

NPIV Estimation through Stochastic Gradients and Kernel Methods

Student: Caio Lins
Advisor: Yuri Saporito

EMAp – FGV

October 19, 2023



Summary

NPIV estimation

Our approach

Results

Next steps

Instrumental Variables

- ▶ Consider a generic regression problem:

$$Y = h^*(X) + \varepsilon,$$

where $\mathbb{E}[\varepsilon] = 0$ and we wish to estimate h^* .

Instrumental Variables

- ▶ Consider a generic regression problem:

$$Y = h^*(X) + \varepsilon,$$

where $\mathbb{E}[\varepsilon] = 0$ and we wish to estimate h^* .

- ▶ What happens if $\varepsilon \not\perp X$? More specifically, if $\mathbb{E}[\varepsilon | X] \neq 0$?

Instrumental Variables

- ▶ Consider a generic regression problem:

$$Y = h^*(X) + \varepsilon,$$

where $\mathbb{E}[\varepsilon] = 0$ and we wish to estimate h^* .

- ▶ What happens if $\varepsilon \not\perp X$? More specifically, if $\mathbb{E}[\varepsilon | X] \neq 0$?
- ▶ Minimizing $\text{MSE}(h) = \mathbb{E}[(Y - h(X))^2]$ over h gives biased results:

Instrumental Variables

- ▶ Consider a generic regression problem:

$$Y = h^*(X) + \varepsilon,$$

where $\mathbb{E}[\varepsilon] = 0$ and we wish to estimate h^* .

- ▶ What happens if $\varepsilon \not\perp X$? More specifically, if $\mathbb{E}[\varepsilon | X] \neq 0$?
- ▶ Minimizing $\text{MSE}(h) = \mathbb{E}[(Y - h(X))^2]$ over h gives biased results:

$$\arg \min_{W=h(X) \text{ for some } h} \mathbb{E}[(Y - W)^2] = \mathbb{E}[Y | X] = h^*(X) + \underbrace{\mathbb{E}[\varepsilon | X]}_{\neq 0}.$$

Instrumental Variables

- ▶ Consider a generic regression problem:

$$Y = h^*(X) + \varepsilon,$$

where $\mathbb{E}[\varepsilon] = 0$ and we wish to estimate h^* .

- ▶ What happens if $\varepsilon \not\perp X$? More specifically, if $\mathbb{E}[\varepsilon | X] \neq 0$?
- ▶ Minimizing $\text{MSE}(h) = \mathbb{E}[(Y - h(X))^2]$ over h gives biased results:

$$\arg \min_{W=h(X) \text{ for some } h} \mathbb{E}[(Y - W)^2] = \mathbb{E}[Y | X] = h^*(X) + \underbrace{\mathbb{E}[\varepsilon | X]}_{\neq 0}.$$

- ▶ We end up estimating $h^*(X) + \mathbb{E}[\varepsilon | X]$. Other problems may occur.

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. $p_{X|Z}(x | z)$ is not constant in z , that is, Z influences X ;

Instrumental Variables

- ▶ Suppose we have access to a variable Z such that
 1. $p_{X|Z}(x | z)$ is not constant in z , that is, Z influences X ;
 2. ε is uncorrelated with Z , that is, $\mathbb{E}[\varepsilon | Z] = 0$.

Instrumental Variables

- ▶ Suppose we have access to a variable Z such that
 1. $p_{X|Z}(x | z)$ is not constant in z , that is, Z influences X ;
 2. ε is uncorrelated with Z , that is, $\mathbb{E}[\varepsilon | Z] = 0$.

Instrumental Variables

- ▶ Suppose we have access to a variable Z such that
 1. $p_{X|Z}(x | z)$ is not constant in z , that is, Z influences X ;
 2. ε is uncorrelated with Z , that is, $\mathbb{E}[\varepsilon | Z] = 0$.

Z is called an *instrumental variable*.

Instrumental Variables

- ▶ Suppose we have access to a variable Z such that
 1. $p_{X|Z}(x | z)$ is not constant in z , that is, Z influences X ;
 2. ε is uncorrelated with Z , that is, $\mathbb{E}[\varepsilon | Z] = 0$.

Z is called an *instrumental variable*.

- ▶ How does it help us?

Instrumental Variables

- Structural equation:

$$Y = h^*(X) + \varepsilon,$$

where $\mathbb{E}[\varepsilon \mid Z] = 0$.

