENGTH: 0.7

| function | AGMM Avg | DeepIV | 2SLSpoly | 2SLS | DirectPoly | DirectNN |
|---|---|---|---|---|---|---|
| abs | **0.80** (0.55, 0.94) | 0.55 (0.34, 0.71) | 0.42 (0.15, 0.55) | -0.13 (-0.44, -0.03) | -2.20 (-2.94, -1.43) | -3.06 (-4.27, -1.86) |
| 2dpoly | **0.93** (0.84, 0.98) | 0.80 (0.71, 0.88) | 0.46 (0.36, 0.58) | 0.57 (0.44, 0.68) | 0.74 (0.65, 0.81) | 0.63 (0.46, 0.73) |
| sigmoid | **0.92** (0.78, 0.97) | 0.75 (0.55, 0.88) | 0.16 (0.10, 0.23) | 0.87 (0.78, 0.94) | -0.28 (-0.65, 0.07) | -1.09 (-1.74, -0.42) |
| step | **0.70** (0.56, 0.78) | 0.54 (0.40, 0.67) | 0.08 (0.04, 0.13) | 0.63 (0.55, 0.70) | -0.20 (-0.48, 0.08) | -0.63 (-1.24, -0.26) |
| 3dpoly | -0.26 (-2.67, 0.56) | -0.58 (-1.75, 0.19) | -0.76 (-1.59, -0.22) | -19.39 (-29.16, -11.07) | **0.38** (-0.06, 0.61) | -0.03 (-0.86, 0.42) |
| sin | **0.88** (0.71, 0.97) | 0.65 (0.48, 0.80) | 0.05 (0.01, 0.11) | 0.70 (0.57, 0.80) | 0.03 (-0.26, 0.37) | -0.70 (-1.47, -0.15) |
| linear | 0.98 (0.93, 0.99) | 0.86 (0.76, 0.94) | 0.34 (0.26, 0.41) | **0.99** (0.98, 1.00) | 0.28 (0.09, 0.44) | -0.03 (-0.30, 0.27) |
| rand_pw | **0.86** (0.02, 0.99) | 0.64 (0.15, 0.88) | 0.26 (-0.13, 0.50) | 0.59 (-0.79, 0.98) | 0.41 (-4.18, 0.92) | 0.40 (-4.12, 0.88) |

*Figure 95.* Instrument strength: 0.7, Number of instruments: 3, Number of samples: 1000, Number of experiments: 100, Training steps: 400, Number of critics: 50, Kernel radius: 50 data points, Critic jitter: 1, Test points: Grid

### G.2.3. INSTRUMENT STRENGTH: 0.9

| function | AGMM Avg | DeepIV | 2SLSpoly | 2SLS | DirectPoly | DirectNN |
|---|---|---|---|---|---|---|
| abs | **0.91** (0.73, 0.98) | 0.82 (0.65, 0.92) | 0.61 (0.52, 0.69) | -0.04 (-0.18, 0.02) | 0.44 (0.12, 0.59) | 0.03 (-0.42, 0.52) |
| 2dpoly | 0.95 (0.88, 0.98) | 0.93 (0.81, 0.97) | 0.48 (0.32, 0.60) | 0.49 (0.35, 0.61) | **0.98** (0.96, 0.99) | 0.94 (0.89, 0.96) |
| sigmoid | **0.93** (0.83, 0.97) | 0.86 (0.69, 0.93) | 0.14 (0.08, 0.21) | 0.86 (0.79, 0.93) | 0.92 (0.83, 0.96) | 0.43 (0.13, 0.65) |
| step | **0.75** (0.59, 0.82) | 0.68 (0.55, 0.75) | 0.08 (0.04, 0.13) | 0.64 (0.56, 0.71) | 0.71 (0.63, 0.75) | 0.42 (0.03, 0.63) |
| 3dpoly | 0.66 (0.12, 0.88) | 0.55 (-0.69, 0.82) | 0.08 (-0.66, 0.40) | -17.18 (-30.54, -9.27) | **0.98** (0.95, 0.99) | 0.82 (-0.15, 0.96) |
| sin | 0.87 (0.70, 0.97) | 0.77 (0.60, 0.89) | 0.01 (-0.02, 0.07) | 0.58 (0.45, 0.71) | **0.97** (0.90, 0.98) | 0.47 (0.19, 0.66) |
| linear | 0.98 (0.95, 1.00) | 0.94 (0.86, 0.98) | 0.33 (0.24, 0.39) | **1.00** (0.99, 1.00) | 0.94 (0.88, 0.98) | 0.76 (0.59, 0.86) |
| rand_pw | **0.91** (0.40, 0.99) | 0.81 (0.31, 0.95) | 0.26 (-0.29, 0.57) | 0.57 (-1.10, 0.98) | 0.76 (-0.40, 0.98) | 0.85 (0.07, 0.97) |

