

Granular Instrumental Variables*

Xavier Gabaix and Ralph S.J. Koijen

Preliminary and incomplete. Comments welcome.

Abstract

In many settings, there is a dearth of instruments, which hampers economists’ ability to investigate causal relations. We propose a quite general way to construct instruments: “granular instrumental variables” (GIVs). In the economies we study, a few large firms or countries account for a large share of economic activity. As they are large, their idiosyncratic shocks affect aggregate outcomes. This makes those idiosyncratic shocks valid instruments for aggregate shocks. We provide a methodology to extract idiosyncratic shocks from the data, this way creating GIVs. Those GIVs allow us to then estimate parameters of interest, including causal elasticities.

We first illustrate the idea in a basic supply and demand framework: we achieve a novel identification of supply and demand elasticities, based on idiosyncratic shocks to supply or demand. We then show how the procedure can be adapted to handle many enrichments. We provide initial illustrations of the procedure with two applications. First, we measure how shocks to domestic banks causally affect sovereign yields. We document how negative shocks to Italian banks adversely affect Italian government bond yields, and vice-versa. This gives the first causal measure of the “doom loop” between banks and sovereign yields. Second, we estimate short-term supply and demand elasticities in the oil market. Our estimates match well existing estimates that use much more complex and labor-intensive (e.g., narrative) methods.

We sketch how GIVs could be useful to estimate a host of other causal parameters in economics, particularly in aggregate macro-finance contexts where instruments are usually very rare.

*xgabaix@fas.harvard.edu, ralph.koijen@chicagobooth.edu. We thank Antonio Coppola, Rishab Guha, Dong Ryeol Lee, and Ashesh Rambachan for excellent research assistance. We thank Isaiah Andrews, Gary Chamberlain, Mark Gertler, Oleg Itskhoki, Serena Ng, Jim Stock, Harald Uhlig, Dacheng Xiu and seminar participants at Chicago Booth, Harvard, NY Fed, NYU Stern, the University of Virginia, and the NBER (conference on “the rise of mega-firms”) for comments.

1 Introduction

In many settings, there is a dearth of instruments, which hampers economists’ understanding of causal relations (Ramey (2016); Nakamura and Steinsson (2018); Stock and Watson (2016, 2018); Chodorow-Reich (2017)). We propose a general way to construct instruments: “granular instrumental variables” (GIVs). Those instruments in turn allow to tease apart causal relations in a wide variety of economic contexts.

In the economies we study, many decisions are taken by a few large actors, e.g. firms, industries or countries. As has been observed before, those actors are often large, and their idiosyncratic shocks (e.g., productivity shocks) affect the aggregate ones.¹

Those idiosyncratic firm- or country-level shocks are valid instruments for the aggregate shocks. Now, we need a methodology to extract those idiosyncratic shocks, and the paper presents such a methodology. This creates GIVs. Those GIVs then allow us to estimate parameters of interest.

We first illustrate the idea in a basic static setup with supply and demand (Section 2). It is a classic case, and we show how GIVs allow for a novel estimation procedure: they yield an instrument that allows us to estimate the elasticities of supply and demand. Indeed, idiosyncratic demand shocks to large firms or countries give a valid instrument for demand change – and thus allow one to estimate the elasticity of supply. They also allow us to estimate the elasticity of demand: the idiosyncratic demand shock of a large firm impacts the price, which changes the demand of other firms. We formalize those ideas, and present a way to “optimally extract” idiosyncratic shocks, this way constructing optimal GIVs.

We will see that some conditions are needed to obtain valid GIVs. Hence, GIVs are not quite a “free lunch” in constructing instruments, but a “very cost-effective lunch” that yields instruments at a modest cost.

Once the ideas are in place, we show in Section 3 how the procedure can be extended to handle many enrichments, such as feedback loops, heterogeneity, and several exogenous factors. We specify the procedure within this general framework.

Empirical illustrations We provide preliminary empirical results for two applications: doom loops and the equilibrium of global crude oil markets. First, we study the link between fragility in the banking sector and sovereign bond markets, so-called “financial doom loops.” The European sovereign debt crisis highlighted the possibility that weakness in the banking sector may spill over to the sovereign, and vice versa. However, identifying contagion in a causal way has been challenging thus far. We study Italy during the European sovereign debt crisis and its aftermath. We show that the Italian banking sector, like many financial sectors in other countries, is highly concentrated. Moreover, as is well known, the correlation between banking returns and sovereign yield changes

¹Hence, economies are “granular,” that is, and their shocks made of incompressible “grains” of economic shocks, those are firm or industry-level shocks. This theme is laid out in Gabaix (2011), and developed in Acemoglu et al. (2012); di Giovanni and Levchenko (2012); Carvalho and Grassi (2019).

turned negative during the crisis.

We show how to construct GIVs in this setting using bank-level stock returns. Our results identify a causal link between the sovereign and the banking sectors during the sovereign debt crisis, but this effect becomes much smaller since December 2014 when the large-scale asset purchase program was announced by the European Central Bank. In terms of magnitudes, we find that, during the peak of the crisis, a “primitive” 40% negative return on banking stocks leads to a 200 basis points increase in the sovereign yield spread (over Germany, considered to be the safe bond) when the doom loop is operative, and via the doom loop, a -60% total return for bank stocks. Those are local effects, which increase as banks get closer to bankruptcy. To the best of our knowledge, this is the first causally estimated measure of doom loops.

Second, we use GIVs to estimate the short-term demand and supply elasticities in the global crude oil market. Since the seminal work by [Kilian \(2009\)](#), who uses an ordered VAR to identify the shocks in a structural VAR, an active literature explores sign restrictions, informative priors, and narrative methods to estimate these elasticities (see for instance [Kilian and Murphy \(2014\)](#), [Baumeister and Hamilton \(2018\)](#), and [Caldara et al. \(2018\)](#)). We use country-level oil supply growth to construct the GIV, after removing common factors using principal components and OPEC membership to construct an OPEC factor. We find that the granularly identified elasticities are in the range documented in the literature. Moreover, given the apparent importance of demand shocks in crude oil markets during the last 15 years, future work can use disaggregated data on (net) imports, inventories, and oil consumption to sharpen the estimates and to estimate more general models to understand price fluctuations in oil markets.

Uses of GIVs GIVs allow to “democratize” and “automatize” instruments, especially in macro-finance where they have been rare. In standard practice, finding an instrument is a heroic and very ingenious affair, e.g. the “China shock” (entry of China in the World Trade Organization, [Autor et al. \(2013\)](#)) depends on detailed historical knowledge, and applies only to a specific time period. With GIVs, we can have a more systematic way to obtain instruments, that can apply more generally and over many time periods.

Once one thinks about causality and GIV procedures, the answers to many interesting questions feel suddenly within reach. We sketch a few here, hoping that they will inspire other researchers to investigate those and related topics with help of GIVs.

If the Turkish Lira (to take a concrete example) appreciates, how does that affect Turkey’s exports and borrowing? One could handle that via idiosyncratic demand shocks by large investment funds for the Turkish Lira ([Kojen and Yogo \(forthcoming\)](#) provide a methodology for demand systems).

If there is an export boom, what’s the impact on the exchange rate, and the rest of the economy? Idiosyncratic shocks to large exporters will be useful for that, as recent research has shown them to be very large ([di Giovanni et al. \(2014\)](#); [Gaubert and Itskhoki \(2018\)](#); [Kramarz et al. \(2016\)](#)).

Do firm-specific hiring, investment and innovations spill over to peer firms operating in the same product market (i.e., what is the sign and magnitude of strategic complementarities)? Idiosyncratic innovation shocks to some large players will help construct the GIV.

How much do constraints of financial intermediaries (e.g., broker dealers) matter for asset prices? The GIV will rely on idiosyncratic shocks to intermediary wealth, which may be related to shocks to other parts of the banks.

How much do international shock propagate (e.g., how does a boom in Germany transmit to the rest of Europe)? Using idiosyncratic shocks (differentiating between productivity and demand shocks) to countries will help us answer that question.

Likewise, how do regional “micro” shocks propagate into macro shocks? GIVs allow to measure that, and also estimate a micro-to-macro multiplier.²

Related literature We relate to a number of literatures. We offer some brief pointers here, while offering a longer discussion in Section 6.3.

Instruments for macro. An active literature discusses identification strategies in macro (Ramey (2016); Nakamura and Steinsson (2018); Stock and Watson (2018); Chodorow-Reich (2017)). We add to it, by proposing to use GIVs, which are quite ubiquitous. There are lots of idiosyncratic shocks, and GIVs allow us to construct them systematically.

Origins of aggregate fluctuations. We provide tools that can tease that apart when you have loops. This literature finds that a sizable amount of volatility is “granular” in nature – coming from idiosyncratic shocks to firms or industries (see Long and Plosser (1983); Gabaix (2011), Acemoglu et al. (2012); di Giovanni and Levchenko (2012); di Giovanni et al. (2014); Baqaee and Farhi (2018); Pasten et al. (2017); Carvalho and Grassi (2019)). Those datasets can be revisited forming GIVs which allow us to investigate causal relations. Gabaix (2011) introduces the notion of “granular residual” – a weighted sum of proxies for idiosyncratic shocks, and shows how idiosyncratic shocks to firms appear to explain about one third of GDP fluctuations. But that paper does not take the crucial step to use this kind of concept as an instrument to measure causal relations, e.g. in a demand and supply setting.

The idea that we propose is, in retrospect, so natural that we suspected that it may have been already proposed in the literature, perhaps in a forgotten paper in the 1940s. However, after searching the literature and consulting with many experts, we could not find it. We are quite sure that this idea has not been systematically implemented in mainstream economic applications. The idea to use idiosyncratic shocks as instruments to estimate spillover effects has been explored in several creative papers, as we discuss in more detail in section 6.3, such as Leary and Roberts (2014), Amiti and Weinstein (2018), and Amiti et al. (2019). However, the typical approach is to use idiosyncratic shocks to other variables that are excluded from the main estimating equation to instrument for the actions of competing firms. We, instead, use the idiosyncratic shocks in

²We are pursuing this last question in ongoing work.

the estimating equation directly. In addition, we allow for more flexible exposures to unobserved common shocks in extracting idiosyncratic shocks.

Plan The reader is encouraged to read Section 2, Section 4 and Section 5 at first, and then go to Section 3, which contains the general procedure, in a second reading. Section 6 presents a number of extensions and robustness checks, and discusses more tightly the link with the rest of the literature. Section 7 presents simulations. Section 8 concludes. Long proofs are in Section 9 and the online appendix.

2 The basics

2.1 Notations

We will throughout use the following notations. For a vector $X = (X_i)_{i=1\dots N}$ and a series of relative weights S_i with $\sum_i S_i = 1$, we let

$$X_E := \frac{1}{N} \sum_{i=1}^N X_i, \quad (1)$$

$$X_S := \sum_{i=1}^N S_i X_i, \quad (2)$$

$$X_\Gamma := X_S - X_E, \quad (3)$$

so that X_E is the equal weighted average of the vector's elements, X_S is the size weighted average, and X_Γ is their difference.

We will also have shocks u_i that are uncorrelated and with variance $\sigma_{u_i}^2$. Then, we will define the “pseudo-equal weights” or “inverse variance weights”

$$\tilde{E}_i := \frac{1/\sigma_{u_i}^2}{\sum_j 1/\sigma_{u_j}^2}, \quad (4)$$

which satisfy $\sum_i \tilde{E}_i = 1$, and are equal to $\tilde{E}_i = \frac{1}{N}$ when all the $\sigma_{u_i}^2$ are equal. Then $X_{\tilde{E}} := \sum_{i=1}^N \tilde{E}_i X_i$. We'll also define

$$\tilde{\Gamma}_i := S_i - \tilde{E}_i. \quad (5)$$

Then, $X_{\tilde{\Gamma}}$ will be the “granular residual” in a number of settings. It is the size-weighted sample average of X minus the “inverse-variance” weighted sample average of X . It will be an optimal proxy for idiosyncratic shocks.

We use the notation β^e for the estimator of a parameter β ; \mathbb{E}_T for the sample temporal mean, $\mathbb{E}_T[Y_t] := \frac{1}{T} \sum_{t=1}^T Y_t$; C_t for a vector of controls; ι for a vector of 1's, I for the identity matrix, of the

appropriate dimension given the context; $V^Y = \mathbb{E}[Y_t Y_t']$ for a variance-covariance matrix; $X \perp Y$ to say that random variables X and Y are uncorrelated.

2.2 A very simple example with no feedback loop

We introduce GIVs by considering a very simple example with supply and demand.

2.2.1 Basic model

For clarity, we lay out a concrete economic model of the equilibrium in, for instance, the oil market. Demand by country i at date t is $D_{it} = \bar{Q} S_i (1 + y_{it})$, where \bar{Q} is the average total world production, y_{it} is a demand disturbance term, and S_i is country i 's share of demand, normalized to follow:

$$\sum_{i=1}^N S_i = 1. \quad (6)$$

The demand disturbance is assumed to be the sum of a common shock η_t , and an idiosyncratic shock u_{it} :

$$y_{it} = \lambda_i \eta_t + u_{it}. \quad (7)$$

For now we consider the case with uniform loadings,

$$\lambda_i = 1, \quad (8)$$

but we will relax that soon.

All shocks are i.i.d. across dates. Then, total world demand is $D_t = \sum_i D_{it} = \bar{Q} (1 + y_{St})$, where $y_{St} := \sum_i S_i y_{it}$ is the size-weighted average demand disturbance. We suppose that supply is $Q_t = \bar{Q} (1 + \frac{p_t - \bar{P}}{\alpha})$, where $p_t = \frac{P_t - \bar{P}}{\bar{P}}$ is the proportional deviation from \bar{P} , which is thus the average price of oil. Then, to equilibrate supply and demand ($D_t = Q_t(p_t)$), we must satisfy: $\bar{Q} (1 + y_{St}) = \bar{Q} (1 + \frac{p_t - \bar{P}}{\alpha})$. That is, the deviation of the price from the average satisfies:

$$p_t = \alpha y_{St} + \varepsilon_t. \quad (9)$$

It depends on the size-weighted average demand shock, $y_{St} = \sum_i S_i y_{it}$.

