

Nonparametric Instrumental Variable Regression through Stochastic Gradients and Kernel Methods

Student: Caio Lins
Advisor: Yuri Saporito

EMAp – FGV

December 15, 2023



Summary

Instrumental Variable Regression

Nonparametric Instrumental Variable Regression

Stochastic Approximate Gradient Descent IV

Instrumental Variables

- Regression: $Y = h^*(X) + \varepsilon$, with $\mathbb{E}[\varepsilon] = 0$. Find h^* .

Instrumental Variables

- Regression: $Y = h^*(X) + \varepsilon$, with $\mathbb{E}[\varepsilon] = 0$. Find h^* .
- Minimizing $\text{MSE}(h) = \mathbb{E}[(Y - h(X))^2]$ gives

$$\hat{h}(X) = \mathbb{E}[Y \mid X] = h^*(X) + \mathbb{E}[\varepsilon \mid X].$$

Instrumental Variables

- Regression: $Y = h^*(X) + \varepsilon$, with $\mathbb{E}[\varepsilon] = 0$. Find h^* .
- Minimizing $\text{MSE}(h) = \mathbb{E}[(Y - h(X))^2]$ gives

$$\hat{h}(X) = \mathbb{E}[Y \mid X] = h^*(X) + \mathbb{E}[\varepsilon \mid X].$$

- What if $\mathbb{E}[\varepsilon \mid X] \neq 0$? That is, if X is *endogenous*?

Instrumental Variables

- Suppose we have access to a variable Z such that

Instrumental Variables

- Suppose we have access to a variable Z such that
 1. Z influences X , that is, $Z \not\perp\!\!\!\perp X$;

Instrumental Variables

- Suppose we have access to a variable Z such that
 1. Z influences X , that is, $Z \not\perp\!\!\!\perp X$;
 2. Z is *exogenous*, that is, $\mathbb{E}[\varepsilon \mid Z] = 0$.

Instrumental Variables

- Suppose we have access to a variable Z such that
 1. Z influences X , that is, $Z \not\perp\!\!\!\perp X$;
 2. Z is *exogenous*, that is, $\mathbb{E}[\varepsilon \mid Z] = 0$.

Z is called an *instrumental variable*.

Instrumental Variables

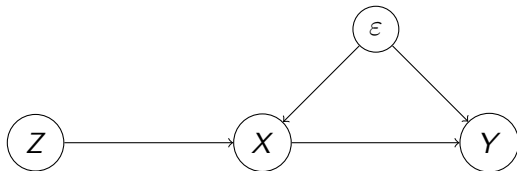
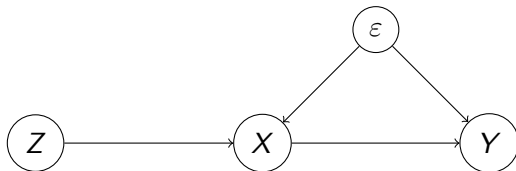


Figure: Causal diagram representing an instrumental variable for an endogenous covariate.

Instrumental Variables

Example:

- X = is smoker?
- Y = general health.
- Z = tax rate on tobacco products.
- ε = depression, self care.



Linear Model

$$\mathbf{Y} = \mathbf{X}\beta^* + \varepsilon.$$

Linear Model

$$\mathbf{Y} = \mathbf{X}\beta^* + \varepsilon.$$

- $\hat{\beta}_{\text{OLS}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}$ is biased and inconsistent;

Linear Model

$$\mathbf{Y} = \mathbf{X}\beta^* + \varepsilon.$$

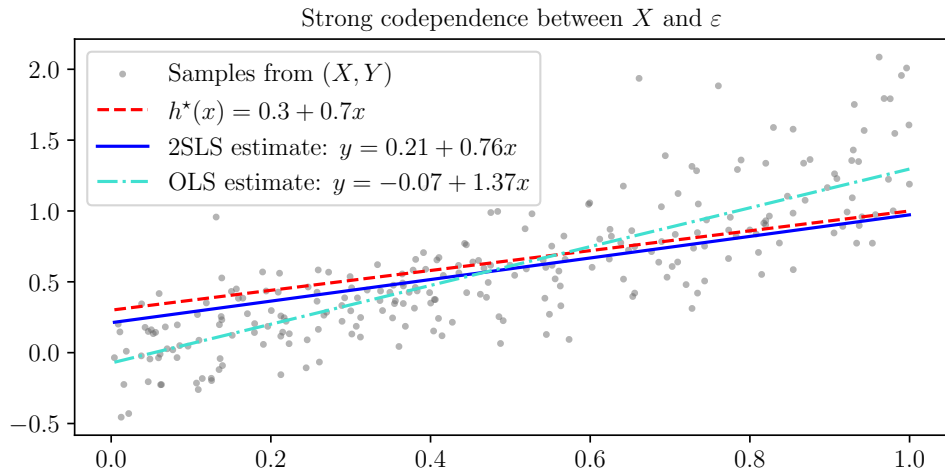
- $\hat{\beta}_{\text{OLS}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}$ is biased and inconsistent;
- $\hat{\beta}_{\text{IV}} = (\hat{\mathbf{X}}^T \hat{\mathbf{X}})^{-1} \hat{\mathbf{X}}^T \mathbf{Y}$, where

