

FUNDAÇÃO GETULIO VARGAS
SCHOOL OF APPLIED MATHEMATICS

CAIO F. LINS PEIXOTO

**NONPARAMETRIC INSTRUMENTAL VARIABLE REGRESSION
THROUGH KERNEL METHODS AND STOCHASTIC GRADIENTS**

Rio de Janeiro
2023

CAIO F. LINS PEIXOTO

**NONPARAMETRIC INSTRUMENTAL VARIABLE REGRESSION
THROUGH KERNEL METHODS AND STOCHASTIC GRADIENTS**

Bachelor's dissertation presented to
the School of Applied Mathematics
(FGV/EMAp) to obtain the Bachelor's
degree in Applied Mathematics.

Area of Study: Nonparametric Regression,
Instrumental Variables, Kernel Methods,
Stochastic Optimization, Machine Learning.

Advisor: Yuri F. Saporito

Rio de Janeiro

2023

Ficha catalográfica elaborada pela BMHS/FGV

Lins, Caio

Nonparametric Instrumental Variable Regression Through Kernel Methods and Stochastic Gradients/ Caio F. Lins Peixoto. – 2023.

22f.

Bachelor's Dissertation (Undergraduate) – School of Applied Mathematics.

Advisor: Yuri F. Saporito.

Includes bibliography.

1. Nonparametric Regression 2. Instrumental Variables 2. Stochastic Optimization I. Saporito, Yuri Fahham II. School of Applied Mathematics. III. Nonparametric Instrumental Variable Regression Through Kernel Methods and Stochastic Gradients

CAIO F. LINS PEIXOTO

NONPARAMETRIC INSTRUMENTAL VARIABLE REGRESSION THROUGH KERNEL METHODS AND STOCHASTIC GRADIENTS

Bachelor's dissertation presented to the School of Applied Mathematics (FGV/EMAp) to obtain the Bachelor's degree in Applied Mathematics.

Area of Study: Nonparametric Regression, Instrumental Variables, Kernel Methods, Stochastic Optimization, Machine Learning.

Approved on December —, 2023
By the organizing committee

Yuri F. Saporito
School of Applied Mathematics

Board Member 1
Institution 1

Board Member 2
Institution 2

I dedicate this thesis to ...

Acknowledgements

Thanks, ...

“ Biped! boped! bum! ”

Albert Einstein

Abstract

Keywords:

Resumo

Palavras-chave:

List of Figures

Figure 1 – Causal diagram for equation (2.1), where X is endogenous and Z is an	
IV. Source: prepared by the author.	15

List of Tables

Contents

1	INTRODUCTION	12
2	INSTRUMENTAL VARIABLE REGRESSION	13
2.1	Endogeneity	13
2.2	Instrumental Variables	14
2.3	Two Stages Least Squares (2SLS)	15
2.4	Nonparametric Instrumental Variable Regression	15
2.4.1	Problem specification	16
2.4.2	Identification	16
3	CONCLUSION	20
	References	21
	APPENDIX	22

1 Introduction

Remember to cite every person ([NEWHEY; POWELL, 2003](#)).

2 Instrumental Variable Regression

This chapter provides an introduction to both parametric and nonparametric instrumental variable regression. Its goal is twofold. Firstly, we want to introduce the subject to readers unfamiliar with it. To make the exposition more fluid, we chose to delay the precise definition of all mathematical objects involved until Section 2.4, which deals with the nonparametric version. Until then, only knowledge of non-measure theoretic probability theory is required. This is sufficient to understand the motivation behind instrumental variables, as well as the most widely used tool to perform regression using them: two stages least squares. The second goal is to precisely state the nonparametric regression problem which will be addressed in the remainder of this thesis and, in doing so, fixate our notation. From that point on, we assume the reader is familiar with measure theory and linear functional analysis. All of the less usual objects and results we make use of will be recalled in a corresponding appendix.