The classic problem is that we cannot estimate α by OLS. Indeed, a direct regression of p_t on y_{St} (that is, a regression of the form $p_t = \alpha y_{St} + \varepsilon_t$) would be biased, as ε_t and η_t (hence ε_t and y_{St}) can be correlated.

However, suppose that we form the GIV:

$$z_t := y_{\Gamma t} = y_{St} - y_{Et} = \sum_{i=1}^N S_i y_{it} - \sum_{i=1}^N \frac{1}{N} y_{it}. \quad (10)$$

Then, we have, using $u_{St} := \sum_{i=1}^N S_i u_{it}$, $u_{Et} := \sum_{i=1}^N \frac{1}{N} u_{it}$, that

$$y_{St} = \sum_{i=1}^N S_i y_{it} = \eta_t + u_{St}, \quad y_{Et} = \sum_{i=1}^N \frac{1}{N} y_{it} = \eta_t + u_{Et},$$

so $z_t := y_{St} - y_{Et} = (\eta_t + u_{St}) - (\eta_t + u_{Et})$ satisfies

$$z_t = u_{St} - u_{Et} =: u_{\Gamma t}. \quad (11)$$

Note that $z_t := y_{St} - y_{Et}$ is just constructed from observables. It is the difference between the size-weighted demand and the equal-weighted demand. Intuitively, it captures the “idiosyncratic demand” by large units, as shown by $z_t = u_{\Gamma t}$.

We assume that the shocks u_{it} are idiosyncratic, in the sense that:

$$\mathbb{E}[u_{it}\varepsilon_t] = 0 \text{ for all } i, t, \quad (12)$$

This “exogeneity” or “exclusion” assumption needs to be discussed in each economic application — as we will below. More minor, for simplicity, the u_{it} are here i.i.d. over time, but the u_{it}, u_{jt} could be correlated.

Then, we have

$$\mathbb{E}[z_t \varepsilon_t] = 0 : \text{Exogeneity}, \quad (13)$$

and

$$\mathbb{E}[z_t y_{St}] \neq 0 : \text{Relevance}. \quad (14)$$

Hence, $z_t = u_{\Gamma t}$ is a valid instrument (and as Proposition [3](#) will show, and optimal one). We call it a “granular instrumental variable” (GIVs).

Given that $p_t - \alpha y_{St} = \varepsilon_t$, we have

$$\mathbb{E}[(p_t - \alpha y_{St}) z_t] = 0, \quad (15)$$

that is, $\mathbb{E}[p_t z_t] - \alpha \mathbb{E}[y_{St} z_t] = 0$, which gives the supply elasticity α , by

$$\alpha = \frac{\mathbb{E}[p_t z_t]}{\mathbb{E}[y_{St} z_t]}. \quad (16)$$

Indeed, in practice, we might estimate α using sample means:

$$\alpha_T^e := \frac{\frac{1}{T} \sum_t p_t z_t}{\frac{1}{T} \sum_t y_{St} z_t}. \quad (17)$$

We now state a formal proposition³

Proposition 1 (Consistency of the GIV estimator in this example). *Suppose that $\mathbb{E}[u_{it}\varepsilon_t] = 0$ (but the u_{it} can have an arbitrary distribution, with mean 0), and we have a succession of i.i.d. dates t . Form the GIV estimator $z_t := y_{\Gamma t}$. Then, z_t identifies the price elasticity, by $\alpha = \frac{\mathbb{E}[p_t z_t]}{\mathbb{E}[y_{\Gamma t} z_t]}$. In other terms, the GIV estimator for the price elasticity α , $\alpha_T^e := \frac{\frac{1}{T} \sum_t p_t z_t}{\frac{1}{T} \sum_t y_{\Gamma t} z_t}$, is a consistent estimator.*

Precision of the GIV estimator We define the excess Herfindahl as:

$$h := \sqrt{-\frac{1}{N} + \sum_{i=1}^N S_i^2}. \quad (18)$$

In the context of industries, for example, a higher $h \in [0, \sqrt{1 - \frac{1}{N}}]$ means that the industry is more concentrated: an industry where all the firms have the same size features $h = 0$.

The quantity h will prove to be analytically useful, since if $(u_i)_{i=1\dots N}$ is a series of uncorrelated random variables with mean 0 and common variance σ_u^2 , then the volatility of the GIV $z_t = u_{\Gamma t}$ is:

$$\sigma_{u_{\Gamma}} = h\sigma_u. \quad (19)$$

The next proposition states the conditions under which we have a precise estimator (its proof is in Section 9).

Proposition 2 (Precision of the GIV estimator in this example). *The above estimator based on the granular instrument variable (GIV) $y_{\Gamma t}$ achieves identification of the elasticity parameter α , at the following rate, for $T \rightarrow \infty$:*

$$\sqrt{T}(\alpha_T^e - \alpha) \sim \mathcal{N}(0, \sigma_\alpha^2),$$

where $\sigma_\alpha = \frac{\sigma_\varepsilon}{\sigma_{u_{\Gamma}}}$. If we assume further than the u_{it} are i.i.d. with variance σ_u^2 ,

$$\sigma_\alpha = \frac{\sigma_\varepsilon}{h\sigma_u}, \quad (20)$$

where $h = \sqrt{-\frac{1}{N} + \sum_{i=1}^N S_i^2}$ is the excess Herfindahl.

So in order to have a precise estimate (low σ_α), we need: some large units (in order to have a large excess Herfindahl h), and that idiosyncratic shocks are large compared to aggregate shocks (large $\sigma_u/\sigma_\varepsilon$).

This simple example illustrates the basic idea. The reader might at this point have in mind a number of questions and objections: What if the factor structure is non-trivial (for instance, we

³It holds under mild regularity conditions on the joint distribution of $u_{it}, \eta_t, \varepsilon_t$ given that the data are i.i.d. across dates.

don't have $\lambda_i = 1$ in (7)? What if the demand is sensitive to price? Is the GIV that we constructed the best instrument we can find? What happens if there are more feedback loops?

The next subsections are devoted to answering them in turn.

2.2.2 Time-varying size weights

Suppose that we have a time-varying size $S_{i,t-1}$, so that the demand increase is $\sum_i S_{i,t-1} y_{it}$, with $S_{i,t-1} \perp (u_{it}, \eta_t, \varepsilon_t)$. Then everything goes through without problems, replacing S_i by $S_{i,t-1}$ throughout. The basic GIV becomes: $z_t = y_{S_{t-1},t} - y_{Et} = \sum_i (S_{i,t-1} - \frac{1}{N}) y_{it}$.

2.2.3 Model with an enriched factor structure

The model might have a richer factor structure, with r factors, i.e. instead of (7) we have:

$$y_{it} = \sum_{f=1}^r \lambda_i^f \eta_t^f + u_{it}, \quad (21)$$

or, in vector form:

$$y_t = \Lambda \eta_t + u_t, \quad (22)$$

where Λ is a $N \times r$ matrix, and $u_t \perp (\eta_t, \varepsilon_t)$

Then, in order to construct a valid GIV we simply run a factor model — for example, via Principal Component Analysis (PCA) — and, in essence, we extract the residuals u_{it} to form the GIV. Let us see that more precisely. Suppose that we have estimated the λ vector (e.g. via PCA, as we will detail later). Then, let Q be a $N \times N$ matrix projecting vectors onto a space orthogonal to Λ , so that $Q\Lambda = 0$.⁴ Then, $Qy_t = Qu_t$. Hence, via factor analysis, we obtain the transformed residuals $\check{u}_t = Qu_t$. Then, the GIV is formed as:

$$z_t := S' \check{u}_t = S' Q y_t = \Gamma' y_t, \quad \Gamma := Q' S, \quad (24)$$

so that

$$z_t = \Gamma' u_t. \quad (25)$$

Then, z_t is a valid instrument, since it is composed of idiosyncratic shocks. Since $p_t - \alpha y_{St} = \varepsilon_t$ and $\mathbb{E}[u_t \varepsilon_t] = 0$, we have $\mathbb{E}[(p_t - \alpha y_{St}) z_t] = 0$, i.e. (15).⁵

⁴For instance, we can take $Q = Q^{\Lambda, W}$, where

$$Q^{\Lambda, W} := I - \Lambda (\Lambda' W \Lambda)^{-1} \Lambda' W, \quad (23)$$

with $W = (V^u)^{-1}$ (optimally) or $W = I$ for simplicity. This choice satisfies $Q^{\Lambda, W} \Lambda = 0$.

⁵This shows that the GIV is valid and possible as long as $\Gamma := Q' S \neq 0$. Fortunately, this is generically true. If Γ were close to 0, that would be picked up by very large standard errors.

This generalizes our basic example (7). In that example, we had $Q = I - \iota E'$, so that $\check{u}_{it} = u_{it} - u_{Et}$, and the GIV was: $z_t = \check{u}_{St} = u_{St} - u_{Et}$. We therefore had $\Gamma = Q'S = S - E$.

2.3 A simple demand and supply example with feedback loops

2.3.1 A simple model

We next enrich the previous example, and consider a simple supply and demand example that features a “loop.” Suppose that demand for some commodity (say, oil) is:

$$y_{it} = \phi^d p_t + \eta_t + u_{it}, \quad (26)$$

and supply is

$$s_t = \phi^s p_t + \varepsilon_t, \quad (27)$$

where η_t, ε_t can be correlated. We can expect that the demand and supply elasticities (respectively ϕ^d and ϕ^s) satisfy $\phi^d < 0 < \phi^s$. Again, to be more formal, y_{it} , s_t , and p_t are understood as percent deviations from the average demand of country i , from supply, and from price, respectively.⁶

In equilibrium, supply equals demand, $y_{St} = s_t$, which gives the price

$$p_t = \frac{u_{St} + \eta_t - \varepsilon_t}{\phi^s - \phi^d}. \quad (28)$$

There is a “loop” because the demand shocks η_t and u_{it} feed into the price p_t , which then in turns affects demand. The equilibrium quantity produced is

$$s_t = y_{St} = \frac{\phi^s u_{St} + \phi^s \eta_t - \phi^d \varepsilon_t}{\phi^s - \phi^d}. \quad (29)$$

The classic problem of estimating supply and demand equilibrium quantity s_t and price p_t is that we cannot regress: $s_t = \beta p_t + \varepsilon_t$, and hope to get $\beta = \phi^s$, as ε_t and p_t are correlated.

However, suppose that we form the GIV, as in (10)

$$z_t := y_{St} - y_{Et}. \quad (30)$$

Given that

$$y_{St} = \phi^d p_t + \eta_t + u_{St}, \quad y_{Et} = \phi^d p_t + \eta_t + u_{Et},$$

we have:

$$z_t = u_{St} - u_{Et} =: u_{\Gamma t}. \quad (31)$$

⁶We take the model of Section 2.2, and simply set $y_{it} = \phi^d p_t + \eta_t + u_{it}$, where $p_t = \frac{P_t - P_*}{P_*}$ is the proportional deviation from the average.

As in the previous example, we assume that the shocks u_{it} are idiosyncratic:

$$\mathbb{E}[u_{it}\eta_t] = \mathbb{E}[u_{it}\varepsilon_t] = 0 \text{ for all } i, t. \quad (32)$$

Then, we have again a valid instrument:

$$\mathbb{E}[z_t\varepsilon_t] = \mathbb{E}[z_t\eta_t] = 0 : \text{ Exogeneity,}$$

$$\mathbb{E}[z_tp_t] \neq 0 : \text{ Relevance.}$$

Estimations of supply and demand elasticities by GIV The supply equation (27) implies

$$\mathbb{E}[(s_t - \phi^s p_t) z_t] = 0, \quad (33)$$

which gives the supply elasticity ϕ^s by

$$\phi^s = \frac{\mathbb{E}[s_t z_t]}{\mathbb{E}[p_t z_t]}. \quad (34)$$

Indeed, in practice, we form the sample average $\phi_T^{s,e} = \frac{\mathbb{E}_T[s_t z_t]}{\mathbb{E}_T[p_t z_t]}$.

Now, we want to estimate demand. For that, we make a stronger assumption: we assume that the shocks u_{it} are i.i.d. across i 's and not just dates (we will relax this later). Then, this implies⁷

$$\mathbb{E}[u_{Et} u_{\Gamma t}] = 0. \quad (35)$$

So, given this, we have: $y_{Et} - \phi^d p_t = \eta_t + u_{Et}$, and $\mathbb{E}[(y_{Et} - \phi^d p_t) z_t] = 0$. This gives an estimate of the demand elasticity ϕ^d ,

$$\phi^d = \frac{\mathbb{E}[y_{Et} z_t]}{\mathbb{E}[p_t z_t]}, \quad (36)$$

and the estimator is $\phi_T^{d,e} = \frac{\mathbb{E}_T[y_{Et} z_t]}{\mathbb{E}_T[p_t z_t]}$.

Estimation by OLS and interpreting it as a first- and second-stage IV estimator Let us recast our GIV in the language of applied microeconomics, and estimate the parameters by OLS

⁷Indeed, in the i.i.d. case we have

$$\mathbb{E}[u_{Et} u_{\Gamma t}] = \mathbb{E}\left[\left(\sum_i \frac{1}{N} u_{it}\right) \left(\sum_i \Gamma_i u_{it}\right)\right] = \frac{1}{N} \sum_i \Gamma_i \sigma_u^2 = 0.$$

as $\sum_i \Gamma_i = 0$. Equation (77) generalizes this to the non-i.i.d. case.