$$\hat{\mathbf{X}} = \mathbf{Z}(\mathbf{Z}^T \mathbf{Z})^{-1} \mathbf{Z}^T \mathbf{X},$$

is unbiased and consistent.

- Two Stages Least Squares (2SLS).

Linear Model



Summary

Instrumental Variable Regression

Nonparametric Instrumental Variable Regression

Stochastic Approximate Gradient Descent IV

Problem formulation

$$Y = h^*(X) + \varepsilon, \quad \text{with } \mathbb{E}[\varepsilon \mid X] \neq 0, Z \not\perp\!\!\!\perp X \text{ and } \mathbb{E}[\varepsilon \mid Z] = 0.$$

- We have access to $\{(X_i, Y_i, Z_i)\}$;

Problem formulation

$$Y = h^*(X) + \varepsilon, \quad \text{with } \mathbb{E}[\varepsilon \mid X] \neq 0, Z \not\perp\!\!\!\perp X \text{ and } \mathbb{E}[\varepsilon \mid Z] = 0.$$

- We have access to $\{(X_i, Y_i, Z_i)\}$;
- We assume $h^* \in L^2(X) = \{h : \mathbb{E}[h(X)^2] < \infty\}$.

Problem formulation

$$Y = h^*(X) + \varepsilon, \quad \text{with } \mathbb{E}[\varepsilon \mid X] \neq 0, Z \not\perp\!\!\!\perp X \text{ and } \mathbb{E}[\varepsilon \mid Z] = 0.$$

- We have access to $\{(X_i, Y_i, Z_i)\}$;
- We assume $h^* \in L^2(X) = \{h : \mathbb{E}[h(X)^2] < \infty\}$.
- Conditioning in Z :

$$\mathbb{E}[Y \mid Z] = \mathbb{E}[h^*(X) \mid Z] \iff r = \mathcal{P}[h^*],$$

where $r(Z) = \mathbb{E}[Y \mid Z]$ and $\mathcal{P} : L^2(X) \rightarrow L^2(Z)$ is the *conditional expectation operator*.

$$\mathcal{P}[h](z) = \mathbb{E}[h(X) \mid Z = z].$$

Problem formulation

$$Y = h^*(X) + \varepsilon, \quad \text{with } \mathbb{E}[\varepsilon \mid X] \neq 0, Z \not\perp\!\!\!\perp X \text{ and } \mathbb{E}[\varepsilon \mid Z] = 0.$$

- We have access to $\{(X_i, Y_i, Z_i)\}$;
- We assume $h^* \in L^2(X) = \{h : \mathbb{E}[h(X)^2] < \infty\}$.
- Conditioning in Z :

$$\mathbb{E}[Y \mid Z] = \mathbb{E}[h^*(X) \mid Z] \iff r = \mathcal{P}[h^*],$$

where $r(Z) = \mathbb{E}[Y \mid Z]$ and $\mathcal{P} : L^2(X) \rightarrow L^2(Z)$ is the *conditional expectation operator*.

$$\mathcal{P}[h](z) = \mathbb{E}[h(X) \mid Z = z].$$

- Ill posed problem.

Classical Approaches

Problem: $r = \mathcal{P}[h^*]$, where $r(Z) = \mathbb{E}[Y \mid Z]$ and $\mathcal{P}[h](Z) = \mathbb{E}[h(X) \mid Z]$.

Classical Approaches

Problem: $r = \mathcal{P}[h^*]$, where $r(Z) = \mathbb{E}[Y \mid Z]$ and $\mathcal{P}[h](Z) = \mathbb{E}[h(X) \mid Z]$.