*Figure 96.* Instrument strength: 0.9, Number of instruments: 3, Number of samples: 1000, Number of experiments: 100, Training steps: 400, Number of critics: 50, Kernel radius: 50 data points, Critic jitter: 1, Test points: Grid

## G.3. Number of Instruments:5

### G.3.1. INSTRUMENT STRENGTH: 0.5

| function | AGMM Avg | DeepIV | 2SLSpoly | 2SLS | DirectPoly | DirectNN |
|---|---|---|---|---|---|---|
| abs | **0.49** (-0.27, 0.85) | 0.07 (-0.15, 0.23) | -0.15 (-0.37, -0.04) | -0.18 (-0.43, -0.07) | -5.54 (-6.45, -4.65) | -5.49 (-6.59, -4.31) |
| 2dpoly | **0.90** (0.78, 0.97) | 0.33 (0.18, 0.45) | 0.30 (0.12, 0.46) | 0.60 (0.47, 0.68) | 0.44 (0.32, 0.53) | 0.43 (0.25, 0.56) |
| sigmoid | 0.87 (0.66, 0.94) | 0.48 (0.26, 0.67) | 0.23 (0.12, 0.34) | **0.91 (0.75, 0.97)** | -2.09 (-2.50, -1.64) | -2.38 (-3.14, -1.69) |
| step | 0.59 (0.39, 0.72) | 0.34 (0.19, 0.48) | 0.13 (0.05, 0.20) | **0.67 (0.51, 0.73)** | -1.57 (-2.05, -1.23) | -1.73 (-2.36, -1.30) |
| 3dpoly | -3.18 (-7.22, -0.32) | -0.64 (-1.75, -0.10) | -1.09 (-2.28, -0.26) | -13.77 (-23.23, -6.88) | **-0.32 (-0.95, 0.18)** | -0.45 (-1.44, 0.12) |
| sin | **0.94 (0.69, 0.98)** | 0.38 (0.19, 0.56) | 0.12 (0.02, 0.22) | 0.78 (0.56, 0.89) | -1.46 (-1.84, -1.08) | -1.93 (-2.48, -1.30) |
| linear | 0.94 (0.78, 0.99) | 0.56 (0.43, 0.73) | 0.41 (0.34, 0.52) | **0.99 (0.96, 1.00)** | -0.57 (-0.77, -0.40) | -0.65 (-1.06, -0.33) |
| rand_pw | **0.63 (-1.03, 0.96)** | 0.27 (-0.19, 0.61) | 0.20 (-0.23, 0.45) | 0.62 (-0.72, 0.98) | -0.08 (-7.10, 0.83) | 0.05 (-5.66, 0.82) |

*Figure 97.* Instrument strength: 0.5, Number of instruments: 5, Number of samples: 1000, Number of experiments: 100, Training steps: 400, Number of critics: 50, Kernel radius: 50 data points, Critic jitter: 1, Test points: Grid