(as we will often do in the general case). Recall that the solutions are:

$$p_t = \frac{1}{\phi^s - \phi^d} u_{St} + \varepsilon_t^p, \quad s_t = y_{St} = \frac{\phi^s}{\phi^s - \phi^d} u_{St} + \varepsilon_t^s,$$

where the $\varepsilon_t^p, \varepsilon_t^s$ are linear combinations of ε_t, η_t . So, if we run the OLS regression, with $z_t = u_{St}$,

$$p_t = b^p z_t + \varepsilon_t^p, \quad (37)$$

we estimate

$$b^p = \frac{1}{\phi^s - \phi^d}, \quad (38)$$

which is the sensitivity of the price to the supply or demand shock. If we run the OLS regression:

$$y_{St} = b^{ys} z_t + \varepsilon_t^s, \quad (39)$$

we estimate

$$b^{ys} = \frac{\phi^s}{\phi^s - \phi^d} = M. \quad (40)$$

In the language of applied microeconomics, one can view the “first stage” as a regression of the price on the GIV (37). The “second stage” is running supply on the instrumented change in the price $b^p z_t$:

$$s_t = \phi^s (b^p z_t) + \varepsilon_t^s, \quad (41)$$

which gives ϕ^s . Alternatively, one can run the “reduced form equation” (39), which estimates estimating M . The supply elasticity is given by:

$$\phi^s = \frac{b^{ys}}{b^p}, \quad (42)$$

The demand elasticity is similar. In the second stage we run equal-weighted demand on the instrumented change in the price, $b^p z_t$:

$$y_{Et} = \phi^d (b^p z_t) + \varepsilon_t^y, \quad (43)$$

which gives the demand elasticity ϕ^d [8]. Alternatively, we can run the reduced form equation $y_{Et} = b^{yE} z_t + \varepsilon_t^y$ which gives $b^{yE} = \phi^d M$, and the demand elasticity is $\phi^d = \frac{b^{yE}}{b^p}$.

In practice, we will add controls to those regressions, including estimates of η_t recovered from factor analysis.

From M and b^p , we can recover the elasticities ϕ^s and ϕ^d . This is exactly the same estimate as the IV estimator, derived earlier in (34), (36).⁹

⁸Here we used (35), which makes the OLS valid.

⁹Indeed, the OLS estimators are $M_T^e = \frac{\mathbb{E}_T[y_{St} z_t]}{\mathbb{E}_T[z_t^2]}$ and $\psi_T^e = \frac{\mathbb{E}_T[p_t z_t]}{\mathbb{E}_T[z_t^2]}$. We have $\phi_T^{s,e} = \frac{M_T^e}{\psi_T^e} = \frac{\mathbb{E}_T[y_{St} z_t]}{\mathbb{E}_T[p_t z_t]}$, which is the

Standard errors: When “weak instruments” are or are not a problem When estimating via OLS (e.g. b^p and M), the standard errors are reliably estimated by the usual OLS method, even in small samples. When a ratio is implicitly performed (e.g. to estimate ϕ^d , ϕ^s), the 2SLS procedure as in (41) will also give correct standard errors when the instrument is strong enough. A good rule of thumb for the strength of the instrument is that the F statistics (which is the squared t-stats on b^p) on the first stage (37) should be greater than 10.¹⁰

2.3.2 A richer model with factor structure

We can also have a richer factor structure, as in Section 2.2.3. The demand could be:

$$y_{it} = \phi^d p_t + \sum_f \lambda_i^p \eta_t^f + u_{it}. \quad (44)$$

This is easy to accommodate. Using the notation $\check{a}_{it} := a_{it} - a_{Et}$, we have: $\check{y}_{it} = \sum_f \check{\lambda}_i^f \eta_t^f + \check{u}_{it}$, which is a plain factor model. Hence, we apply the procedure of Section 2.2.3. For instance, we get the vector of recovered idiosyncratic shocks $v_{it} := Q\check{u}_{it}$ (which again is $v_{it} = u_{it} - u_{Et}$ in our basic model), and set the GIV as $z_t = v_{St}$. We proceed as above (Section 2.3.1).¹¹

2.3.3 Misspecified models

The reader might have other concerns, e.g. what happens if there’s more heterogeneity, e.g. in demand elasticities? we if we don’t control for all common factors? We defer their discussion to Section 6.2, after we know the general model and procedure.

2.4 Optimality of the GIV

We come back to the simplest case of Section 2.2.1, for ease of exposition.¹² Above, we have shown that $z_t = y_{\Gamma t}$ allows for identification, for a specific $\Gamma = S - E$. It is easily verified that GIV with weights Γ such that $\sum_i \Gamma_i = 0$ would work. Hence, we can ask for an optimal Γ . The Γ we proposed initially was actually optimal, as we formalize below.

Proposition 3 (Optimal weights Γ for the GIV $y_{\Gamma t}$). *Consider the GIV $z_t = y_{\Gamma t} = \sum_i \Gamma_i y_{it}$, with some weights Γ_i with $\sum_i \Gamma_i = 0$. The idiosyncratic shocks u_{it} ’s are assumed to be i.i.d. across time, and have variance-covariance matrix V^u . Then, in the basic supply and demand model of section*

same as (34), as $s_t = y_{St}$ in equilibrium.

¹⁰See the literature on weak instruments, e.g. Stock and Yogo (2005); Andrews et al. (forthcoming).

¹¹More generally, if we know that λ^p lives in a given vector space Λ^p (i.e., $\lambda^p = \Lambda^p K$ for some K and specified Λ^p), we use $\check{y}_t = Q^{\Lambda^p} y_t$, where Q^{Λ^p} is some matrix such that $Q^{\Lambda^p} \Lambda^p = 0$, as in (23). We run a factor model on the \check{y}_t , get the residual v_t , and use the GIV $z_t = \Gamma' v_t$.

¹²But it will be clear to the reader that the result in the present subsection hold much more generally.

[2.2](#), the asymptotic variance of the estimator α_T^e in [17](#) (which is $\sigma_\alpha^2 = \lim_{T \rightarrow \infty} T \text{var}(\alpha_T^e - \alpha)$) satisfies $\sigma_\alpha^2(\Gamma) = \frac{\sigma_\varepsilon^2 \mathbb{E}[y_{\Gamma t}^2]}{\mathbb{E}[y_{St} y_{\Gamma t}]^2}$. The value

$$\tilde{\Gamma} = S - \tilde{E}, \quad \tilde{E} := \frac{(V^u)^{-1} \iota}{\iota' (V^u)^{-1} \iota}, \quad (45)$$

gives the optimal GIV estimator, in the sense that for any other Γ that is not collinear to $\tilde{\Gamma}$, the asymptotic variance $\sigma_\alpha^2(\Gamma)$ is larger. When the shocks are i.i.d., this implies $\tilde{E}_i = \frac{1}{N}$, and when they are uncorrelated, this implies $\tilde{E}_i := \frac{1/\sigma_{u_i}^2}{\sum_j 1/\sigma_{u_j}^2}$, so that \tilde{E} may be called the “precision-weighted quasi-equal” weights.

Hence, the “essence” of the GIV is not to be “size weighted minus value weighted” idiosyncratic shocks, but rather “size weighted minus precision weighted” (i.e. inverse-variance weighted when shocks are uncorrelated) idiosyncratic shocks.

There are two more ways in which the GIV is optimal. First, it is the optimally-weighted GMM estimator.[13](#) This implies that other combinations of idiosyncratic shocks (besides weighing by Γ) would not help the precision of the estimator. Second, one can show that it is the maximum likelihood estimator, if we assume that all shocks are Gaussian (see Section [14](#)). Still, the optimality formulation of Proposition [3](#) is the simplest to use in other contexts.

2.5 Interpreting and diagnosing idiosyncratic shocks

What is an idiosyncratic shock? Mathematically, an idiosyncratic shock is plainly a random variable u_{it} such that $\mathbb{E}_{t-1}[(\eta_t, \varepsilon_t) u_{it}] = 0$. But it may be useful to discuss different types of economic settings that map into that definition.

In some cases it is quite clear – e.g. a random productivity shock, or demand shock. But there are more subtle types of idiosyncratic shocks. One is an “unexpected change in the loading on a common shock”. For instance, suppose that OPEC decided to cut down production, which in the language of our example is an aggregate η_t shock. However, if Saudi Arabia cuts down production by more than anticipated, that is an idiosyncratic shock. Formally, if supply is $y_{it} = \phi^s p_t + (\lambda + \check{\lambda}_{it}) \eta_t + v_{it}$, with $\mathbb{E}_{t-1}[(1, \eta_t^2) \check{\lambda}_{it}] = 0$, then $u_{it} = \check{\lambda}_{it} \eta_t + v_{it}$ is an bona fide idiosyncratic shock. To take another example, suppose that we hear about a change in real estate prices in the economy, η_t , but that a bank i was more exposed to it than anticipated: the market thought the bank’s equity would move by $\lambda_i \eta_t$, but it moved by $r_{it} = (\lambda_i + \check{\lambda}_{it}) \eta_t$ for an expectational surprise $\check{\lambda}_{it}$ with $\mathbb{E}_{t-1}[(1, \eta_t^2) \check{\lambda}_{it}] = 0$. Then, the bank will have an idiosyncratic shock $u_{it} = \check{\lambda}_{it} \eta_t$ as part of its total return r_{it} .

Likewise suppose that the news is that a bank failed a stress test (while it was anticipated it would pass the test). This is an idiosyncratic shock. However, the bank could have failed the test

¹³Any moment $\mathbb{E}_T[(p_t - \alpha y_{St})(u_{it} - u_{Et})] = 0$ is a valid GMM moment to identify α . It is easy to check that the optimal GMM weighted estimator is our GIV, $\mathbb{E}_T[(p_t - \alpha y_{St})(u_{St} - u_{Et})] = 0$

because of some development in the macroeconomy η_t . Then, provided that the factor model allows for a rich enough structure in η_t , the latter will be controlled for.

Thresholded and narrative GIVs In applications, it is possible to make further progress by assessing the drivers of the top shocks narratively. One procedure is to simply select the top K shocks, by $S_i |\check{u}_{it}|$ (where \check{u}_{it} is the residual from factor analysis, e.g. $\check{u}_{it} = u_{it} - u_{Et}$), and check in the news what happened on that day (and check that the shocks are idiosyncratic indeed). We did that for some of our applications. Formally, that means that we formulate a “thresholded” GIV,

$$z_t^\tau = \sum_i \tau(S_i \check{u}_{it}), \quad (46)$$

using the thresholding function $\tau(x) = x1_{|x| \geq b}$, which only keeps granular shocks bigger than $b > 0$.¹⁴ Then, the GIV procedure works using that “thresholded” GIV (see Section 12.8). This thresholded GIV might also be useful to assess non-linear effects, for instance, in case of demand or supply curves.

After examining those largest shocks by looking at the news, some shocks might be eliminated as not idiosyncratic; we can call $I_t^\mathcal{N}$ the set of shocks that are “narratively certified” to be idiosyncratic by this procedure, and form the alternative instrument

$$z_t^\mathcal{N} = \sum_{i \in I_t^\mathcal{N}} S_i \check{u}_{it}. \quad (47)$$

This is roughly what the “narrative” approach in the literature (e.g. Caldara et al. (2018)) does. Here, in addition, we have a systematic way to select the candidate large shocks (by top values of $S_i |\check{u}_{it}|$), and get the controls η_t^e for the idiosyncratic shocks, when we run the regressions.

Quasi-experimental instruments and identification by functional form A large literature explores identification by functional form, where consistency of the estimator depends on functional form or distributional assumptions. Classic examples include the Heckman (1978) selection model, identification via heteroskedasticity, see Rigobon (2003) and Lewbel (2012), and Arellano and Bond (1991) and Blundell and Bond (1998) in the context of dynamic panel data models. The typical concern with these approaches, compared to quasi-experimental instruments that are outside of the model, is that the estimators are inconsistent when the model is misspecified.

In case of GIVs, we generally start from a structural model that motivates the estimating equation, as for instance in the model of doom loops in section 4, which prescribes the definition of the size vector, S , and, in some cases, the characteristics that determine the exposures, x_{it} . To extract idiosyncratic shocks, we rely on statistical factor models.¹⁵

¹⁴We adjust b to select a pre-specified number K of shocks that survive the thresholding.

¹⁵We discuss the robustness of GIVs to various forms of misspecification in Section 6.2.

Instead of viewing this last step as merely a statistical exercise that is hard to validate externally, GIVs provide an empirical strategy to understand the economic drivers of the instrument by screening the top shocks narratively. By understanding the nature of the shock based on news coverage (as in the narrative examination we just discussed), for instance, we can ensure that the shocks are truly idiosyncratic and interpretable. For instance, a large negative return associated with a failed stress test of a bank in the context of doom loops or a negative supply shock in Kuwait and Iraq during the first gulf war and positive demand shocks in China in the early 2000s in the context crude oil markets, are valid instruments. While alternative identification methods might rely on functional form assumptions only, GIVs, by being able to screen the shocks economically, provide a systematic way to construct instruments more in the spirit of quasi-experimental instruments.

3 General setup and procedure

The previous section introduced the GIV in a simple context, with no loops or a single loop. We now propose a more general setup with potentially several factors, arbitrary loop structure, and rich heterogeneity.

3.1 Framework

Consider the following model of stationary “actions” y_{it} (e.g. employment, investment, TFP shock, return, etc.) by “actor” i (e.g., a firm i in a closed-economy setting, or a country i in an international setting):

$$y_{it} = \sum_f \lambda_{it}^f F_t^f + u_{it} + k_{y_i} + m C_{it}^y, \quad (48)$$

where each F_t^f is a factor, λ_{it}^f are factor loadings, u_{it} is an idiosyncratic shocks, k_{y_i} is a constant, C_{it}^y is a vector of controls that may include lagged demands and other characteristics. Factor f follows:

$$F_t^f = \alpha^f y_{St} + \eta_t^f + k^f + m^f C_t^F. \quad (49)$$

It depends on an exogenous shock η_t^f , and some on the mean action y_{St} , and on a set of controls C_t^F (different from C_{it}^y). Those controls may include, for instance, lagged values. We assume that the “size” weights have been normalized to add to one, $\sum_i S_i = 1$. We partition the factors into “exogenous factors” \mathcal{F}^{Exo} , where we know $\alpha^f = 0$, and “endogenous” factors $\mathcal{F}^{\text{Endo}}$, where α^f may be non-zero.