- Nonlinear Two Stages Least Squares [3]:

- $h^*(x) \approx \sum_{j=1}^J \gamma_j p_j(x);$
- $\mathbb{E}[p_j(X) \mid Z = z] \approx \sum_{i=1}^n a_{ji} q_i(z).$

Classical Approaches

Problem: $r = \mathcal{P}[h^*]$, where $r(Z) = \mathbb{E}[Y \mid Z]$ and $\mathcal{P}[h](Z) = \mathbb{E}[h(X) \mid Z]$.

- Nonlinear Two Stages Least Squares [3]:

- $h^*(x) \approx \sum_{j=1}^J \gamma_j p_j(x);$
- $\mathbb{E}[p_j(X) \mid Z = z] \approx \sum_{i=1}^n a_{ji} q_i(z).$

- Iterated Tikhonov regularization [1]:

- $\arg \min_h \|\mathcal{P}[h] - r\|_{L^2(Z)}^2 + \alpha \|h\|_{L^2(X)}^2 = (\mathcal{P}^* \mathcal{P} + \alpha I)^{-1} \mathcal{P}^*[r];$
- $h_{k+1} = (\mathcal{P}^* \mathcal{P} + \alpha I)^{-1} [\mathcal{P}^* r + h_k]$

Summary

Instrumental Variable Regression

Nonparametric Instrumental Variable Regression

Stochastic Approximate Gradient Descent IV

Risk measure

We know that h^* satisfies

$$r = \mathcal{P}[h^*].$$

Risk measure

We know that h^* satisfies

$$r = \mathcal{P}[h^*].$$

- Define the risk

$$\mathcal{R}(h) = \mathbb{E}[\ell(r(Z), \mathcal{P}[h](Z))].$$

E.g., if $\ell(y, y') = (y - y')^2$:

$$\begin{aligned}\mathcal{R}(h) &= \mathbb{E} \left[(r(Z) - \mathcal{P}[h](Z))^2 \right] \\ &= \mathbb{E} \left[(\mathcal{P}[h - h^*](Z))^2 \right].\end{aligned}$$

Gradient Descent

$$\mathcal{R}(h) = \mathbb{E}[\ell(r(Z), \mathcal{P}[h](Z))].$$

Gradient Descent

$$\mathcal{R}(h) = \mathbb{E}[\ell(r(Z), \mathcal{P}[h](Z))].$$

- We showed that

$$\begin{aligned}\nabla \mathcal{R}(h)(x) &= \mathcal{P}^*[\partial_2(r(Z), \mathcal{P}[h](Z))] \\ &= \mathbb{E}[\Phi(x, Z) \partial_2(r(Z), \mathcal{P}[h](Z))],\end{aligned}$$

where $\Phi(x, z) = \frac{p_{XZ}(x, z)}{p_X(x)p_Z(z)}$.

Gradient Descent

$$\mathcal{R}(h) = \mathbb{E}[\ell(r(Z), \mathcal{P}[h](Z))].$$

- We showed that

$$\begin{aligned}\nabla \mathcal{R}(h)(x) &= \mathcal{P}^*[\partial_2(r(Z), \mathcal{P}[h](Z))] \\ &= \mathbb{E}[\Phi(x, Z) \partial_2(r(Z), \mathcal{P}[h](Z))],\end{aligned}$$

where $\Phi(x, z) = \frac{p_{XZ}(x, z)}{p_X(x)p_Z(z)}$.

- The term in blue is a stochastic gradient.

Gradient Descent

$$\mathcal{R}(h) = \mathbb{E}[\ell(r(Z), \mathcal{P}[h](Z))].$$

- We showed that

$$\begin{aligned}\nabla \mathcal{R}(h)(x) &= \mathcal{P}^*[\partial_2(r(Z), \mathcal{P}[h](Z))] \\ &= \mathbb{E}[\Phi(x, Z) \partial_2(r(Z), \mathcal{P}[h](Z))],\end{aligned}$$

where $\Phi(x, z) = \frac{p_{XZ}(x, z)}{p_X(x)p_Z(z)}$.

- The term in blue is a stochastic gradient.
- Problem: do not know Φ, \mathcal{P} nor r . Only have access to $\{(X_i, Y_i, Z_i)\}$.