Instrumental Variables

- ▶ Structural equation:

$$Y = h^*(X) + \varepsilon,$$

where $\mathbb{E}[\varepsilon \mid Z] = 0$.

- ▶ Consider minimizing $\mathcal{R}(h) = \mathbb{E} \left[(\mathbb{E}[Y - h(X) \mid Z])^2 \right]$ over h .

Instrumental Variables

- Structural equation:

$$Y = h^*(X) + \varepsilon,$$

where $\mathbb{E}[\varepsilon \mid Z] = 0$.

- Consider minimizing $\mathcal{R}(h) = \mathbb{E} \left[(\mathbb{E}[Y - h(X) \mid Z])^2 \right]$ over h .
- Compare:

$$\text{MSE}(h) = \mathbb{E}[(Y - h(X))^2] \quad \text{v.s.} \quad \mathcal{R}(h) = \mathbb{E} \left[(\mathbb{E}[Y - h(X) \mid Z])^2 \right].$$

Instrumental Variables

► Since

$$\mathbb{E}[Y | Z] = \mathbb{E}[h^*(X) + \varepsilon | Z] = \mathbb{E}[h^*(X) | X] + \cancel{\mathbb{E}[\varepsilon | Z]} \stackrel{0}{=} \mathbb{E}[h^*(X) | Z],$$

Instrumental Variables

► Since

$$\mathbb{E}[Y | Z] = \mathbb{E}[h^*(X) + \varepsilon | Z] = \mathbb{E}[h^*(X) | X] + \cancel{\mathbb{E}[\varepsilon | Z]} \overset{0}{=} \mathbb{E}[h^*(X) | Z],$$

We have

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

Instrumental Variables

► Since

$$\mathbb{E}[Y | Z] = \mathbb{E}[h^*(X) + \varepsilon | Z] = \mathbb{E}[h^*(X) | X] + \cancel{\mathbb{E}[\varepsilon | Z]} \overset{0}{=} \mathbb{E}[h^*(X) | Z],$$

We have

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

► $\mathcal{R}(h) = 0 \iff \mathbb{E}[(h^* - h)(X) | Z] = 0 \iff \mathbb{E}[h^*(X) | Z] = \mathbb{E}[h(X) | Z].$

Example

$$\underbrace{\text{Grades}}_Y = h^*(\underbrace{\text{Attends tutoring sessions?}}_X) + \varepsilon.$$

Example

$$\underbrace{\text{Grades}}_Y = h^*(\underbrace{\text{Attends tutoring sessions?}}_X) + \varepsilon.$$

- Natural ability is a confounding variable: maybe only people who struggle a lot go to tutoring sessions.

Example

$$\underbrace{\text{Grades}}_Y = h^*(\underbrace{\text{Attends tutoring sessions?}}_X) + \varepsilon.$$

- ▶ Natural ability is a confounding variable: maybe only people who struggle a lot go to tutoring sessions.
- ▶ Z = Lives close to school?

Example

$$\underbrace{\text{Grades}}_Y = h^*(\underbrace{\text{Attends tutoring sessions?}}_X) + \varepsilon.$$

- ▶ Natural ability is a confounding variable: maybe only people who struggle a lot go to tutoring sessions.
- ▶ $Z = \text{Lives close to school?}$
 1. $Z \not\perp\!\!\!\perp X$,
 2. $\varepsilon \perp\!\!\!\perp Z$.

NPIV estimation

- ▶ Stands for “Nonparametric Instrumental Variable estimation”.

NPIV estimation

- ▶ Stands for “Nonparametric Instrumental Variable estimation”.
- ▶ No assumptions about some parametric form for h^* .

Summary

NPIV estimation

Our approach

Results

Next steps

Problem Formulation

- We have

$$Y = h^*(X) + \varepsilon, \tag{1}$$

where $\mathbb{E}[\varepsilon \mid Z] = 0$, and want to estimate h^* .

Problem Formulation

- We have

$$Y = h^*(X) + \varepsilon, \quad (1)$$

where $\mathbb{E}[\varepsilon \mid Z] = 0$, and want to estimate h^* .

- Let $\mathcal{P}[h] = \mathbb{E}[h(X) \mid Z = \cdot]$, that is

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

Problem Formulation

- We have

$$Y = h^*(X) + \varepsilon, \quad (1)$$

where $\mathbb{E}[\varepsilon \mid Z] = 0$, and want to estimate h^* .