We model the exposures to factors as either non-parametric (unrestricted λ_i^f) or as parametric:

$$\lambda_{it}^f = \lambda_0^f + \lambda_1^f x_{it}^f, \quad (50)$$

where $x_{it}^f \in \mathbb{R}^x$ is an observable (e.g. x_{it} is the de-meaned log size of entity i). For endogenous

factors, we treat here the parametric case, and defer the non-parametric case to Section 10.1. We can also treat the semi-parametric case where

$$\lambda_{it}^f = \lambda_0^f + \lambda_1^f x_{it}^f + \zeta_i^f, \quad (51)$$

where ζ_i^f is an extra non-parametric case.

We make the following identifying assumptions, for all f, i , the noise u_{it} are idiosyncratic:

$$u_{it} \perp (\eta_t^f, C_t^y, C_t^F, x_{it}^f), \quad (52)$$

but the η_t^f may be correlated across f 's, and η_t^f may be correlated with the controls, C_t^y and C_t^F . The u_{it} may have some correlation across i 's and can be heteroskedastic, as we discuss later. For expositional simplicity we assume that all dates are i.i.d.

We rewrite model (48) in vector form:

$$y_t = \Lambda_t F_t + u_t + C_t^y m + k_y, \quad (53)$$

with $F_t = (F_t^{\text{Endo}}, F_t^{\text{Exo}})'$, $\Lambda_t = (\Lambda_t^{\text{Endo}}, \Lambda_t^{\text{Exo}})'$, Λ_t a $N \times r$ matrix, F_t a $r \times 1$ vector, C_t^y an $N \times c$ matrix, m is $c \times 1$, where c is the dimension of the controls, and k_y is a $N \times 1$ vector. The endogenous and exogenous factors, and their loadings, are defined by:¹⁶

$$F_t^j = \alpha^j y_{St} + \eta_t^j + k_F^j + C_t^F \phi^j,$$

with $j = \text{Exo}, \text{Endo}$, with $\alpha^{\text{Exo}} = 0$. We write $\alpha = (\alpha^j)_{j=\text{Exo}, \text{Endo}} = (\alpha^f)_{f=1 \dots r}$ is an $r \times 1$ vector.

Multipliers Solving for the model gives, $y_{St} = \Lambda_{St} F_t + u_{St} + k_{yS} + C_{St}^y m$, that is,

$$y_{St} = \Lambda_{St} \alpha y_{St} + u_{St} + b_t, \quad (54)$$

where b_t satisfies $b_t \perp u_{St}$.¹⁷ So, we can solve for the aggregate outcome y_{St} as $y_{St} = \frac{u_{St} + b_t}{1 - \Lambda_{St} \alpha}$, that is,

$$y_{St} = M_t (u_{St} + b_t), \quad (55)$$

where the multiplier $M_t = \frac{dy_{St}}{du_{St}}$ is

$$M_t = \frac{1}{1 - \Lambda_{St} \alpha} = \frac{1}{1 - \sum_{f:\text{Endo}} \Lambda_{St}^f \alpha^f} \quad (56)$$

¹⁶Our initial examples are particular cases of the general procedure, as detailed in Section 12.10.

¹⁷We have $b_t = k_{yS} + C_{St}^y m + \Lambda_{St}^{\text{Exo}} F_t^{\text{Exo}} + \Lambda_{St}^{\text{Endo}} (\eta_t^{\text{Endo}} + k_F^{\text{Endo}} + C_t^F \phi^{\text{Endo}})$, with $k_{yS} := \sum_i S_i k_{y_i}$.

measures the total impact of shocks, after going through all feedback loops. Hence, an idiosyncratic shock has an impact that M_t times bigger than its direct effect. Also, the total impact of an idiosyncratic shock on factor f is:

$$F_t^f = M_t \alpha^f u_{St} + b_t^f \quad (57)$$

for another expression $b_t^f \perp u_{St}$. Our regressions will allow to identify M and $M\alpha^f$.¹⁸

In some cases, we may not observe all endogenous factors, F_t^{Endo} . In this case, we still recover the correct multiplier, M_t , and it should be interpreted as accounting for all feedback loops in the economy, including those operating via the unobservable, endogenous factors. However, we can obviously not estimate α^f for those unobserved factors.

3.2 A step-by-step user's guide

We outline the benchmark procedure, alongside several extensions. We summarize the model as

$$\begin{aligned} y_t &= \Lambda_t F_t + u_t + k_y + C_t^y m, \\ F_t^j &= \alpha^j y_{St} + \eta_t^j + k_F^j + C_t^F \phi^j, \end{aligned}$$

where we focus on the case in which $\Lambda_t^{\text{Endo}} = \lambda_0^{\text{Endo}} \iota$, where λ_0^{Endo} is a scalar.¹⁹ Loadings can be semi-parametric, $\Lambda_t^{\text{Exo}} = \Lambda_0^{\text{Exo}} + \Lambda_1^{\text{Exo}} x_t + \zeta_t$, where $\mathbb{E}[\zeta_t \check{x}_t] = 0$, or non-parametric.

We assume that y_t is observed, and we run regressions on the observed factors F_t^f .

1. Define $\check{a}_{it} = a_{it} - a_{Et}$ for a generic variable a_t , and estimate the panel regression

$$\check{y}_t = c + \check{\Lambda}_t^{\text{Exo}} C_t^{F, \text{Exo}} \phi + \check{C}_t^y m + e_t,$$

which removes endogenous factors and estimates the coefficients on the controls. The vector of residuals equals $e_t = \check{\Lambda}_t^{\text{Exo}} \eta_t^{\text{Exo}} + \check{u}_t$.

2. We can estimate the factors, η_t^{Exo} , and denote the estimates by η_t^e , in one of two ways:

(a) In the case of semi-parametric loadings, $\Lambda_t^{\text{Exo}} = \Lambda_0^{\text{Exo}} + \Lambda_1^{\text{Exo}} x_t + \zeta_t$, we estimate the

¹⁸The more advanced impact of a shock to the variable $j \neq i$:

$$\frac{dy_{it}}{du_{St, -i}} = M \sum_f \lambda_i^f \alpha^f, \quad (58)$$

which can also be estimated, see (126). The above equation is the reaction dy_{it} of to a shock du_{St} , when we also impose that $du_{it} = 0$.

¹⁹The procedure in this section extends to the model with multiple endogenous factors, when the factor exposures, λ_i , depend on a vector of characteristics, and to the case where the characteristics may vary over time, x_{it} . Then, the generalized version of $\check{a}_{it} = a_{it} - a_{Et}$ in Step 1 becomes $\check{a}_t = Q^{X, W} a_t$ where $Q^{X, W}$ is defined in (23) (if $X = \iota$, then we recover $\check{a}_{it} = a_{it} - a_{Et}$ for $W = I$, and $\check{a}_{it} = a_{it} - a_{\bar{E}t}$ for $W = (V^u)^{-1}$).

factors using

$$e_t = \tilde{x}_t \eta_t^{\text{Exo}} + \epsilon_t,$$

which is equivalent to a series of cross-sectional regressions of e_t on x_t to estimate η_t^{Exo} . By considering semi-parametric loadings, we do not need to know exactly how the loadings depend on characteristics, x_t , and a noisy signal of exposures suffices to estimate the factors, η_t^{Exo} .²⁰

- (b) In the case of non-parametric loadings, we estimate the factors using factor analysis, in case of small N , or principal components analysis (PCA), in case of large N , from e_t . We estimate the number of factors using the methods developed in Bai and Ng (2002).

3. Estimate $(M, \alpha^f M)$, using OLS, with $Z_t = y_{\Gamma t}$,

$$\begin{aligned} y_{St} &= MZ_t + \gamma'_y \eta_t^e + k_y + C_{St}^y m + e_t^y, \\ F_t^{\text{Endo}} &= \alpha^{\text{Endo}} MZ_t + \gamma'_F \eta_t^e + k_F^{\text{Endo}} + C_t^F \phi^{\text{Endo}} + e_t^{\text{Endo}}. \end{aligned}$$

4. Estimate $(\alpha^{\text{Endo}})^e := \frac{(\alpha^{\text{Endo}} M)^e}{M^e}$. Similarly, we can recover from $\Lambda_S^{\text{Endo}} \alpha^{\text{Endo}}$ from $M = \frac{1}{1 - \Lambda_S^{\text{Endo}} \alpha^{\text{Endo}}}$. In both cases, M^e and $(\alpha^{\text{Endo}} M)^e$ need to be sufficiently precisely estimated. This is analogous to the weak instruments problem in the context of IV estimators.

When the idiosyncratic shocks are heteroskedastic, we may be able to improve the finite-sample properties. We start from $E_i = \frac{1}{N}$. After Step 2, we use the residuals to obtain an estimate of $\sigma_{u_i}^2$ and update the weights to $\tilde{E}_i = \sigma_{u_i}^{-2} / \sum_{i=1}^N \sigma_{u_i}^{-2}$. In all cases, our estimators are consistent for large N and T . However, to obtain consistency with finite N and large T , we need the precision-weighted \tilde{E} .²¹ This implies a standard bias-efficiency trade-off if estimating volatilities reduces the efficiency of the estimator. If the idiosyncratic volatilities are related to size, we can parameterize them and estimate

$$\ln \sigma_{u_i}^2 = \sigma_0 + \sigma_1 \ln S_i + \epsilon_i,$$

and use $\sigma_{u_i}^2 = \exp(\sigma_0) S_i^{\sigma_1}$.

First and second stages Let us see how to estimate α^f in that first and second stage language. The *first stage* is the regression:

$$y_{St} = bz_t + \beta_{\eta}^{ys} \eta_t^e + \beta_C^{ys} C_{St} + \varepsilon_t^{ys}, \quad (59)$$

²⁰We refer to Fan et al. (2016) for a related approach in the context of principal components analysis.

²¹Here we are talking about consistency in the estimation of M and $\alpha^f M$. It is achieved even if we do not consistently estimate the underlying factors η_t . This may be surprising, but this is already the case in the simple supply and demand case of Section 2.

where we regress on z_t , using our recovered factors η_t as controls. From the model, we know that the regression coefficient identifies $b = M$, the multiplier.

The *second stage* is the regression:

$$F_t^f = \alpha^f (b^e z_t) + \beta^{F^f} \eta_t^e + \varepsilon_t^{F^f}, \quad (60)$$

which gives the estimator for α^f . Alternatively, we can regress F_t^f on z_t (with controls) and the coefficient is $\alpha^f M$. When estimating the influence coefficient, we can also view the second stage as estimating

$$y_{Et} = \gamma M^e z_t + \beta_{\eta}^{ys} \eta_t^e + \beta_C^{ys} C_{St} + \varepsilon_t^{ys},$$

which gives γ .

3.3 A formal identifiability result

We encourage the reader to skip this section at the first reading. We provide here formal conditions for identification, completing the simpler case of Section 2.

We study the “semi-parametric” case. We have some characteristics x_{it} of actors (e.g. countries or and firms): for instance, depending on the application we know that the loading is an affine function of log market capitalization, or the stock market beta of a bank, or OPEC membership. We also have a priori knowledge that $\lambda_{it}^f = \lambda_0^f + \lambda_1^f x_{it}^f$, for some parameters λ_1^f in the parametric case (50), and something similar in the non-parametric case (51). This is consistent with the practice in modern finance in which risk exposures (betas) align with characteristics (see e.g. Fama and French (1993)), so that parametric approaches are preferred, in particular because they are more stable than non-parametric approaches. Section 10.1 develops the full non-parametric version, estimating the factors. We don’t have a priori information about the η_t , nor their covariance V^η .

To obtain identification, we shall make the following two assumptions.

Assumption 1 (Condition for identification with GIV) *The vector $V^u S$ is not spanned by the factors loadings λ^f (where V^u is the covariance matrix of u_t).*

Assumption 2 (Restriction on the admissible variance-covariance matrix of residual u_t) *The variance-covariance on u_t is diagonal.*

Assumption 1 is essential and could not be relaxed. It ensures that the GIV is not identically 0. Economically, this assumption seems like a mild restriction. It is generically satisfied.²²

Assumption 2 could be relaxed in number of ways.²³ Other sufficient condition for identification

²²One case that does prevent this assumption to hold is the case where the variance would be inversely proportional to size: then, GIV would fail, as then $V^u S = aI$ for some scalar a . Fortunately, in most contexts, variance may decay a bit with size S_i , but less violently than in $1/S_i$ (see e.g. Lee et al. (1998) and the discussion in Gabaix (2011)).

²³However, relaxations of Assumption 2 will still need to ensure some restrictions on the space of variance-covariances allowed.

might be that V^u is k -sparse, e.g. has at most k non-zero off-diagonal elements, for some k , e.g. $N - r^2$ (see also Zou et al. (2006)). Another is to allow for some correlation that depends on the distance between entities i and j , perhaps via Gaussian processes (Rasmussen and Williams (2005)). We conjecture that this proposition could be generalized in a number of ways, including in the large T, N domain, using material such as Bai and Ng (2006). Doing this would however take us too far afield.

We assume that all shocks are i.i.d. over time, though this would be easy to relax.

We next state a formal identification result, which is proven in Section 9.

Proposition 4 (Sufficient condition for identification with GIV) *Consider the factor model above, when N is fixed but $T \rightarrow \infty$, and makes Assumptions 1 and 2. Then, the parametric (and semi-parametric) procedure of Section 3.2 identifies M , α^f by GIV. Furthermore, the standard errors on M and $\alpha^f M$ returned by OLS in this procedure are valid.*

This completes our “abstract” development of GIV. We now turn to two initial applications.

4 Doom loops: Spillovers from banks to sovereign yields and vice-versa

4.1 Financial doom loops

As a first application of GIVs, we estimate a model of “financial doom loops,” which refers to the mechanism that a shock to the banking sector impacts the state’s balance sheet (e.g. as the state is more likely to have to bail out the banking sector), hence the sovereign yield; and that, in turn, weakening of the sovereign may impact the banking sector, creating a vicious circle (Farhi and Tirole (2017)). Doom loops have received significant attention during the European sovereign debt crisis. Although correlations between banks’ equity returns and sovereign yields turned strongly negative during this episode in Greece, Italy, and Spain, there are no causal estimates and it cannot be ruled out that common shocks, for instance related to European Central Bank’s (ECB) policies, are responsible for the comovement. We use GIVs to estimate the full doom loop for Italy. Section 11 provides a model microfoundation and we use the structure of that model to guide our empirical analysis. The results in this section are preliminary and future versions of this paper will extend the methodology in this section to other European countries.