### G.3.2. INSTRUMENT STRENGTH: 0.7

| function | AGMM Avg | DeepIV | 2SLSpoly | 2SLS | DirectPoly | DirectNN |
|---|---|---|---|---|---|---|
| abs | **0.76 (0.40, 0.92)** | 0.52 (0.35, 0.66) | 0.41 (0.17, 0.56) | -0.14 (-0.34, -0.03) | -2.32 (-3.03, -1.48) | -2.98 (-4.23, -1.72) |
| 2dpoly | **0.92 (0.82, 0.97)** | 0.78 (0.64, 0.87) | 0.47 (0.28, 0.58) | 0.58 (0.40, 0.68) | 0.74 (0.62, 0.81) | 0.63 (0.45, 0.78) |
| sigmoid | **0.91 (0.69, 0.98)** | 0.76 (0.55, 0.90) | 0.17 (0.10, 0.24) | 0.89 (0.77, 0.95) | -0.30 (-0.69, 0.12) | -1.15 (-1.80, -0.43) |
| step | **0.68 (0.48, 0.79)** | 0.58 (0.43, 0.70) | 0.09 (0.05, 0.14) | 0.64 (0.52, 0.71) | -0.18 (-0.54, 0.14) | -0.66 (-1.41, -0.13) |
| 3dpoly | -1.11 (-4.76, 0.28) | -0.68 (-1.86, 0.23) | -0.96 (-1.88, 0.02) | -19.70 (-36.03, -10.71) | **0.32 (-0.11, 0.66)** | -0.15 (-1.24, 0.51) |
| sin | **0.90 (0.64, 0.98)** | 0.68 (0.51, 0.80) | 0.07 (0.01, 0.12) | 0.73 (0.54, 0.84) | 0.03 (-0.32, 0.39) | -0.69 (-1.51, -0.17) |
| linear | 0.97 (0.92, 0.99) | 0.88 (0.73, 0.95) | 0.33 (0.26, 0.42) | **0.99 (0.97, 1.00)** | 0.25 (0.07, 0.44) | -0.03 (-0.37, 0.30) |
| rand_pw | **0.84 (-0.08, 0.98)** | 0.66 (-0.02, 0.86) | 0.27 (-0.17, 0.48) | 0.57 (-0.84, 0.98) | 0.38 (-2.49, 0.92) | 0.36 (-2.94, 0.87) |

*Figure 98.* Instrument strength: 0.7, Number of instruments: 5, Number of samples: 1000, Number of experiments: 100, Training steps: 400, Number of critics: 50, Kernel radius: 50 data points, Critic jitter: 1, Test points: Grid

### G.3.3. INSTRUMENT STRENGTH: 0.9

| function | AGMM Avg | DeepIV | 2SLSpoly | 2SLS | DirectPoly | DirectNN |
|---|---|---|---|---|---|---|
| abs | **0.89 (0.71, 0.97)** | 0.80 (0.64, 0.90) | 0.61 (0.51, 0.72) | -0.05 (-0.15, 0.02) | 0.40 (0.07, 0.59) | 0.04 (-0.69, 0.44) |
| 2dpoly | 0.92 (0.86, 0.96) | 0.92 (0.77, 0.97) | 0.49 (0.36, 0.62) | 0.50 (0.30, 0.61) | **0.98 (0.96, 0.99)** | 0.94 (0.88, 0.97) |
| sigmoid | 0.89 (0.68, 0.97) | 0.84 (0.73, 0.93) | 0.15 (0.09, 0.21) | 0.87 (0.78, 0.93) | **0.92 (0.82, 0.96)** | 0.42 (0.01, 0.67) |
| step | 0.68 (0.51, 0.77) | 0.68 (0.52, 0.77) | 0.08 (0.05, 0.12) | 0.66 (0.55, 0.71) | **0.72 (0.62, 0.75)** | 0.44 (0.05, 0.69) |
| 3dpoly | 0.33 (-1.22, 0.72) | 0.47 (-0.86, 0.78) | -0.04 (-0.73, 0.49) | -18.54 (-32.32, -8.94) | **0.98 (0.94, 0.99)** | 0.78 (0.16, 0.96) |
| sin | 0.80 (0.57, 0.95) | 0.78 (0.58, 0.91) | 0.02 (-0.02, 0.07) | 0.61 (0.44, 0.74) | **0.97 (0.91, 0.98)** | 0.47 (0.13, 0.71) |
| linear | 0.98 (0.94, 1.00) | 0.94 (0.85, 0.98) | 0.32 (0.25, 0.40) | **1.00 (0.98, 1.00)** | 0.94 (0.88, 0.98) | 0.74 (0.61, 0.86) |
| rand_pw | **0.86 (-0.17, 0.99)** | 0.82 (0.29, 0.95) | 0.28 (-0.29, 0.56) | 0.51 (-1.20, 0.98) | 0.73 (-0.47, 0.98) | 0.82 (-0.42, 0.98) |