2.1 Endogeneity

We start by introducing the problem of endogenous covariates. The structural equation we consider is the following:

$$Y = h^*(X) + \varepsilon, \quad (2.1)$$

where X is a vector of explanatory variables, Y is the scalar response, ε is a zero mean noise and the function h^* is the structural parameter we would like to estimate. The simplest estimation method for this model specification — and, therefore, one we would like to be able to use — is ordinary least squares (OLS), which works by finding, within a given class of functions \mathcal{H} , the element which minimizes the mean squared error:

$$\hat{h} = \arg \min_{h \in \mathcal{H}} \mathbb{E}[(Y - h(X))^2]. \quad (2.2)$$

A reasonable and ample choice for \mathcal{H} is the set of all square-integrable functions of X , that is, such that $\mathbb{E}[h(X)^2] < \infty$. Under this choice, we recover the conditional expectation of Y given X , i.e., $\hat{h}(X) = \mathbb{E}[Y \mid X]$. Expanding Y through (2.1), we find that $\hat{h}(X) = h^*(X) + \mathbb{E}[\varepsilon \mid X]$. Hence, if $\mathbb{E}[\varepsilon \mid X]$ is not identically null, we have introduced bias in our estimation.

This is one of the problems which appear when $\mathbb{E}[\varepsilon \mid X] \neq 0$, or, more generally, when X and ε are correlated in some way. When this happens, we say that X is *endogenous*. There are several causes for endogenous covariates, the most common of which are (WOOLDRIDGE, 2001):

Omitted Variables This means ε can be decomposed as $g^*(W) + \eta$, where $\mathbb{E}[\eta \mid X, W] = 0$ a.s. and X and W are correlated. Hence, when we don't observe W and leave it to the error term, we end up estimating $h^*(X) + \mathbb{E}[g^*(W) \mid X]$. For example, if we want to regress a person's wage solely on her number of schooling years, there are other variables unaccounted for which influence both wages and schooling, such as natural ability. Innately skilled people may tend to be successful in school — and, therefore, pursue higher levels of education — as well as show higher performance in their future jobs, resulting in better wages.

Measurement Error If we are unable to exactly measure one of the covariates, X_k , and instead measure X'_k subject to some stochastic error, by using X'_k in our regression instead of X_k we are delegating to ε some measure of the difference between X_k and X'_k . Depending on how these two variables are related, we may introduce endogeneity. For example, X_k may be a marginal tax rate, but we may only have access to an average tax rate X'_k .

Simultaneity Simultaneity arises when one covariate X_k is determined simultaneously with Y . For example, if we are regressing neighborhood murder rates using the size of the local task force as a covariate, there is a simultaneity problem, since larger murder rates in a place cause a larger task force to be allocated there.

Bias in the estimation procedure is only one of the problems which arise when there are endogenous covariates. It's well known that the OLS estimate for linear regression fails to be consistent if any one of the covariates is endogenous (WOOLDRIDGE, 2001). To overcome endogeneity a few approaches exist, but by far the one most used by empirical economic research is instrumental variable estimation (WOOLDRIDGE, 2001).

2.2 Instrumental Variables

2.1 Definition An *instrumental variable* for regression problem (2.1) is a random variable Z such that

- (i) There is some influence of Z upon X , that is, the marginal distribution of X is not the same as the distribution of X conditioned on Z ;
- (ii) The conditional mean of ε given Z is almost surely null, i.e., $\mathbb{E}[\varepsilon \mid Z] = 0$.

The idea behind an instrumental variable is that it is exogenous (ii) while still influencing Y through X (i). An exogenous covariate, in contrast to an endogenous one, is a variable that is determined outside of the system described by (2.1).

Condition (ii) is only one of the possible meanings for the statement that Z is exogenous. Two possible alternatives are requiring that Z be (1) independent from, or (2)

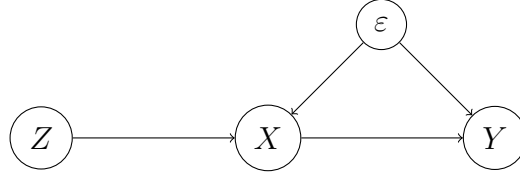


Figure 1 – Causal diagram for equation (2.1), where X is endogenous and Z is an IV.
Source: prepared by the author.

uncorrelated with ε . Of course, (1) is a much more strict requirement which implies (ii), while (2) is a softer condition, implied by (ii). Independence is almost always impossible to verify in real scenarios, so (1) is not a good option. In contrast, there are situations where condition (2) is enough for ensuring good properties of IV estimators, including one we will present shortly, the linear model (WOOLDRIDGE, 2001). However, in order to prepare grounds for the nonparametric methods that will come later, we chose to use the definition which serves both.

Instrumental variables are also studied in the context of causal inference, where the conditions above are presented differently, in terms of causal diagrams. In this field, instrumental variables are also required that to satisfy a third condition, phrased in terms of the causal diagram describing the relations between variables of interest (HERNÁN; ROBINS, 2020):

- (iii) All paths from Z to Y must pass through X , that is, Z *only* influences Y through X .