As the model has predictions for credit spreads, we use the spread between Italian and German yields in the analysis below, where we interpret the German yield as the risk-free yield.²⁴

²⁴When we use the Italian CDS spread instead of the spread between Italian and German yields, we obtain very similar results.

We model the return of bank i as

$$r_{it} = \kappa_i + (\lambda_0^F + \lambda_1^F x_{i,t-1}) F_t + (\lambda_0^y + \lambda_1^y x_{i,t-1}) \Delta y_t + u_{it}, \quad (61)$$

where $x_{i,t-1}$ measures the heterogeneous exposure to common shocks and yields. F_t captures common shocks, such as broad movements to the stock market, industry-wide shocks to the banking sector or ECB policies. Δy_t is the change in the sovereign yield spread and $\lambda_0^y + \lambda_1^y x_{i,t-1}$ captures the exposure of bank i to sovereign yield spread changes.

The sovereign yield spread is in turn impacted by the health of the banking sector

$$\Delta y_t = k_y + \gamma_{t-1} r_{St} + \eta_t^y. \quad (62)$$

In the model of section [11](#), we find that the multiplier is $\gamma_{t-1} = c p_{G,t-1}$, where c is a constant, and $p_{G,t-1}$ is the (market-based) probability of default of the sovereign.^{[25](#)}

If we solve the model, then we find for the value-weighted banking return, with $\lambda_{S,t-1}^f = \lambda_0^f + \lambda_1^f x_{S,t-1}$,

$$r_{St} = M_{t-1} u_{St} + M_{t-1} \lambda_{S,t-1}^F F_t + M_{t-1} \lambda_{S,t-1}^y \eta_t^y + k^r,$$

where M_t is the “multiplier” for the doom loop, $M_{t-1} = \frac{1}{1 - \lambda_{S,t-1}^y \gamma_{t-1}}$. For changes in sovereign yield spreads, we have

$$\Delta y_t = \gamma_{t-1} M_{t-1} u_{St} + \xi_t + k^y.$$

4.2 Interpreting the coefficients and multipliers

The coefficient $\lambda_0^y + \lambda_1^y x_{i,t-1}$ in [\(61\)](#) measures the overall impact of changes in sovereign spreads on stock returns. The spread reflects the fiscal position of the sovereign, which may impact banks via its ability to (i) bail out the financial sector and to (ii) provide economic stimulus more broadly that may impact the Italian economy. It may also be the case that banks hold significant positions in sovereign debt and capital losses due to increasing yields negatively impact banks; see also [Altavilla et al. \(2017\)](#).

The banking sector may impact the sovereign via different channels. A failure of a large bank may directly impact the sovereign if it needs to be bailed out. Second, a weakening of the banking sector may also lead to a contraction in credit supply and slow growth. Third, to the extent that banks are important holders of sovereign debt, a weakened banking sector may lower the demand for sovereign debt. The coefficient γ in [\(62\)](#) measures the overall impact.

4.3 GIV estimation

We estimate the model using the following steps

²⁵That is, the probability under the risk-neutral probability \mathbb{Q} , not the physical probability \mathbb{P} .

1. Run a panel regression of returns with bank and time fixed effects

$$r_{it} = \kappa_i + a_t + \check{e}_{it},$$

where we weigh the observations using the inverse of the variance of returns, r_{it} , estimated using the three most recent months to adjust for heteroskedasticity.

2. Use \check{e}_{it} to estimate η_t^x (as in Section 3.2, Step 2) using $x_{i,t-1}$, the (cross-sectionally) demeaned log market capitalization of the bank. We define the vector of controls as $C_t = (\eta_t^x, p_{Gt}\eta_t^x)$, that is, the extracted factor and its interaction with Italy's CDS spread. Recall that p_{Gt} is the (market-based) probability of default of the sovereign.
3. Estimate the multiplier M_t (with $Z_t := r_{rt}$ which is our z_t plus some linear function of η_t^x , which is anyway controlled for in the regression) using the OLS regression

$$r_{St} = M_t Z_t + \beta^{r'} C_t + k^r + e_t^r, \quad (63)$$

where we model M_t as a function of Italy's CDS spread for parsimony, $M_t = M_0 + M_1 p_{Gt}$.

4. Estimate $\beta_t = \gamma_t M_t$ using the OLS regression

$$\Delta y_t = \beta_t z_t + \beta^y C_t + k^y + e_t^y, \quad (64)$$

where $\beta_t = \gamma_t p_{Gt}$. Our micro-founded model suggests a linear, instead of an affine, specification of $\gamma_t M_t$ in p_{Gt} . After all, if there is no risk that the sovereign defaults, the return on the banking sector should not impact the spread either.

4.4 Data

We estimate the model using data from the Italian banking sector and government bond market. The sample is from from September 2009 to June 2018. For the equity data, we construct a list of banks using ORBIS from Bureau van Dijk. We include banks that have been listed (but may delist during the sample). We merge the banks using ISIN codes to the research lists for equities from Thomson Reuters' Datastream²⁶ and only keep the primary security. We include all banks with assets above EUR2.5bn.²⁷

We use sovereign 10-year benchmark yields from Thomson Reuters and CDS data from IHS Markit. We compute arithmetic returns for stocks and simple differences in yields as our measure

²⁶For Italy, the codes for these lists are FITA and DEADIT, and include both alive and dead securities.

²⁷We also remove a small fraction of observations with daily returns over 100% or for which the total return index falls below one to avoid problems with discreteness.

Table 1: List of Italian banks. The table lists the banks in our sample from September 2009 to June 2018 in the first column and the average assets of the banks in the second column.

Bank name	Mean assets (EUR bn)
UniCredit SpA	863.7
Intesa Sanpaolo	667.5
Banca Monte dei Paschi di Siena SpA	201.3
Banco BPM SpA	130.3
Banche Popolari Unite	122.2
Mediobanca Banca Di Credito Finanz	72.3
BPER Banca S.P.A.	62.0
Banca Popolare Di Milano	50.0
Banca Mediolanum SpA	39.0
Banca Carige SpA	35.9
Credito Emiliano SpA	34.6
Banca Popolare di Sondrio	33.6
Banca Piccolo Credito Valtellinese	26.7
Credito Bergamasco	14.5
Banco Di Sardegna	12.8
Banco di Desio e della Brianza SpA	10.7
Credito Artigiano	9.1
Banca Generali SpA	6.6
Banca Popolare di Spoleto SpA	3.3
IW Bank SpA	2.8

Figure 1: Yield dynamics in Germany and Italy. The figure plots the 10-year government bond yields from January 2000 until July 2018 in Germany and Italy.

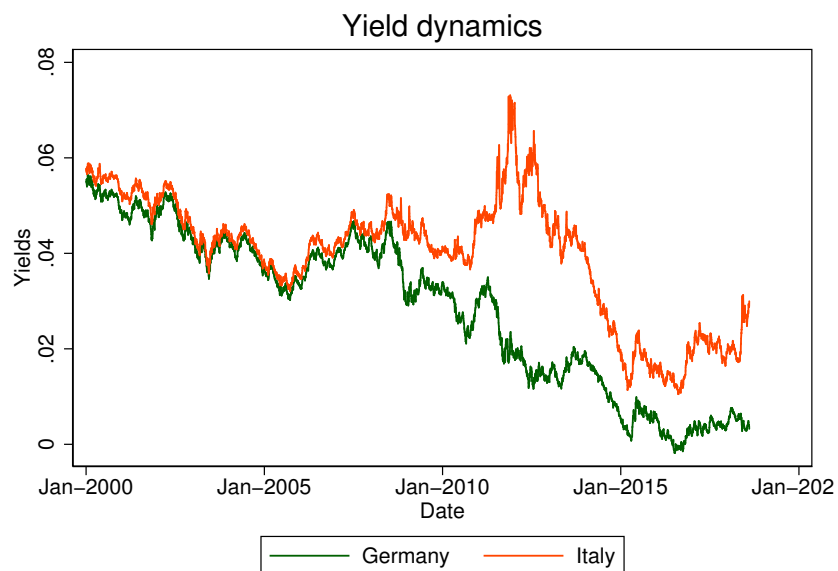
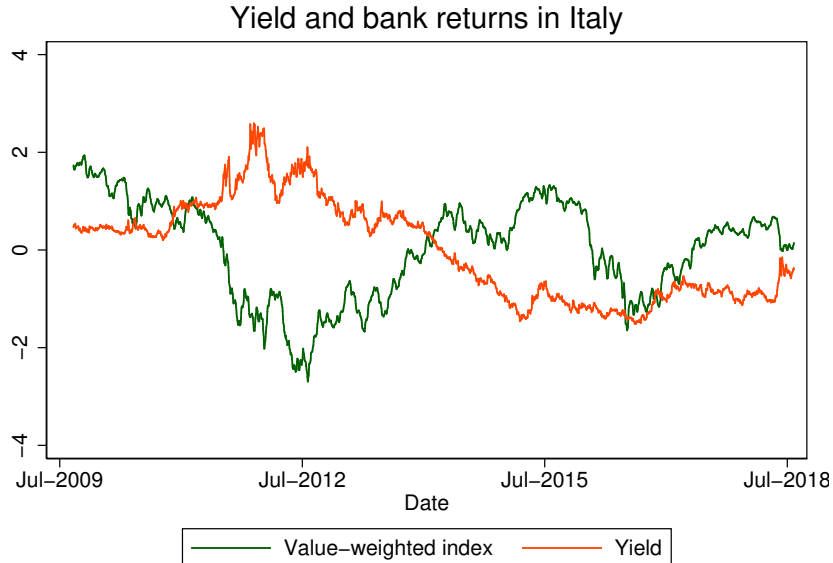


Figure 2: Sovereign yield and banking returns. The figure plots the 10-year government bond yield from September 2009 until July 2018 in Italy alongside the equity return on the value-weighted banking index.



of Δy_t . We use weekly returns at a daily frequency (using overlapping data) to mitigate the impact of non-synchronous trading.

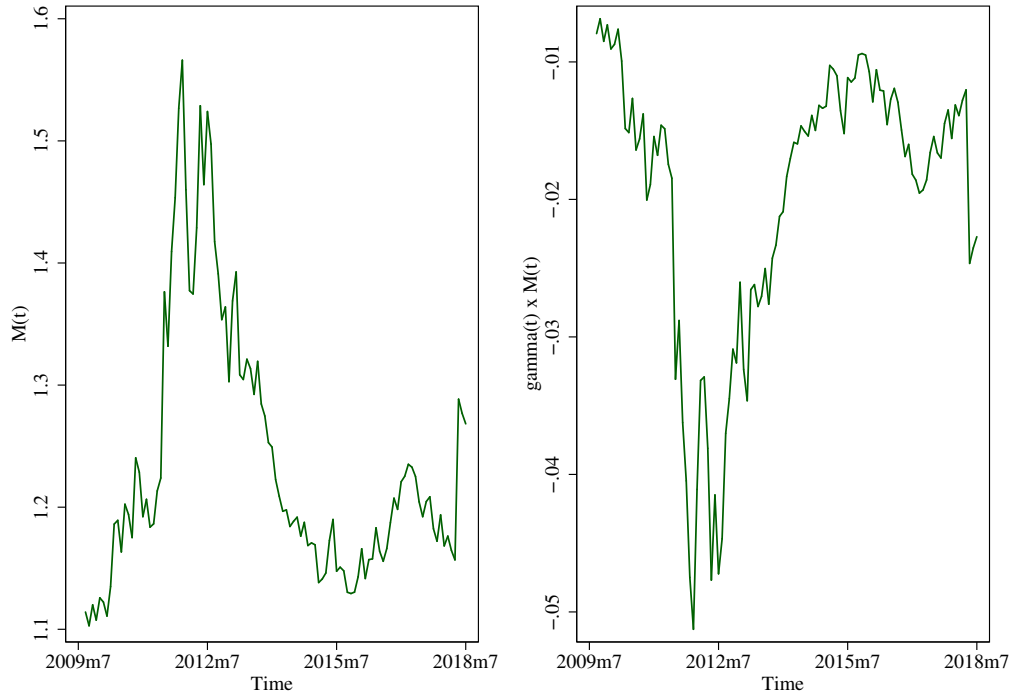
In Table [1](#), we list the banks in our sample alongside the average assets on the banks' balance sheets. As is clear from the table, the Italian banking sector is highly concentrated. In Figure [1](#), we plot the dynamics of German and Italian yields since the introduction of the Euro. During the first period, Italian yields and Germany yields move in lockstep following the introduction of the Euro. This changed drastically following the financial crisis and in particular the European sovereign debt crisis when yields diverged.

In Figure [2](#), we plot the dynamics of yields and value-weighted returns. The figure clearly illustrates the strong negative correlation between the sovereign yield and the returns on the banking sector during the European sovereign debt crisis. The same happened in recent months following the heightened political uncertainty in Italy in May 2018.

4.5 Empirical results

The estimation results are summarized in Figure [3](#). The left panel plots M_t , the impact of shocks on the banking sector, and the right panel plots $\gamma_t M_t$, which measures the impact of banking shocks on the sovereign. The impact of idiosyncratic shocks on the banking sector is close to one before the European sovereign debt crisis and thereafter. At the peak of the crisis, the multiplier equals $M_t \simeq 1.5$. The impact on sovereign yields, as depicted in the right panel, is close to zero most of the time but drops to $\gamma_t M_t \simeq -0.05$ at the peak of the crisis. It implies $\gamma_t \simeq -0.033$. All estimates

Figure 3: Estimated dynamics of M_t and $\gamma_t M_t$. The left panel plots M_t , which measures the impact of idiosyncratic shocks to banks on the banking sector, see equation (63). The right panel plots $\gamma_t M_t$, which measures the ultimate impact of idiosyncratic shocks to banks on sovereign yields, see equation (64). The multipliers are estimated using weekly returns at a daily frequency from September 2009 to June 2018.



are significant at a 5% significance level. These estimates together imply a financial doom loop that was active during the European sovereign debt crisis.