*Figure 99.* Instrument strength: 0.9, Number of instruments: 5, Number of samples: 1000, Number of experiments: 100, Training steps: 400, Number of critics: 50, Kernel radius: 50 data points, Critic jitter: 1, Test points: Grid

## G.4. Number of Instruments:10

### G.4.1. INSTRUMENT STRENGTH: 0.5

| function | AGMM Avg | DeepIV | 2SLSpoly | 2SLS | DirectPoly | DirectNN |
|---|---|---|---|---|---|---|
| abs | -1.34 (-3.28, -0.26) | **0.08 (-0.16, 0.23)** | -0.16 (-0.39, -0.04) | -0.21 (-0.44, -0.06) | -5.64 (-6.31, -4.72) | -5.37 (-6.49, -4.18) |
| 2dpoly | **0.69 (0.44, 0.86)** | 0.36 (0.18, 0.48) | 0.30 (0.09, 0.44) | 0.59 (0.45, 0.69) | 0.44 (0.33, 0.52) | 0.40 (0.25, 0.54) |
| sigmoid | 0.10 (-0.83, 0.69) | 0.74 (0.43, 0.91) | 0.26 (0.14, 0.34) | **0.91 (0.72, 0.97)** | -1.99 (-2.50, -1.67) | -2.48 (-3.20, -1.70) |
| step | 0.09 (-0.62, 0.54) | 0.56 (0.34, 0.70) | 0.15 (0.07, 0.21) | **0.66 (0.49, 0.74)** | -1.50 (-1.99, -1.25) | -1.76 (-2.41, -1.19) |
| 3dpoly | -0.99 (-3.83, 0.15) | -0.85 (-1.91, -0.10) | -1.23 (-2.28, -0.36) | -13.94 (-21.87, -8.72) | **-0.37 (-1.00, 0.15)** | -0.52 (-1.40, 0.08) |
| sin | 0.48 (-0.31, 0.88) | 0.66 (0.43, 0.83) | 0.16 (0.04, 0.23) | **0.78 (0.54, 0.90)** | -1.38 (-1.84, -1.10) | -1.87 (-2.63, -1.27) |
| linear | 0.43 (-0.18, 0.75) | 0.77 (0.62, 0.91) | 0.44 (0.35, 0.52) | **0.99 (0.96, 1.00)** | -0.58 (-0.77, -0.40) | -0.65 (-1.00, -0.28) |
| rand_pw | 0.46 (-2.56, 0.90) | 0.36 (-0.39, 0.75) | 0.21 (-0.33, 0.48) | **0.64 (-0.86, 0.98)** | -0.05 (-5.53, 0.81) | 0.11 (-6.03, 0.83) |

*Figure 100.* Instrument strength: 0.5, Number of instruments: 10, Number of samples: 1000, Number of experiments: 100, Training steps: 400, Number of critics: 50, Kernel radius: 50 data points, Critic jitter: 1, Test points: Grid