- Let $\mathcal{P}[h] = \mathbb{E}[h(X) \mid Z = \cdot]$, that is

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

- Let $r_0(Z) = \mathbb{E}[Y \mid Z]$. Then (1) is equivalent to:

$$\mathbb{E}[Y \mid Z] = \mathbb{E}[h^*(X) \mid Z] \iff r_0 = \mathcal{P}[h^*].$$

Problem Formulation

- ▶ We have

$$Y = h^*(X) + \varepsilon, \quad (1)$$

where $\mathbb{E}[\varepsilon \mid Z] = 0$, and want to estimate h^* .

- ▶ Let $\mathcal{P}[h] = \mathbb{E}[h(X) \mid Z = \cdot]$, that is

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

- ▶ Let $r_0(Z) = \mathbb{E}[Y \mid Z]$. Then (1) is equivalent to:

$$\mathbb{E}[Y \mid Z] = \mathbb{E}[h^*(X) \mid Z] \iff r_0 = \mathcal{P}[h^*].$$

- ▶ We wish to “invert” \mathcal{P} .

- ▶ Risk measure:

$$\mathcal{R}(h) = \mathbb{E} \left[\frac{1}{2} (\mathbb{E}[Y - h(X) \mid Z])^2 \right] = \mathbb{E} \left[\frac{1}{2} (r_0(Z) - \mathcal{T}[h](Z))^2 \right].$$

Problem Formulation

- Our risk measure:

$$\begin{aligned}\mathcal{R}(h) &= \frac{1}{2} \mathbb{E} \left[(\mathbb{E}[Y - h(X) \mid Z])^2 \right] \\ &= \frac{1}{2} \mathbb{E} \left[(r_0(Z) - \mathcal{P}[h](Z))^2 \right] \\ &= \frac{1}{2} \mathbb{E} \left[(r_0 - \mathcal{P}[h])(Z)^2 \right] .\end{aligned}$$

Stochastic Gradients

- We showed that $\nabla \mathcal{R}(h) = \mathcal{P}^*[\mathcal{P}[h] - r_0]$, where

$$\mathcal{P}^*[f](X) = \mathbb{E}[f(Z) \mid X].$$

Stochastic Gradients

- ▶ We showed that $\nabla \mathcal{R}(h) = \mathcal{P}^*[\mathcal{P}[h] - r_0]$, where

$$\mathcal{P}^*[f](X) = \mathbb{E}[f(Z) \mid X].$$

- ▶ Immediate idea:

$$\begin{cases} h_0 \equiv 0, \\ h_t \leftarrow h_{t-1} - \alpha_t \nabla \mathcal{R}(h_{t-1}) \quad \text{for } t \geq 1. \end{cases}$$

Stochastic Gradients

- ▶ We showed that $\nabla \mathcal{R}(h) = \mathcal{P}^*[\mathcal{P}[h] - r_0]$, where

$$\mathcal{P}^*[f](X) = \mathbb{E}[f(Z) \mid X].$$

- ▶ Immediate idea:

$$\begin{cases} h_0 \equiv 0, \\ h_t \leftarrow h_{t-1} - \alpha_t \nabla \mathcal{R}(h_{t-1}) \quad \text{for } t \geq 1. \end{cases}$$

- ▶ There are problems...

Stochastic Gradients

- ▶ Another idea: notice that

$$\begin{aligned}\mathcal{P}^*[f](x) &= \mathbb{E}[f(Z) \mid X = x] \\ &= \int_{\mathcal{Z}} f(z) p(z \mid x) \, dz \\ &= \int_{\mathcal{Z}} f(z) \frac{p(x, z)}{p(x)p(z)} p(z) \, dz \\ &= \mathbb{E}_Z[f(Z)\Phi(x, Z)],\end{aligned}$$

where $\Phi(x, z) = \frac{p(x, z)}{p(x)p(z)}$.

Stochastic Gradients

- ▶ Another idea: notice that

$$\begin{aligned}\mathcal{P}^*[f](x) &= \mathbb{E}[f(Z) \mid X = x] \\ &= \int_{\mathcal{Z}} f(z) p(z \mid x) \, dz \\ &= \int_{\mathcal{Z}} f(z) \frac{p(x, z)}{p(x)p(z)} p(z) \, dz \\ &= \mathbb{E}_Z[f(Z)\Phi(x, Z)],\end{aligned}$$

where $\Phi(x, z) = \frac{p(x, z)}{p(x)p(z)}$.