We inspect some of the largest shocks during our sample and they, for instance, correspond to a large negative shock to UniCredit in January 2012 in which they announced they would raise capital following writedowns in all of their divisions, even abroad, at a price that represented a 43% discount to the theoretical ex-rights price. Other examples include outcomes of stress tests in which investors learn that the exposures of some banks to aggregate shocks are larger than expected. This, in our framework, is a valid idiosyncratic shock.

To interpret these estimates, consider a “primary” impulse of -40% to the Italian banking sector. In response, Italian yields increase by $0.033 \times 0.4 = 132\text{bp}$. As a result of the increase in sovereign yields, banking stocks fall further, and so on. In the end, bank stocks fall by -60% and bond yields spike by 200bp.

5 Demand and supply elasticities in crude oil markets

5.1 Data

Our data construction follows the existing literature (Kilian (2009), Caldara et al. (2018), Baumeister and Hamilton (2018), henceforth BH). The data on oil supply and prices are from the U.S. Energy Information Administration (EIA). We observe the monthly oil supply for 20 countries (both OPEC and non-OPEC) from January 1985 until December 2015.²⁸ As we also observe the total non-OPEC production, we also construct a fictitious country which produces the residual non-OPEC supply. The real oil price series is obtained based on the refiner acquisition cost of imported crude oil and deflated using the US CPI to obtain the real price of oil as in Kilian (2009).

We focus on estimating short-run (monthly) demand and supply elasticities, consistent with the literature. To construct innovations, we use a state vector X_t that includes lagged (i) monthly price changes, (ii) world supply growth, (iii) changes in inventories, and (iv) growth in industrial production.²⁹ We use the data of BH for the latter two series.

5.2 Model

We model the supply growth of country i in period t as

$$\Delta y_{it} = \phi^s \Delta p_t + \lambda_i \eta_t + u_{it} + \gamma'_y X_{t-1},$$

²⁸We follow Caldara et al. (2018) and remove Gabon from the sample due to concerns about data quality. In addition, we scale the supply of the USSR using the ratio of supply of the USSR to the supply of Russia to obtain a continuous series and to avoid a sudden jump in the non-OPEC supply.

²⁹The results do not change significantly if we use 12 lags instead of one lag.

and model changes in aggregate oil demand (both in use and inventories) as

$$\Delta d_t = \phi^d \Delta p_t + \gamma'_d X_{t-1} + \epsilon_t.$$

Market clearing, $\Delta y_{St} = \Delta d_t$, implies

$$\begin{aligned} \Delta p_t &= \frac{M}{\phi^d} u_{St} + \gamma^{p'} \mathcal{C}_t, \\ \Delta y_{St} &= M u_{St} + \gamma^{y'} \mathcal{C}_t \end{aligned}$$

where

$$M = -\frac{\phi^d}{\phi^s - \phi^d} \in [0, 1]$$

is the multiplier, and γ^p, γ^y are loadings on $\mathcal{C}_t = (\eta_t, \varepsilon_t, X_{t-1})$, and whose precise value does not matter here.

Our goal is to estimate the short-run supply and demand elasticities, ϕ^s and ϕ^d (with presumably $\phi^d < 0 < \phi^s$). The equations for aggregate supply and price changes are part of the VAR models that are commonly used in the recent literature on oil prices and their impact on economic growth.

5.3 GIV estimation

The supply changes in some periods are extreme for some countries during supply disruptions, and we therefore winsorize the growth rates at 2.5% and 97.5% across all countries and periods to estimate Δy_{Et} .³⁰ We use Δy_{it}^W to denote the winsorized supply growth. We then estimate the model using the following steps:

1. Run a panel regression with country and time fixed effects.³¹

$$\Delta y_{it} = k_i + a_t + \check{e}_{it}.$$

2. Use \check{e}_{it} to estimate η_t^x and η_t^{PCA} (as in Section 3.2, Step 2), using which we define a new vector of controls $C_t = (\eta_t^x, \eta_t^{PCA}, X_{t-1})'$.
3. Estimate $\frac{M}{\phi^d}$ using (with $Z_t := \Delta y_{\Gamma t}$ which is our z_t plus some linear function of C_t , which is anyway controlled for in the regression):

$$\Delta p_t = \frac{M}{\phi^d} Z_t + \beta^{p'} C_t + e_t^p, \tag{65}$$

³⁰To ensure growth rates are always defined, we set supply to one in case it drops to zero, which happens in seven country-months.

³¹Note that the time fixed effects absorb the controls, X_{t-1} , in this case.

Table 2: Multiplier estimates in the oil market. The first column reports the estimate of M , see (65), and the second column of $\frac{M}{\phi^d}$, see (66). The third column reports the Two Stage Least Square (2SLS) estimate of the demand elasticity ϕ^d , see (67), and the fourth column the 2SLS estimate of the supply elasticity ϕ^s , see (68). We suppress the coefficients on the controls, X_{t-1} , that include lagged (i) monthly price changes, (ii) world supply growth, (iii) changes in inventories, and (iv) growth in industrial production. The t -statistics, which are reported in parentheses, are based on OLS and 2SLS standard errors. The sample is from January 1985 to December 2015.

	y_{St}	Δp_t	y_{St}	y_{Et}
Z_t	0.878 (15.12)	-2.328 (-4.42)		
$\eta_{PCA,t}$	-0.117 (-13.84)	0.176 (2.29)	-0.0508 (-1.58)	-0.126 (-12.54)
$\eta_{OPEC,t}$	0.391 (12.91)	-0.188 (-0.69)	0.320 (2.99)	0.401 (11.97)
Δp_t			-0.377 (-4.30)	0.0524 (1.91)
N	370	370	370	370
R^2	0.542	0.263		

and $M = -\frac{\phi^d}{\phi^s - \phi^d}$ using

$$\Delta y_{St} = M Z_t + \beta^{y'} C_t + e_t^y. \quad (66)$$

4. We can recover the supply and demand elasticities using the estimates of $\frac{M}{\phi^d}$ and M , where $\phi^s = \frac{\phi^d}{M}(M - 1)$. However, to get the standard errors on the elasticities as well, we use the 2SLS estimator based on the first stage, which corresponds to (65), and denote the fitted value by $\Delta \hat{p}_t := \left(\frac{M}{\phi^d}\right)^e Z_t$. The second stage estimator for the demand elasticity corresponds to

$$\Delta y_{St} = \phi^d \Delta \hat{p}_t + \beta_d' C_t + e_t^d, \quad (67)$$

and for the supply elasticity to

$$\Delta y_{Et} = \phi^s \Delta \hat{p}_t + \beta_s' C_t + e_t^s. \quad (68)$$

5.4 Empirical results

We report the estimation results of the multipliers $M = 0.88$ and $\frac{M}{\phi^d} = -2.3$ in Table 2 alongside both elasticities. We estimate a demand elasticity of $\phi^d = -38\%$ (with a standard error of 9%) and a supply elasticity of $\phi^s = 5\%$ (with a standard error of 3%). Changes in demand also include

changes in inventories, which respond more elastically to changes in prices (Kilian and Murphy (2014)).

To put these estimates in perspective, we compare them to recent estimates in the literature. Baumeister and Hamilton (2018) use sign restrictions in combination with a Bayesian estimator to find supply and demand elasticities of 15% and -35%, respectively, with 68% credibility intervals of (9%, 22%) for the supply elasticity and (-51%, -24%) for the demand elasticity. Caldara et al. (2018) use a narrative approach and estimate a supply elasticity of 8% (with a standard error of 3.7%) and a demand elasticity of -8% (with a standard error of 8%). Kilian and Murphy (2014) also combine sign restrictions and a Bayesian estimator with short-run supply elasticities bounded at 2.5%, 5%, and 10%, and corresponding demand elasticities range from -44% to -47%.

We construct our instrument as the residual from a regression of $y_{\Gamma t}$ on X_{t-1} and the two factors and refer to it as $u_{\Gamma t}$. If we regress it on the instrument of Caldara et al. (2018), which is non-zero only during 14 months in this sample, we get a slope coefficient of 0.93, with a t-statistic of 12.7 and an R-squared of 93%. Moreover, if we restrict ourselves to more extreme episodes by only using data when $u_{\Gamma t}$, in absolute value, exceeds a threshold of 0.5% (370 observations), 0.75%, ..., 1.25% (19 observations), then the first-stage estimate declines monotonically from -2.3, with a t-statistic of -4.4, to -5.2, with a t-statistic of 3.5. This highlights that by focusing on the more extreme events, $\frac{M}{\phi d}$ becomes more negative. Intuitively, in case of more extreme shocks, the role of inventories diminishes and the demand curve becomes more inelastic. This reconciles our estimates with those of Caldara et al. (2018).

If we inspect the largest shocks in terms of contribution to the instrument, $S_{i,t-1}u_{it}$, then many of the extreme shocks are as described in Caldara et al. (2018). However, in some cases, the GIV identifies shocks that are not included in the narrative approach. An example includes a reduction in supply by Saudi Arabia in January 1989. Per the description of Caldara et al. (2018), OPEC agreed upon a reduction in supply in November of 1988 but reports in subsequent months were interpreted as “indicating that the OPEC member country was seriously attempting to cut back production based on the new agreement.” One possible interpretation of this shock is that markets learn about the exposure to the common OPEC shock, $\eta_{OPEC,t}$. In January 1989, the real price of oil jumped up by 13.2%. In the context of GIV, those are valid idiosyncratic shocks that can be used as instruments.

In summary, the GIV estimator results in estimates that are in the range of estimates documented in the recent literature, thereby providing some external validation of GIVs as an approach to estimating demand and supply elasticities. At the same time, the GIV procedure arguably requires less domain-specific ingenuity than the previous studies we mentioned.

In future work, granular country-level data on (net) imports and oil consumption can be used to construct a second instrument that can be used to both sharpen the estimates and to test for overidentifying restrictions. This instrument may be particularly powerful given the apparent importance of demand shocks during the last 15 years.

6 Discussion and extensions of the framework

6.1 Extensions

There are many ways to increase the number of setups in which the GIV idea can be applied.

Multidimensional GIV One can handle multidimensional “actions”: for instance, a firm could have a shock that affect both productivity and labor demand. A country could have a shock that affects both productivity and oil demand. Formally, the actions y_{it} and shocks u_{it} are now multidimensional. The GIV idea goes through, and this is developed in Section 12.1. We have seen that with one GIV, we can estimate $1 + d_F$ parameters ($M, M\alpha^f$), where d_F is the number of endogenous, observed factors. With q -dimensional actions, we have q GIVs, and we can estimate $q^2 + qd_F$ parameters, which correspond to M and α^f .³² So, potential many parameters can be recovered with multidimensional “actions” by firms or countries.

GIV with different size weights This framework can be extended with size weights that vary across factors, $F_t^f = \eta_t^f + \alpha^f y_{S^f,t} + k^f + \phi^f C_t^F$. Then, we can identify more parameters, as each $z_{S^f,t} = y_{S^f,t} - y_{E,t}$ is an instrument (see Section 12.2), for each distinct and useful weight S^f . Indeed, then we can not only identify M and $M\alpha^f$ as in the regular GIV, but also all the λ^f .³³

GIV with a more complex matrix of influences The GIV can also be extended to non-homogeneous influences. Suppose a model:

$$y_{it} = \gamma \sum_j G_{ij} y_{jt} + \lambda_i \eta_t + u_{it}, \quad (69)$$

i.e., in vector form

$$y_t = \gamma G y_t + \Lambda \eta_t + u_t, \quad (70)$$

with a given “influence” matrix G (in our baseline model, $G = \iota S'$). We’d like to identify γ , the strength of linkages.

A simple generalization of our GIV is to define a “size” vector $S := G'E$. Then, left-multiplying (70) by E' , we get

$$y_{Et} = \gamma y_{St} + \Lambda_S \eta_t + u_{St}.$$

The key moment is still $\mathbb{E}[(y_{Et} - \gamma y_{St}) z_t] = 0$, where the GIV is again $z_t = y_{St} - y_{Et}$ in the simple case where $\Lambda = \iota$ and $G\iota = \iota$; see Section 12.6 for the general case. Hence, GIV generalizes to “spatial” models with common shocks (most spatial models do not have latent common shocks).

³²As u_{St} and y_{St} are q -dimensional, $M = \frac{dy_{St}}{du_{St}}$ $q \times q$ dimensional, and each of the $\frac{dF_t^f}{du_{St}} = M\alpha^f$ is also q -dimensional.

³³On the other hand, the difference between different size weights may be small to the estimation will be more fragile.

Bayesian GIV Another extension is a Bayesian interpretation of the GIVs. This way, we can interpret do GIVs in a Bayesian framework – see Section 14. In particular, this opens the possibility of marrying GIV estimation with priors on other parameters. In the simplest cases with Gaussian shocks, the maximum likelihood estimate is our GIV – confirming its optimality properties. At the same time, the basic GIV doesn't actually use normality assumptions.

Furthermore, many econometric extensions might be useful, e.g. with stochastic volatility, and various dimensions of autocorrelations. We leave those extensions to future research.

6.2 Discussion: Robustness to misspecification

The GIV procedure is robust to some forms of misspecification, and more fragile to others.

We may keep only the shocks to some actors (in a space I_t), i.e. set $z_t = \sum_{i \in I_t} S_i (u_{it} - u_{Et})$, selecting for example the shocks to the top K entities, the shocks for which we have data, or some subset of the entities based on size. Then again, everything goes through.³⁴ The estimator is still valid, just not the optimal GIV estimator.

Suppose that we misspecify the vector S of sizes, for example by defining $z_t = \sum_i S_i^\circ (u_{it} - u_{Et})$ using a wrong vector S° . Then, the IV is still valid, but the OLS can be biased. In our basic example of Section 2.2, we still have $\mathbb{E}[(p_t - \alpha y_{St}) z_t] = 0$, so that the IV procedure (16) still works. Likewise, in the more complex supply and demand case, the IV relations (34) and (36) still hold. But the OLS relations are slightly biased.³⁵

If we assume homogeneous coefficients (e.g. on the elasticities of demand or supply), while in fact they are truly heterogeneous, then again (assuming that η_t was well estimated in the cross-section) the IV estimates are correct, and so are the OLS estimates.³⁶

If we misspecify the variance of the u_{it} (but keeping them uncorrelated), things are essentially fine: as $u_E = O_p\left(\frac{1}{\sqrt{N}}\right)$, we do not need $\mathbb{E}[u_{\Gamma t} u_{Et}] = 0$, but the term $\mathbb{E}[u_{\Gamma t} u_{Et}]$ will still be small ($O\left(\frac{1}{\sqrt{N}}\right)$), and vanishes for large N .