In this sense, a typical causal diagram for an IV problem is the one in Figure 1.

2.3 Two Stages Least Squares (2SLS)

In this section, we restrict the structural function h^* in (2.1) to be affine:

$$h^*(x) = \beta_0 + \beta_1 x_1 + \cdots + \beta_{d_X} x_{d_X}, \quad (2.3)$$

and assume to have access to a random variable Z , taking values in \mathbf{R}^{d_Z} , satisfying conditions 2.1 (i) and (ii), so that Z is a valid instrumental variable. Hence, our data is composed of n independent joint samples $\{(X_i, Z_i, Y_i)\}_{i=1}^n$. Let $\mathbf{X} \in \mathbf{R}^{n \times (d_X+1)}$ and $\mathbf{Z} \in \mathbf{R}^{n \times (d_Z+1)}$ be the experiment design matrices with 1's in the first column, and let $\mathbf{Y} \in \mathbf{R}^n$ be the vector with all observations of Y .

2.4 Nonparametric Instrumental Variable Regression

In nonparametric regression, we do not specify *a priori* a finite dimensional parametric form for the structural function (such as restricting it to be affine), and so we allow

our search space to potentially be infinite dimensional. However, in doing this, we must still precisely define the infinite dimensional space where the solution will be searched for. Hence, we start by precisely defining the nonparametric regression problem given by (2.1)

2.4.1 Problem specification

Let $(\Omega, \mathcal{A}, \mathbb{P})$ be the underlying probability space. Assume that $X : (\Omega, \mathcal{A}) \rightarrow (\mathbf{R}^{d_x}, \mathcal{B}(\mathbf{R}^{d_x}))$ and $\varepsilon : (\Omega, \mathcal{A}) \rightarrow (\mathbf{R}, \mathcal{B}(\mathbf{R}))$ are measurable¹ and, furthermore, that $\varepsilon \in L^1(\Omega, \mathcal{A}, \mathbb{P})$ with $\mathbb{E}[\varepsilon] = 0$. We also assume that $\mathbb{E}[\varepsilon \mid X]$ is *not* almost surely null and, hence, X is endogenous. Denote by \mathbb{P}_X the distribution of the random variable² X , that is, the pushforward measure $\mathbb{P} \circ X^{-1}$ defined on $\mathcal{B}(\mathbf{R}^{d_x})$. We write $L^2(X)$ as a shorthand for the space $L^2(\mathbf{R}^{d_x}, \mathcal{B}(\mathbf{R}^{d_x}), \mathbb{P}_X)$ of real and square integrable (equivalence classes of) measurable functions defined on the measure space $(\mathbf{R}^{d_x}, \mathcal{B}(\mathbf{R}^{d_x}), \mathbb{P}_X)$. It's important to recall that the inner product and norm in $L^2(X)$ are given by $\langle h, g \rangle_{L^2(X)} = \mathbb{E}[h(X)g(X)]$ and $\|h\|_{L^2(X)}^2 = \langle h, h \rangle_{L^2(X)} = \mathbb{E}[h(X)^2]$.

We assume there exists $h^* \in L^2(X)$ such that (2.1) holds, that is, $Y = h^*(X) + \varepsilon$. Finally, we assume there exists a random variable $Z : (\Omega, \mathcal{A}) \rightarrow (\mathbf{R}^{d_z}, \mathcal{B}(\mathbf{R}^{d_z}))$ such that Z qualifies as an instrumental variable, i.e., Z satisfies conditions 2.1 (i) and (ii). We define \mathbb{P}_Z and $L^2(Z)$ in an manner analogous to \mathbb{P}_X and $L^2(X)$. Our goal is to estimate h^* based on i.i.d. samples from the joint distribution of X, Z and Y .

2.4.2 Identification

An important question to ask after specifying the problem is whether the function h^* is identified. The answer is negative without further assumptions, which will be presented in this subsection. This discussion was inspired on the second section of (NEWHEY; POWELL, 2003).

Suppose there exists $\delta \in L^2(X)$ such that $\delta \neq 0$, but $\mathbb{E}[\delta(X) \mid Z] = 0$. Without loss of generality, we can assume³ $\delta(X) \neq \mathbb{E}[\varepsilon \mid X]$. Defining $g \triangleq h^* + \delta$ and $\eta \triangleq \varepsilon - \delta(X)$, we have

$$Y = g(X) + \eta,$$

where $\mathbb{E}[\eta \mid Z] = 0$ and $\mathbb{E}[\eta \mid X] \neq 0$. Hence, $g \neq h^*$ but they are indistinguishable from the data generating process' perspective. Reciprocally, suppose that the only member of $L^2(X)$ which has null mean conditioned on Z is the null function. Then, given $g \in L^2(X)$

¹ We denote by $\mathcal{B}(\mathbf{R}^k)$ the Borel σ -algebra in \mathbf{R}^k .