### G.4.2. INSTRUMENT STRENGTH: 0.7

| function | AGMM Avg | DeepIV | 2SLSpoly | 2SLS | DirectPoly | DirectNN |
|---|---|---|---|---|---|---|
| abs | 0.26 (-0.87, 0.74) | **0.53 (0.35, 0.65)** | 0.43 (0.17, 0.56) | -0.16 (-0.36, -0.04) | -2.10 (-3.07, -1.55) | -2.88 (-4.26, -1.96) |
| 2dpoly | **0.83 (0.58, 0.92)** | 0.75 (0.59, 0.84) | 0.46 (0.32, 0.56) | 0.56 (0.42, 0.69) | 0.74 (0.64, 0.82) | 0.63 (0.44, 0.74) |
| sigmoid | 0.82 (0.34, 0.96) | 0.82 (0.67, 0.94) | 0.17 (0.10, 0.24) | **0.88 (0.76, 0.95)** | -0.23 (-0.67, 0.10) | -0.96 (-1.93, -0.46) |
| step | **0.65 (0.34, 0.79)** | 0.64 (0.49, 0.72) | 0.09 (0.04, 0.14) | 0.64 (0.52, 0.71) | -0.13 (-0.46, 0.11) | -0.55 (-1.24, -0.17) |
| 3dpoly | -0.36 (-2.70, 0.44) | -0.83 (-2.43, 0.07) | -0.94 (-1.93, -0.11) | -19.94 (-33.56, -10.56) | **0.35 (-0.09, 0.61)** | -0.10 (-1.21, 0.41) |
| sin | **0.91 (0.66, 0.98)** | 0.75 (0.62, 0.87) | 0.07 (0.02, 0.13) | 0.71 (0.54, 0.84) | 0.10 (-0.31, 0.37) | -0.70 (-1.36, -0.01) |
| linear | 0.85 (0.56, 0.98) | 0.91 (0.80, 0.97) | 0.34 (0.25, 0.41) | **0.99 (0.98, 1.00)** | 0.29 (0.04, 0.46) | 0.03 (-0.47, 0.34) |
| rand_pw | **0.72 (-1.23, 0.95)** | 0.65 (-0.01, 0.89) | 0.23 (-0.19, 0.53) | 0.64 (-1.05, 0.98) | 0.39 (-2.39, 0.93) | 0.40 (-3.76, 0.90) |

*Figure 101.* Instrument strength: 0.7, Number of instruments: 10, Number of samples: 1000, Number of experiments: 100, Training steps: 400, Number of critics: 50, Kernel radius: 50 data points, Critic jitter: 1, Test points: Grid

### G.4.3. INSTRUMENT STRENGTH: 0.9

| function | AGMM Avg | DeepIV | 2SLSpoly | 2SLS | DirectPoly | DirectNN |
|---|---|---|---|---|---|---|
| abs | **0.82 (0.43, 0.95)** | 0.81 (0.58, 0.91) | 0.61 (0.50, 0.72) | -0.04 (-0.16, 0.01) | 0.43 (0.18, 0.61) | 0.07 (-0.52, 0.50) |
| 2dpoly | 0.90 (0.80, 0.96) | 0.92 (0.71, 0.96) | 0.49 (0.32, 0.58) | 0.49 (0.34, 0.60) | **0.98 (0.96, 0.99)** | 0.94 (0.88, 0.97) |
| sigmoid | 0.90 (0.70, 0.98) | 0.86 (0.72, 0.94) | 0.15 (0.09, 0.21) | 0.86 (0.75, 0.93) | **0.92 (0.83, 0.96)** | 0.44 (0.05, 0.70) |
| step | 0.72 (0.51, 0.86) | 0.69 (0.55, 0.77) | 0.08 (0.04, 0.13) | 0.65 (0.55, 0.72) | **0.72 (0.63, 0.76)** | 0.43 (-0.01, 0.70) |
| 3dpoly | 0.47 (-0.71, 0.80) | 0.31 (-0.79, 0.64) | 0.06 (-0.79, 0.48) | -17.28 (-34.57, -8.63) | **0.98 (0.94, 0.99)** | 0.85 (-0.07, 0.96) |
| sin | 0.82 (0.55, 0.97) | 0.76 (0.64, 0.88) | 0.01 (-0.02, 0.05) | 0.59 (0.40, 0.73) | **0.96 (0.92, 0.98)** | 0.51 (0.30, 0.76) |
| linear | 0.97 (0.90, 0.99) | 0.94 (0.85, 0.98) | 0.32 (0.24, 0.40) | **1.00 (0.98, 1.00)** | 0.94 (0.89, 0.98) | 0.77 (0.61, 0.89) |
| rand_pw | **0.86 (-0.22, 0.98)** | 0.80 (0.28, 0.95) | 0.26 (-0.29, 0.57) | 0.58 (-1.10, 0.98) | 0.75 (-0.31, 0.98) | 0.83 (-0.22, 0.97) |

*Figure 102.* Instrument strength: 0.9, Number of instruments: 10, Number of samples: 1000, Number of experiments: 100, Training steps: 400, Number of critics: 50, Kernel radius: 50 data points, Critic jitter: 1, Test points: Grid