A threat is that we might not control properly for common factors. Indeed, $z_t = u_{\Gamma t} + \lambda_{\Gamma} \eta_t - \lambda_{\Gamma}^e \eta_t^e$, so there is a danger that, even after controlling for η_t^e in the regression we will not completely eliminate the $\lambda_{\Gamma} \eta_t - \lambda_{\Gamma}^e \eta_t^e$ error.³⁷ This danger is greater when $|\lambda_{\Gamma}|$ is greater, i.e. when loadings are correlated with size. This is a small sample problem (with a large enough T, N we measure η_t, λ accurately). One solution to it is to do a Monte-Carlo simulation (as in Section 7) to evaluate the size of the residual error and potential bias, and to correct the point estimates accordingly.

³⁴For instance, we still have $u_{St} = z_t + \varepsilon_t^{uS}$ with $z_t \perp \varepsilon_t^{uS}$.

³⁵Calling $\psi = \frac{\mathbb{E}[z_t u_{St}]}{\mathbb{E}[z_t^2]}$ (which is 1 when $S^\circ = S$), then the OLS above gives (in population) $b^{p,e} = b^p \psi$ and $M^e = M \psi$. For some selection procedures (e.g. selecting the shocks to some pre-specified entities as we discussed), we still have that $\psi = 1$, so that OLS is still valid.

³⁶Just, the IV estimates yield ϕ^s , ϕ_E^d , and the OLS coefficients are those corresponding to the interpretation that the elasticity of demand is ϕ_E^d rather than ϕ_S^d (see Section 12.9).

³⁷As we do control for η_t^e in the regression, the bias is due to the residual of $\lambda_{\Gamma} \eta_t - \lambda_{\Gamma}^e \eta_t^e$ after controlling for η_t^e .

6.3 Comparison with Bartik instruments and other procedures

Comparison with Bartik instruments The GIV estimator shares some similarities with Bartik instruments, also known as “shift-share estimators,” that were first introduced in [Bartik \(1991\)](#). To put it simply, Bartik instruments allow to estimate the cross-sectional (or micro) sensitivities to shocks, but not aggregate sensitivities; whereas GIVs are mostly designed to estimate aggregate (or macro) elasticities. Hence, they are complementary.

To see this, let us use the notation established earlier. Shift-share estimators aim to estimate the coefficients λ_1^f in the structural equation $y_{it} = \sum_f (\lambda_0^f + \lambda_1^f x_{it}) F_t^f + \eta_t^y + u_{it}$ using $x_{it} g_t^f$ as an instrument for $x_{it} F_t^f$. In this notation, the “shares” are x_{it} and the “shifters” are g_t^f (for instance g_t^f could be the China shock, and be correlated with η_t^y). Shift-share estimators have been the study of much recent econometric work including [Goldsmith-Pinkham et al. \(2018\)](#); [Adao et al. \(2018\)](#); [Borusyak et al. \(2018\)](#). [Borusyak et al. \(2018\)](#) lay out sufficient identifying conditions for the shift-share estimator to estimate the structural parameter of interest λ_1^f and show that the key orthogonality condition is that the shifters g_t^f are orthogonal to the share-weighted structural disturbances. That is, the shifters are as-good-as-randomly assigned. [Goldsmith-Pinkham et al. \(2018\)](#) provide alternative identifying conditions for shift-shares but these are less relevant for the GIV estimator.

Returning to the GIV estimator and the notation we established earlier, recall that the shares S_i are either held fixed throughout the analysis or set in the previous period, e.g. $S_i = S_{i,t-1}$ provided that the previous period shares are orthogonal to u_{it} . Therefore, a critical orthogonality condition for the GIV estimator is that the idiosyncratic errors u_{it} are orthogonal to the disturbances in the structural equation of interest. In this sense, the orthogonality condition for the GIV estimator is similar to the condition provided in [Borusyak et al. \(2018\)](#), where we now think of the idiosyncratic errors u_{it} as the shifter. The GIV estimator then provides a very general strategy for constructing valid instruments based upon the underlying granular economic structure and as shown earlier, these granular instruments are optimal instruments. However, this does not fully capture the contribution of the GIV estimator. Shift-share estimators are unable to estimate the mean effect λ_0^f . Moreover, as we also show earlier, the GIV approach identifies multiple parameters in a system of simultaneous equations $(M_t, M_t \alpha^f)$ and therefore it additionally enables the researcher to identify multipliers. This is generally not true in shift-share settings, which typically consider single-equation systems.

Procedures containing elements of GIVs A few papers have explored the idea of using idiosyncratic shocks as instruments to estimate spillover effects, such as [Leary and Roberts \(2014\)](#) in the context of capital structure choice of firms and [Amiti et al. \(2019\)](#) in the context of price setting decisions of firms. The structure of the estimating equation is similar to the model that we

consider in this paper³⁸

$$y_t = \lambda y_{wt} + mC_t + u_t,$$

where $y_{wt} = w'y_t$ can be equally-weighted (Leary and Roberts (2014)) or size-weighted (Amiti et al. (2019)), depending in the weights w . Both papers use industry and/or year fixed effects, which can be viewed as a choice of controls or exogenous factors, η_t , to which all firms in a given industry have the same exposure.

There are two main differences compared to GIV. First, both papers use idiosyncratic shocks to another variable than y_t , say g_t , to construct an instrument for y_{wt} . Leary and Roberts (2014) use idiosyncratic stock returns and Amiti et al. (2019) use shocks to competitors' marginal cost, exchange rates, or export prices. To estimate λ , g_{wt} is used as an instrument for y_{wt} . We, instead, propose to use idiosyncratic shocks to y_t rather than another instrument (this way requiring fewer times series). Second, and related, we control for heterogeneous exposures to common factors to extract the idiosyncratic shocks, which is important in asymptotic theory and practice in realistic samples (see Section 7).

A third difference is specific to Leary and Roberts (2014). GIVs crucially depend on the difference between size- and equal-weighted averages of variables. If the estimating equation depends on equal-weighted averages, GIV cannot be applied. In most models, however, not all competitors receive equal weight and larger firms, or perhaps firms that are closer in product space, receive a larger weight.

Lastly, the use of model-based idiosyncratic shocks has some similarities with Amiti and Weinstein (2018), who extract bank supply shocks from Japanese data using a panel of fixed effects, and then estimate the sensitivity of aggregate investment to these shocks. However, unlike our model, Amiti and Weinstein (2018) assume a uniform sensitivity to the aggregate shocks ($\lambda_i \eta_t$ with $\lambda_i = 1$ for all i), and do not allow for feedback loops: shocks to banks affect aggregate investment, but aggregate investment does not circle back around to affect individual bank behavior (so, they assume $\alpha^f = 0$ in our notations). This is the key source of endogeneity in many of the models we consider, and by tackling it we are able to estimate a richer set of parameters.

Other methods to estimate aggregate elasticities Rigobon (2003) introduces another method that can be used to estimate spillover effects and aggregate multipliers using time-variation in second moments. If shocks are heteroskedastic and the structural parameters are stable across regimes, then the different volatility regimes add additional equations to the system so that the structural parameters can be identified. GIV does not require heteroskedasticity, but can accommodate it, and is therefore complementary to identification methods that rely on heteroskedasticity.

³⁸Amiti et al. (2019) study the price setting decision of firms. In their model, the pricing equation features two endogenous variables, namely the same firm's marginal cost and the size-weighted average of competitors' prices. We focus on the spillover effects of competitors' prices in our discussion in this section.

Influence and the “reflection problem” We finish by another example, known as “contagion” or the “reflection problem” (Manski (1993); Kline and Tamer (Forthcoming)). Suppose that actions follow:

$$y_{it} = \gamma y_{St} + \eta_t + u_{it}, \quad (71)$$

where η_t is uncorrelated with the u_{it} . This equation means that y_{it} is influenced by the aggregate action of other agents (γy_{St}), and in addition to the usual aggregate shocks η_t , and idiosyncratic shocks u_{it} (which we assume to be uncorrelated to η_t). The “influence” or “contagion” parameter γ is of high interest.

The GIV approach works as follows. Taking the size-weighted average of (71), we have $y_{St} = \gamma y_{St} + \eta_t + u_{St}$, so that

$$y_{St} = M (\eta_t + u_{St}), \quad M = \frac{1}{1 - \gamma}. \quad (72)$$

We form $z_t := y_{\Gamma t}$, which by (71) will give $y_{\Gamma t} = u_{\Gamma t}$. Hence, if we estimate M^e by OLS:

$$y_{St} = M^e y_{\Gamma t} + \varepsilon_t^y,$$

then we have a consistent estimator of the multiplier M , and therefore of γ .³⁹ We have a simple GIV approach to the “reflection problem”. To the best of our knowledge, this approach is new. Indeed, it may seem to contradict earlier impossibility results. Section 12.7 solves the apparent contradiction. The short summary is that Manski (1993) and Bramoullé et al. (2009) do not consider anything like a GIV, as they immediately reason on averages based on observables, eschewing any exploration of the noise.⁴⁰ In contrast, GIVs are all about exploring some structure in the noise — the idiosyncratic shocks of large entities.⁴¹

In a tangentially related recent paper, Sarto (2018) uses factor analysis to extract values of η^f (much as we do when we “recover” a factor η^f). Take the basic example in our paper. Then, Sarto does not identify α : even if η (the aggregate shock to demand) were perfectly identified, that would not allow to estimate p . In the supply and demand example, Sarto would identify the demand elasticity ϕ^d , but not the supply elasticity ϕ^s .

Spatial econometrics. In some applications of GIVs we have considered separately, growth in a region affects that of the other regions. So there is a similarity between our setup and that of spatial econometrics (e.g. Kelejian and Prucha (1999); Blasques et al. (2016); Shi and Lee (2017); Kuersteiner and Prucha (2018)). However, the estimators are quite different. The reason is that spatial econometrics studies the “local” influence (e.g. of neighboring cities on a city), while GIVs

³⁹And we will have $\varepsilon_t^y = \eta_t + u_{Et}$.

⁴⁰Somewhat related, Graham (2008) explores the identification of peer effects using conditional variance restrictions on the outcomes by exploiting differences in the sizes of the peer group. Intuitively, smaller peer group sizes leads to a larger contribution of each individual peer on the peer component.

⁴¹Economically, the idiosyncratic shocks to “big influencers” (e.g. large firms, or perhaps famous people in the networks) affect the aggregate, hence they allow to estimate the social or economic multiplier. This is why they can be handled with GIVs.

study the global influence. Hence, the sources of variation, identifiability conditions and methods are quite different. Certainly, the spatial literature has not identified, as we do, the GIVs as a simple way to estimate elasticities in contexts such as supply and demand problems, and models with feedback loops from banks to sovereign yields (and vice versa). Still, some of the sophisticated techniques of the spatial literature might be used one day to enrich a GIV-type of analysis.

7 Simulations

We illustrate the precision of granularly identified parameters depending on the size of the sample (both N and T), the degree of concentration, and the volatility of idiosyncratic shocks relative to aggregate shocks.

7.1 Model

We start from the standard supply, y_{it}^s , and demand, y_t^d , model

$$y_{it}^s = \phi^s p_t + \lambda'_i \eta_t + u_{it}, \quad y_t^d = \phi^d p_t + \epsilon_t,$$

where $\phi^d < 0 < \phi^s$, implying, with $M = -\frac{\phi^d}{\phi^s - \phi^d}$,

$$p_t = \frac{M}{\phi^d} (u_{St} + \lambda'_S \eta_t - \epsilon_t), \quad y_{St}^s = M u_{St} + M \lambda'_S \eta_t + (1 - M) \epsilon_t.$$

7.2 Estimators and standard errors

To estimate M and $\frac{M}{\phi^d}$, we can use standard OLS. To estimate M , we use

$$y_{St}^s = a + M y_{\Gamma t} + \theta' \eta_t^e + e_t, \tag{73}$$

and to estimate $\frac{M}{\phi^d}$, we use

$$p_t = a_p + \frac{M}{\phi^d} y_{\Gamma t} + \theta'_p \eta_t^e + e_t^p. \tag{74}$$

All standard OLS results apply if we observe the factors, η_t . However, we often do not directly observe all factors. We consider the case in which we know the factor loadings, λ_i^η , and where the loadings are unobserved and estimated using PCA. To provide a point of reference, we also consider the case where we do not control for factors and impose that $\theta = \theta_p = 0$. In all cases, we report the OLS standard errors to assess to what extent the OLS standard errors need to be adjusted for the fact that we use η_t^e instead of η_t .

To estimate the demand and supply elasticities, we can recover them from the estimates of M and $\frac{M}{\phi^d}$. However, as discussed before, this is equivalent to a 2SLS estimator using $y_{\Gamma t}$ as instrument

for price, while controlling in this case for the factors. Hence, the first stage corresponds to

$$p_t = a_p + \xi y_{\Gamma t} + \theta'_p \eta_t^e + e_t,$$

and the second stage to estimate the demand elasticity is, with $\hat{p}_t = a_p^e + \xi^e y_{\Gamma t} + \theta'_p \eta_t^e$,

$$y_t^d = a_d + \phi^d \hat{p}_t + \theta'_d \eta_t^e + e_t^d,$$

and for the supply elasticity

$$y_{Et}^s = a_s + \phi^s \hat{p}_t + \theta'_s \eta_t^e + e_t^s.$$

The standard weak instrument tests can be used to assess whether $y_{\Gamma t}$ is a sufficiently strong instrument for price (Section 2.3.1). In this case, we report the 2SLS standard errors to assess whether their accuracy is impacted by the fact that we estimate the common factors.