Projected Gradient Descent

- Our regularization: look for solutions in $\mathcal{H} \subseteq L^2(X)$;

Projected Gradient Descent

- Our regularization: look for solutions in $\mathcal{H} \subseteq L^2(X)$;
- \mathcal{H} is convex, closed and bounded.

Projected Gradient Descent

- Our regularization: look for solutions in $\mathcal{H} \subseteq L^2(X)$;
- \mathcal{H} is convex, closed and bounded.
- E.g. for $A > 0$:

$$\mathcal{H} = \{h \in L^2(X) : |h(x)| \leq A \ \forall x\}.$$

Stochastic Approximate Gradient Descent

- $\nabla \mathcal{R}(h)(x) = \mathbb{E}[\Phi(x, Z) \partial_2(r(X), \mathcal{P}[h](Z))]$, but we do not know Φ, r, \mathcal{P} ;

Stochastic Approximate Gradient Descent

- $\nabla \mathcal{R}(h)(x) = \mathbb{E}[\Phi(x, Z) \partial_2(r(X), \mathcal{P}[h](Z))]$, but we do not know Φ, r, \mathcal{P} ;
- Assume we have estimators $\hat{\Phi}, \hat{r}$ and $\hat{\mathcal{P}}$, so that

$$u_h(x) = \hat{\Phi}(x, Z) \partial_2(\hat{r}(Z), \hat{\mathcal{P}}[h](Z))$$

is an *approximate stochastic gradient* for \mathcal{R} at h .

Kernel Methods

- Kernel Ridge Regression to compute $\hat{\Phi}$, \hat{r} and $\hat{\mathcal{P}}$;

Kernel Methods

- Kernel Ridge Regression to compute $\hat{\Phi}$, \hat{r} and $\hat{\mathcal{P}}$;
- Reproducing Kernel Hilbert Space (RKHS) as a class of approximating functions
 \implies closed form solutions;

Kernel Methods

- Kernel Ridge Regression to compute $\hat{\Phi}$, \hat{r} and $\hat{\mathcal{P}}$;
- Reproducing Kernel Hilbert Space (RKHS) as a class of approximating functions
 \implies closed form solutions;
- $\hat{\mathcal{P}}$ is tricky. We used KIV's first stage [4].

Algorithm 1: SAGD-IV

input : Samples $\{(\mathbf{z}_m)_{m=1}^M\}$. Estimators $\hat{\Phi}, \hat{r}$ and $\hat{\mathcal{P}}$. Sequence of learning rates $(\alpha_m)_{m=1}^M$.

output: \hat{h}

for $1 \leq m \leq M$ **do**

 Set $u_m = \hat{\Phi}(\cdot, \mathbf{z}_m) \partial_2 \ell \left(\hat{r}(\mathbf{z}_m), \hat{\mathcal{P}}[\hat{h}_{m-1}](\mathbf{z}_m) \right)$;

 Set $\hat{h}_m = \text{proj}_{\mathcal{H}} \left[\hat{h}_{m-1} - \alpha_m u_m \right]$;

end

Set $\hat{h} = \frac{1}{M} \sum_{m=1}^M \hat{h}_m$;

Theory

Theorem (SAGD–IV convergence rate)

Under suitable assumptions on $\ell, \mathcal{H}, \mathcal{P}$ and the estimators $\hat{\Phi}, \hat{r}, \hat{\mathcal{P}}$, we have

$$\mathbb{E}_{\mathbf{z}_{1:M}} \left[\mathcal{R}(\hat{h}) - \mathcal{R}(h^*) \right] \leq \frac{D^2}{2M\alpha_M} + \frac{\xi}{M} \sum_{m=1}^M \alpha_m + \tau \sqrt{\zeta},$$

Theory

Theorem (SAGD–IV convergence rate)

Under suitable assumptions on $\ell, \mathcal{H}, \mathcal{P}$ and the estimators $\hat{\Phi}, \hat{r}, \hat{\mathcal{P}}$, we have

$$\mathbb{E}_{\mathbf{z}_{1:M}} \left[\mathcal{R}(\hat{h}) - \mathcal{R}(h^*) \right] \leq \frac{D^2}{2M\alpha_M} + \frac{\xi}{M} \sum_{m=1}^M \alpha_m + \tau \sqrt{\zeta},$$