- ▶ In the spirit of SGD, $f(Z)\Phi(x, Z)$ is a stochastic estimate for $\mathcal{P}^*[f](x)$.

Stochastic Gradients

- Substitute $\nabla \mathcal{R}(h) = \mathcal{P}^*[\mathcal{P}[h] - r_0]$ for

$$\hat{\Phi}(\cdot, Z) \left(\hat{\mathcal{P}}[h](Z) - \hat{r}_0(Z) \right).$$

Stochastic Gradients

- ▶ Substitute $\nabla \mathcal{R}(h) = \mathcal{P}^*[\mathcal{P}[h] - r_0]$ for

$$\hat{\Phi}(\cdot, Z) \left(\hat{\mathcal{P}}[h](Z) - \hat{r}_0(Z) \right).$$

- ▶ New algorithm:

$$\begin{cases} h_0 \equiv 0, \\ h_t \leftarrow h_{t-1} - \alpha_t \hat{\Phi}(\cdot, z_t) \left(\hat{\mathcal{P}}[h](z_t) - \hat{r}_0(z_t) \right) \quad \text{for } t \geq 1. \end{cases}$$

Stochastic Gradients

- ▶ Substitute $\nabla \mathcal{R}(h) = \mathcal{P}^*[\mathcal{P}[h] - r_0]$ for

$$\hat{\Phi}(\cdot, Z) \left(\hat{\mathcal{P}}[h](Z) - \hat{r}_0(Z) \right).$$

- ▶ New algorithm:

$$\begin{cases} h_0 \equiv 0, \\ h_t \leftarrow h_{t-1} - \alpha_t \hat{\Phi}(\cdot, z_t) \left(\hat{\mathcal{P}}[h](z_t) - \hat{r}_0(z_t) \right) \quad \text{for } t \geq 1. \end{cases}$$

- ▶ We chose to use RKHS (kernel) methods for computing $\hat{\Phi}$, $\hat{\mathcal{P}}$ and \hat{r}_0 .

Stochastic Gradients

- ▶ Substitute $\nabla \mathcal{R}(h) = \mathcal{P}^*[\mathcal{P}[h] - r_0]$ for

$$\hat{\Phi}(\cdot, Z) \left(\hat{\mathcal{P}}[h](Z) - \hat{r}_0(Z) \right).$$

- ▶ New algorithm:

$$\begin{cases} h_0 \equiv 0, \\ h_t \leftarrow h_{t-1} - \alpha_t \hat{\Phi}(\cdot, z_t) \left(\hat{\mathcal{P}}[h](z_t) - \hat{r}_0(z_t) \right) \quad \text{for } t \geq 1. \end{cases}$$

- ▶ We chose to use RKHS (kernel) methods for computing $\hat{\Phi}$, $\hat{\mathcal{P}}$ and \hat{r}_0 .
- ▶ One dataset with samples from (X, Z, Y) to compute $\hat{\Phi}$, $\hat{\mathcal{P}}$, \hat{r}_0 , and one dataset with samples from Z to conduct SGD loop.

Summary

NPIV estimation

Our approach

Results

Next steps

Practical Results

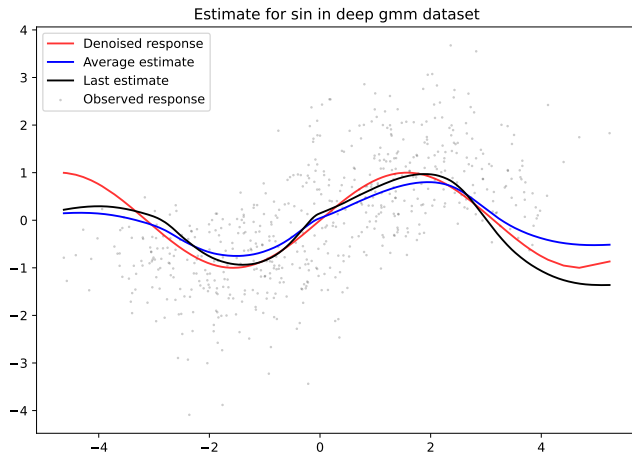


Figure: In red we have $h^* = \sin$, in black we have \hat{h}_N and in blue, $h = \frac{1}{N} \sum_{t=1}^N h_N$. Results produced with 600 joint samples of (X, Z, Y) and 2000 more samples of Z only.