² We use the term "random variable" when referring to scalar or vector valued measurable functions defined on (Ω, \mathcal{A}) .

³ If it happens to be the case that $\mathbb{E}[\mathbb{E}[\varepsilon \mid X] \mid Z] = 0$ (which is *not* implied by our assumptions so far), we can simply take $\delta(X) = \lambda \mathbb{E}[\varepsilon \mid X]$, for some $\lambda \in \mathbf{R} \setminus \{0, 1\}$. Since, by hypothesis, $\mathbb{E}[\varepsilon \mid X] \neq 0$, this satisfies our requirements and is different from $\mathbb{E}[\varepsilon \mid X]$.

such that

$$Y = g(X) + \eta$$

with $\mathbb{E}[\eta \mid Z] = 0$, we have

$$0 = (g - h^*)(X) + \eta - \varepsilon.$$

Conditioning on Z , we get

$$\mathbb{E}[(g - h^*)(X) \mid Z] = 0,$$

which, by assumption, implies $h^* = g$.

Therefore, a necessary and sufficient condition for identification of our regression problem is the following:

Assumption (Identification) If $\delta \in L^2(X)$ satisfies $\mathbb{E}[\delta(X) \mid Z] = 0$, then $\delta = 0$.

This condition has an interpretation in terms of the conditional expectation operator, which will be a key object in the construction of our estimator for h^* . Let $h \in L^2(X)$. Notice that, by Jensen's inequality,

$$\mathbb{E}[(\mathbb{E}[h(X) \mid Z])^2] \leq \mathbb{E}[\mathbb{E}[h(X)^2 \mid Z]] = \mathbb{E}[h(X)^2] < +\infty. \quad (2.4)$$

Furthermore, since $\mathbb{E}[h(X) \mid Z]$ is a $\sigma(Z)$ -measurable⁴ random variable, by the Doob-Dynkin Lemma there exists a measurable function $f_h : (\mathbf{R}^{dz}, \mathcal{B}(\mathbf{R}^{dz})) \rightarrow (\mathbf{R}, \mathcal{B}(\mathbf{R}))$ such that

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

Under these conditions, we write $\mathbb{E}[h(X) \mid Z = z]$ for $f_h(z)$. The computation in (2.4) shows that $f_h \in L^2(Z)$ for every $h \in L^2(X)$. Therefore, we can define the operator $\mathcal{P} : L^2(X) \rightarrow L^2(Z)$ given by $\mathcal{P}[h] = f_h = \mathbb{E}[h(X) \mid Z = \cdot]$. This operator, called the *conditional expectation operator*, is clearly linear and, again by (2.4), also bounded, satisfying $\|\mathcal{P}\|_{\text{op}} \leq 1$. The identification assumption thus amounts to saying that the kernel of \mathcal{P} is trivial, i.e., \mathcal{P} is injective.

It is hard to quantify how restrictive this condition is for arbitrary X and Z , so we analyze it in a more familiar setting, the exponential family. We will use a classic completeness result for statistics in this family of distributions to reformulate the identification assumption in terms of more familiar objects. We first define completeness:

2.2 Definition (LEHMANN, 1959) We say that a family \mathcal{P} of probability distributions on a measurable space (E, \mathcal{E}) is *complete* if

$$\int_E f(x) P(dx) = 0 \quad \text{for all } P \in \mathcal{P}$$

implies $f(x) = 0$ \mathcal{P} -a.e.⁵.

⁴ We denote by $\sigma(Z)$ the smallest σ -algebra in Ω with respect to which Z is measurable.

⁵ We say that a statement $Q(x)$ is true \mathcal{P} -a.e. if there exists a set $N \in \mathcal{E}$ such that $Q(x)$ is true for $x \in E \setminus N$ and $P(N) = 0$ for all $P \in \mathcal{P}$.