Where

$$\zeta = \left\| \Phi - \hat{\Phi} \right\|_{L^2(\mathbb{P}_X \otimes \mathbb{P}_Z)}^2 + \|r - \hat{r}\|_{L^2(Z)}^2 + \left\| \mathcal{P} - \hat{\mathcal{P}} \right\|_{\text{op}}^2,$$

$$\xi = \frac{3}{2} \left\| \hat{\Phi} \right\|_{\infty}^2 \left(C_0^2 + L^2 \|\hat{r}\|_{L^2(Z)}^2 + L^2 D^2 \left\| \hat{\mathcal{P}} \right\|_{\text{op}}^2 \right),$$

$$\tau = 2D \max \left\{ 3(C_0^2 + L^2 \mathbb{E}[Y^2] + L^2 D^2), 2L^2 \left\| \hat{\Phi} \right\|_{\infty}^2, 2L^2 D^2 \left\| \hat{\Phi} \right\|_{\infty}^2 \right\}.$$

Theory

- The bound

$$\mathbb{E}_{\mathbf{z}_{1:M}} \left[\mathcal{R}(\hat{h}) - \mathcal{R}(h^*) \right] \leq \frac{D^2}{2M\alpha_M} + \frac{\xi}{M} \sum_{m=1}^M \alpha_m + \tau \sqrt{\zeta},$$

suggests (α_m) should satisfy

$$M\alpha_M \rightarrow \infty \quad \text{and} \quad \frac{1}{M} \sum_{i=1}^M \alpha_m \rightarrow 0.$$

Theory

- The bound

$$\mathbb{E}_{\mathbf{z}_{1:M}} \left[\mathcal{R}(\hat{h}) - \mathcal{R}(h^*) \right] \leq \frac{D^2}{2M\alpha_M} + \frac{\xi}{M} \sum_{m=1}^M \alpha_m + \tau \sqrt{\zeta},$$

suggests (α_m) should satisfy

$$M\alpha_M \rightarrow \infty \quad \text{and} \quad \frac{1}{M} \sum_{i=1}^M \alpha_m \rightarrow 0.$$

- Additional samples from just Z can already increase estimator's quality.

Practice

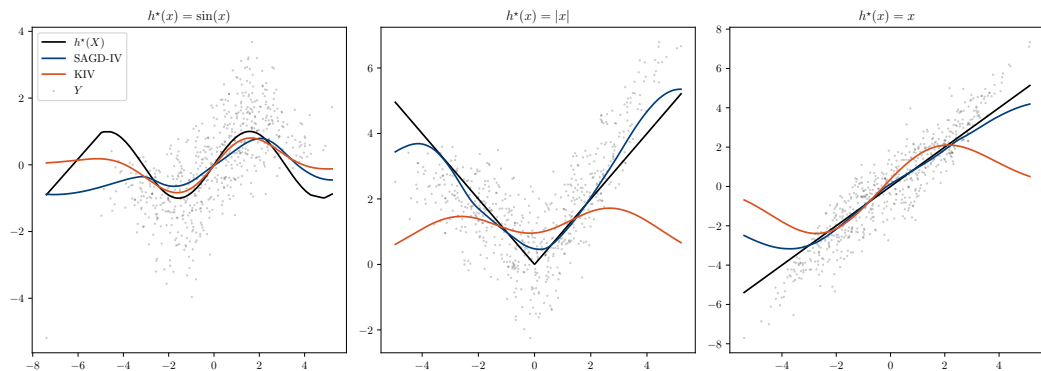


Figure: Benchmark against Kernel Instrumental Variable (KIV) [4]

Future Work

- Robust benchmarks against other recent methods;
- Application to discrete outcome models: $Y = \mathbf{1}\{h^*(X) + \varepsilon > 0\}$.

References

- [1] S. Darolles et al. “Nonparametric Instrumental Regression”. In: *Econometrica* 79.5 (2011), pp. 1541–5165.
- [2] Yuri R. Fonseca and Yuri F. Saporito. *Statistical Learning and Inverse Problems: A Stochastic Gradient Approach*. 2022. arXiv: 2209.14967 [stat.ML].
- [3] Whitney K. Newey and James L. Powell. “Instrumental Variable Estimation of Nonparametric Models”. In: *Econometrica* 71.5 (2003), pp. 1565–1578. ISSN: 00129682, 14680262. URL: <http://www.jstor.org/stable/1555512> (visited on 07/03/2023).
- [4] Rahul Singh, Maneesh Sahani, and Arthur Gretton. “Kernel Instrumental Variable Regression”. In: *Advances in Neural Information Processing Systems*. Vol. 32. Curran Associates, Inc., 2019.

Thank You!