Practical Results

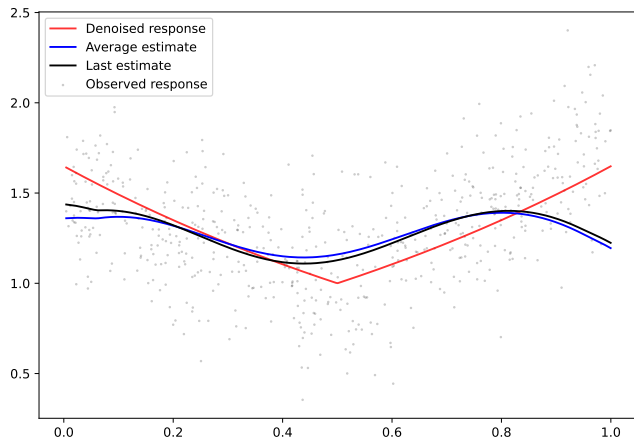


Figure: In red we have $h^*(x) = \exp(|x|)$, in black we have \hat{h}_N and in blue, $h = \frac{1}{N} \sum_{t=1}^N h_N$. Results produced with 600 joint samples of (X, Z, Y) and 2000 more samples of Z only.

Practical Results

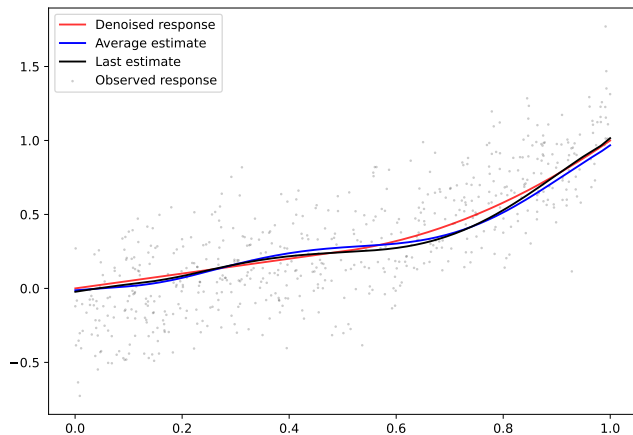


Figure: In red we have $h^*(x) = 2((x - 1/2)^+)^2 + x/2$, in black we have \hat{h}_N and in blue, $h = \frac{1}{N} \sum_{t=1}^N h_N$. Results produced with 600 joint samples of (X, Z, Y) and 2000 more samples of Z only.

Theoretical properties

Letting $\bar{h} = \frac{1}{N} \sum_{t=1}^N h_t$, we have

$$\begin{aligned} \mathbb{E}_{z_{1:N}} [\mathcal{R}(\bar{h})] &\leq \frac{D^2}{2N\alpha_N} + \mathcal{O}_p(1) \frac{1}{N} \sum_{t=1}^N \alpha_t \\ &\quad + \mathcal{O}_p(1) \left(\left\| \hat{\Phi} - \Phi \right\|^2 + \left\| \hat{r}_0 - r_0 \right\|^2 + \left\| \hat{\mathcal{P}} - \mathcal{P} \right\|^2 \right)^{1/2}. \end{aligned}$$

Theoretical properties

Letting $\bar{h} = \frac{1}{N} \sum_{t=1}^N h_t$, we have

$$\begin{aligned} \mathbb{E}_{z_{1:N}} [\mathcal{R}(\bar{h})] &\leq \frac{D^2}{2N\alpha_N} + \mathcal{O}_p(1) \frac{1}{N} \sum_{t=1}^N \alpha_t \\ &\quad + \mathcal{O}_p(1) \left(\left\| \hat{\Phi} - \Phi \right\|^2 + \left\| \hat{r}_0 - r_0 \right\|^2 + \left\| \hat{\mathcal{P}} - \mathcal{P} \right\|^2 \right)^{1/2}. \end{aligned}$$

Choose $(\alpha_t)_{t=1}^\infty$ so that $N\alpha_N \rightarrow \infty$ but $\frac{1}{N} \sum_{t=1}^N \alpha_t \rightarrow 0$ as $N \rightarrow \infty$.

Summary

NPIV estimation

Our approach

Results

Next steps

Next steps

- ▶ Benchmark against current methods.
- ▶ Discrete outcome models.

References

- [1] Yuri R. Fonseca and Yuri F. Saporito. *Statistical Learning and Inverse Problems: A Stochastic Gradient Approach*. 2022. arXiv: 2209.14967 [stat.ML].
- [2] 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).

Thank You!