Then, a remarkable fact about the exponential family is the completeness of the natural statistics under a mild condition on the set of parameters:

2.3 Theorem (LEHMANN, 1959) Let Ξ be a subset of an Euclidian space with nonempty interior. Let X be a random vector with distribution P^θ parametrized by $\theta \in \Xi$ in the following manner:

$$P^\theta(dx) = C(\theta) \exp \left\{ \sum_{i=1}^s \theta_i T_i(x) \right\} \mu(dx),$$

where μ is the underlying measure. Then, the family \mathcal{P}_T , formed by the distributions of the random vector $T(X) = (T_1(X), \dots, T_s(X))$ as θ ranges through Ξ , is complete.

Using this result and under suitable hypotheses, we can reformulate the identification condition.

2.4 Theorem For $z \in \mathbf{R}^{d_z}$, let $\mathbb{Q}_z : \mathcal{B}(\mathbf{R}^{d_x}) \rightarrow [0, 1]$ denote the conditional distribution of X given $Z = z$. Assume there exists $U \in \mathcal{B}(\mathbf{R}^{d_z})$ such that $\mathbb{P}_Z(U) = 1$ and for all $z \in U$ we have

$$\mathbb{Q}_z(dx) = C(z) \exp(\alpha(z)^\top T(x)) \mu(dx)$$

for an underlying measure μ on $\mathcal{B}(\mathbf{R}^{d_x})$ and some functions $\alpha : \mathbf{R}^{d_z} \rightarrow \mathbf{R}^s$ and $T : \mathbf{R}^{d_x} \rightarrow \mathbf{R}^s$. Assume that T is injective and that the image of α restricted to U contains an open set. Then, h^* is identified.

Proof. Taking $\Xi = \alpha(U)$ and $\theta = \alpha(z)$, we see that the hypotheses of Theorem 2.3 are satisfied and, hence,

$$\mathcal{P} \triangleq \{\mathbb{Q}_z \circ T^{-1} : z \in U\}$$

is a complete family of probability distributions. Let $h \in L^2(X)$ be such that $\mathcal{P}[h] = 0$, i.e., $\mathbb{E}[h(X) \mid Z] = 0$ \mathbb{P}_Z -a.s. This means that the function

$$z \mapsto \int_{\mathbf{R}^{d_x}} h(x) \mathbb{Q}_z(dx)$$

is null \mathbb{P}_Z -a.s. Without loss of generality, we may assume that its null on U . But notice that, since T is injective, we can rewrite this integral as

$$\begin{aligned} 0 &= \int_{\mathbf{R}^{d_x}} h(x) \mathbb{Q}_z(dx) = \int_{\mathbf{R}^{d_x}} (h \circ T^{-1})(T(x)) \mathbb{Q}_z(dx) \\ &= \int_{\mathbf{R}^s} (h \circ T^{-1})(t) (\mathbb{Q}_z \circ T^{-1})(dt) \end{aligned}$$

for all $z \in U$ and some left inverse T^{-1} of T . By completeness of \mathcal{P} , this implies $h \circ T^{-1}(t) = 0$ \mathcal{P} -a.s. which, in turn, means that for all $z \in U$ we have

$$\begin{aligned} 1 &= (\mathbb{Q}_z \circ T^{-1})[(h \circ T^{-1})(t) = 0] = \mathbb{Q}_z[(h \circ T^{-1})(T(x)) = 0] \\ &= \mathbb{Q}_z[h(x) = 0]. \end{aligned}$$

Now, by the definition of conditional probability we have

$$\begin{aligned}\mathbb{P}_X[h(x) = 0] &= \int_{\mathbf{R}^{d_z}} \mathbb{Q}_z[h(x) = 0] \mathbb{P}_Z(dz) \\ &= \int_U \mathbb{Q}_z[h(x) = 0] \mathbb{P}_Z(dz) \\ &= 0.\end{aligned}$$

Therefore, h is the null function, which means h^* is identified. □

Check
multi-
variate
normal
case?

3 Conclusion

References

HERNÁN, Miguel A.; ROBINS, James M. **Causal Inference: What If**. [S.l.]: Chapman & Hall/CRC, 2020.

LEHMANN, E. L. **Testing Statistical Hypotheses**. [S.l.]: John Wiley & Sons, 1959.

NEWY, Whitney K.; POWELL, James L. Instrumental Variable Estimation of Nonparametric Models. **Econometrica**, v. 71, n. 5, p. 1565–1578, 2003. DOI: <http://dx.doi.org/10.1111/1468-0262.00459>.

WOOLDRIDGE, Jeffrey M. **Econometric Analysis of Cross Section and Panel Data**. [S.l.]: The MIT Press, 2001. ISBN 9780262232197.

Appendix