7.3 Calibration

In calibrating the model, we target (i) concentration, as measured by the excess Herfindahl, $h = \sqrt{\sum_i S_i^2 - 1/N}$, and (ii) the ratio of the volatility of idiosyncratic shocks to the volatility of aggregate supply shocks. In all cases, we estimate the number of common factors using the procedure in Bai and Ng (2002) by minimizing their $IC_{p2}(k)$ criterion.

We set $\phi^d = -0.3$, $\phi^s = 0.1$, and $\sigma_\epsilon = 3\%$. The size weights are generated as $k_i = i^{-1/\zeta}$, $S_i = k_i / \sum_i k_i$, where ζ is chosen so that $h \in \{0.2, 0.3\}$.⁴² In the benchmark case, we assume a single common factor, which follows a standard normal distribution, and uniformly distributed loadings. We consider two cases, namely where $\text{Corr}(\lambda, S) = 0$ and $\text{Corr}(\lambda, S) = -20\%$. We scale the loadings so that the variance of aggregate supply shocks follows $V(\lambda'_S \eta_t) = \lambda_S^2 = 0.03^2$. Lastly, we select $\sigma_u = \tau(\lambda'_S \lambda_S)^{1/2}$ to target the ratio τ of idiosyncratic shock volatility to aggregate shock volatility. We vary $\tau \in \{3, 4\}$, $N \in \{25, 50\}$, and $T \in \{120, 360\}$.

The cases considered are summarized in Table 3. The final column reports the fraction of price volatility that is due to idiosyncratic shocks, which ranges approximately from 10% to 30%, in line with the recent literature on granularity in terms of how much of aggregate fluctuations can be traced back to idiosyncratic shocks.

7.4 Simulation results

The simulation results when $\text{Corr}(\lambda, S) = 0$ are reported in Table 4. We consider four estimators. In the case of M1, we assume that the loadings are known in estimating the factors; this is an ideal case taken as a benchmark. In the case of M2, we use PCA to estimate the factors. In the case of M3, we control for the factors estimated using the known loadings and PCA. In the case of

⁴²Here ζ is the power law exponent of the size distribution. See Gabaix (2009) and Section 12.4

Table 3: Cases considered in simulations. We calibrate the supply-and-demand model in Section 7 under seven alternative parameterizations. The parameters are the following: N is the cross-sectional sample size; T is the number of simulated i.i.d. time periods; h is the excess Herfindahl that we target in our simulation of the size weights (as described in Section 7.3); τ is the targeted ratio of the volatility of idiosyncratic shocks to the volatility of aggregate supply shocks; the multipliers $M = -\frac{\phi^d}{\phi^s - \phi^d}$ and $\frac{M}{\phi^d}$ are functions of the elasticities ϕ^d and ϕ^s of demand and supply with respect to price. The final column reports the share of the price volatility that is due to idiosyncratic shocks under each of the seven parameterizations.

Case	N	T	h	τ	M	$\frac{M}{\phi^d}$	% price vol. idiosyncratic
1	25	360	0.2	3	0.75	-2.5	12.6%
2	25	360	0.2	4	0.75	-2.5	20.4%
3	25	360	0.3	3	0.75	-2.5	19.1%
4	25	360	0.3	4	0.75	-2.5	29.5%
5	25	120	0.2	4	0.75	-2.5	20.4%
6	50	120	0.2	4	0.75	-2.5	16.1%
7	50	360	0.2	4	0.75	-2.5	16.1%

M4, we use no factors and just use $y_{\Gamma t}^s$ without any factors. Note that we do not advocate M4 in practice: M4 is there simply to illustrate what goes wrong if we don't control for factors. The first four columns correspond to the estimates of M , the next four columns to estimates of $\frac{M}{\phi^d}$, the next four columns to estimates of ϕ^d , and the last four columns to estimates of ϕ^s .

For each of the estimators, we report the median, the mean, and the 2.5% and 97.5% percentiles. We also compute the fraction of estimates that fall within the 95% confidence intervals constructed using OLS standard errors (columns 1 to 8) or the 2SLS standard errors (columns 9 to 16). We refer to this as the “coverage.”

As is clear from all cases, the estimators are mean- and median-unbiased. Moreover, confidence intervals tighten when concentration increases (case 3 relative to case 1 and case 4 relative to case 2) and when the volatility of idiosyncratic shocks increases (case 2 relative to case 1 and case 4 relative to case 3). Naturally, the confidence interval tightens when we increase N and T . The coverage is generally accurate and OLS standard errors only slightly overstate the precision in the case of M2 in estimating M ; the 2SLS standard errors are somewhat small in small samples in estimating ϕ^s .

It is tempting to conclude that using $y_{\Gamma t}^s$ as instrument, even without estimating the factors, results in accurate and unbiased estimates of the parameters of interest. However, this is only the case when $\text{Corr}(\lambda, S) = 0$. To illustrate this, we consider a negative correlation between size and exposures, $\text{Corr}(\lambda, S) = -20\%$.

The results are presented in Table 5. Now we find a large bias in case of M4, both in terms of the mean and median. The coverage estimates are also heavily distorted. Intuitively, $y_{\Gamma t}^s$ does not filter out aggregate shocks and the exogeneity restriction is violated. This is why factor estimates are required when loadings may be correlated with size. Even in the case where we have no information

Table 4: Simulation results when $Corr(\lambda, S) = 0$ based on 10,000 replications. The parameters used in the different cases are summarized in Table 3. In particular, the data are generated from a model in which $M = 0.75$, $\frac{M}{\phi^d} = -2.5$, $\phi^d = -0.3$, and $\phi^s = 0.1$. GIV estimators M1,..., M4 are described at the beginning of Section 7.4. For each estimator, we report the median, the mean, and percentiles 2.5% (P2.5) and 97.5% (P97.5) in the simulated distribution of estimates. “Coverage” is the fraction of estimates falling within the 95% confidence intervals constructed using OLS standard errors (columns 1 through 8) or the 2SLS standard errors (columns 9 through 16).

Case	Statistic	M				$\frac{M}{\phi^d}$				ϕ^d				ϕ^s			
		M1	M2	M3	M4	M1	M2	M3	M4	M1	M2	M3	M4	M1	M2	M3	M4
1	Median	0.75	0.75	0.75	0.75	-2.50	-2.50	-2.50	-2.49	-0.30	-0.30	-0.30	-0.30	0.10	0.10	0.10	0.10
	Mean	0.75	0.75	0.75	0.75	-2.50	-2.50	-2.50	-2.50	-0.33	-0.34	-0.32	-0.34	0.11	0.12	0.11	0.11
	P2.5	0.53	0.49	0.52	0.51	-3.81	-3.90	-3.86	-3.83	-0.65	-0.67	-0.67	-0.65	0.01	0.00	0.00	0.00
	P97.5	0.98	1.01	0.99	0.99	-1.20	-1.11	-1.15	-1.17	-0.17	-0.17	-0.17	-0.17	0.27	0.32	0.28	0.30
	Coverage	0.95	0.93	0.95	0.95	0.95	0.94	0.95	0.95	0.93	0.93	0.93	0.93	0.96	0.94	0.95	0.96
2	Median	0.75	0.75	0.75	0.75	-2.51	-2.50	-2.51	-2.51	-0.30	-0.30	-0.30	-0.30	0.10	0.10	0.10	0.10
	Mean	0.75	0.75	0.75	0.75	-2.51	-2.51	-2.51	-2.51	-0.31	-0.31	-0.32	-0.31	0.10	0.11	0.10	0.10
	P2.5	0.60	0.56	0.59	0.58	-3.48	-3.54	-3.50	-3.52	-0.52	-0.52	-0.53	-0.52	0.04	0.02	0.04	0.03
	P97.5	0.90	0.94	0.90	0.92	-1.54	-1.47	-1.51	-1.49	-0.20	-0.19	-0.19	-0.19	0.19	0.22	0.19	0.21
	Coverage	0.95	0.90	0.95	0.95	0.95	0.93	0.95	0.95	0.94	0.94	0.94	0.94	0.96	0.91	0.96	0.95
3	Median	0.75	0.75	0.75	0.75	-2.51	-2.50	-2.51	-2.51	-0.30	-0.30	-0.30	-0.30	0.10	0.10	0.10	0.10
	Mean	0.75	0.75	0.75	0.75	-2.51	-2.51	-2.51	-2.51	-0.31	-0.31	-0.31	-0.31	0.10	0.10	0.10	0.10
	P2.5	0.62	0.59	0.62	0.61	-3.27	-3.34	-3.29	-3.29	-0.44	-0.45	-0.45	-0.44	0.05	0.03	0.04	0.04
	P97.5	0.88	0.92	0.88	0.89	-1.74	-1.68	-1.72	-1.72	-0.22	-0.22	-0.21	-0.22	0.17	0.21	0.18	0.18
	Coverage	0.95	0.90	0.95	0.95	0.94	0.93	0.95	0.94	0.94	0.94	0.95	0.95	0.95	0.90	0.95	0.95
4	Median	0.75	0.75	0.75	0.75	-2.50	-2.50	-2.50	-2.51	-0.30	-0.30	-0.30	-0.30	0.10	0.10	0.10	0.10
	Mean	0.75	0.75	0.75	0.75	-2.51	-2.50	-2.51	-2.51	-0.31	-0.31	-0.31	-0.31	0.10	0.10	0.10	0.10
	P2.5	0.65	0.63	0.65	0.64	-3.15	-3.18	-3.17	-3.18	-0.42	-0.42	-0.42	-0.42	0.06	0.05	0.06	0.05
	P97.5	0.85	0.87	0.85	0.87	-1.86	-1.81	-1.84	-1.84	-0.22	-0.22	-0.22	-0.22	0.15	0.17	0.15	0.16
	Coverage	0.95	0.90	0.95	0.95	0.94	0.93	0.95	0.95	0.95	0.95	0.95	0.95	0.95	0.89	0.95	0.95
5	Median	0.75	0.75	0.75	0.75	-2.51	-2.50	-2.50	-2.50	-0.30	-0.30	-0.30	-0.30	0.10	0.10	0.10	0.10
	Mean	0.75	0.75	0.75	0.75	-2.50	-2.50	-2.50	-2.50	-0.30	-0.30	-0.30	-0.30	0.10	0.10	0.10	0.10
	P2.5	0.66	0.64	0.66	0.66	-3.01	-3.05	-3.02	-3.02	-0.38	-0.38	-0.39	-0.38	0.06	0.05	0.06	0.06
	P97.5	0.84	0.86	0.84	0.84	-1.99	-1.95	-1.99	-1.98	-0.24	-0.24	-0.24	-0.24	0.14	0.16	0.15	0.15
	Coverage	0.95	0.90	0.95	0.95	0.94	0.93	0.95	0.95	0.95	0.95	0.95	0.95	0.95	0.89	0.95	0.95
6	Median	0.75	0.75	0.75	0.75	-2.50	-2.50	-2.50	-2.50	-0.30	-0.30	-0.30	-0.30	0.10	0.10	0.10	0.10
	Mean	0.75	0.75	0.75	0.75	-2.50	-2.49	-2.50	-2.50	-0.32	-0.33	-0.33	-0.33	0.11	0.11	0.11	0.11
	P2.5	0.56	0.52	0.55	0.52	-3.71	-3.78	-3.74	-3.77	-0.64	-0.66	-0.65	-0.64	0.02	0.01	0.02	0.01
	P97.5	0.94	0.98	0.95	0.97	-1.24	-1.16	-1.24	-1.20	-0.17	-0.17	-0.17	-0.17	0.22	0.26	0.22	0.26
	Coverage	0.95	0.93	0.95	0.95	0.95	0.95	0.95	0.95	0.93	0.93	0.93	0.93	0.96	0.95	0.96	0.96
7	Median	0.75	0.75	0.75	0.75	-2.50	-2.50	-2.50	-2.50	-0.30	-0.30	-0.30	-0.30	0.10	0.10	0.10	0.10
	Mean	0.75	0.75	0.75	0.75	-2.50	-2.50	-2.50	-2.50	-0.31	-0.31	-0.31	-0.31	0.10	0.10	0.10	0.10
	P2.5	0.64	0.62	0.64	0.62	-3.21	-3.27	-3.23	-3.24	-0.44	-0.44	-0.45	-0.44	0.06	0.05	0.06	0.05
	P97.5	0.86	0.88	0.86	0.88	-1.78	-1.75	-1.77	-1.75	-0.22	-0.22	-0.22	-0.22	0.15	0.17	0.16	0.17
	Coverage	0.95	0.91	0.95	0.95	0.95	0.94	0.95	0.95	0.94	0.95	0.95	0.95	0.95	0.91	0.95	0.95

Table 5: Simulation results when $Corr(\lambda, S) = -20\%$ based on 10,000 replications. The parameters used in the different cases are summarized in Table 3. In particular, the data are generated from a model in which $M = 0.75$, $\frac{M}{\phi^d} = -2.5$, $\phi^d = -0.3$, and $\phi^s = 0.1$. GIV estimators M1,..., M4 are described at the beginning of Section 7.4. For each estimator, we report the median, the mean, and percentiles 2.5% (P2.5) and 97.5% (P97.5) in the simulated distribution of estimates. “Coverage” is the fraction of estimates falling within the 95% confidence intervals constructed using OLS standard errors (columns 1 through 8) or the 2SLS standard errors (columns 9 through 16).

Case	Statistic	M				$\frac{M}{\phi^d}$				ϕ^d				ϕ^s			
		M1	M2	M3	M4	M1	M2	M3	M4	M1	M2	M3	M4	M1	M2	M3	M4
1	Median	0.75	0.69	0.75	0.56	-2.49	-2.29	-2.50	-1.85	-0.30	-0.30	-0.30	-0.30	0.10	0.13	0.10	0.24
	Mean	0.75	0.69	0.75	0.56	-2.50	-2.30	-2.50	-1.85	-0.33	-0.34	-0.33	-0.30	0.11	0.16	0.11	0.21
	P2.5	0.52	0.41	0.51	0.29	-3.84	-3.72	-3.87	-3.23	-0.66	-0.74	-0.67	-0.97	0.01	0.01	0.00	0.07
	P97.5	0.98	0.97	0.99	0.82	-1.17	-0.88	-1.13	-0.50	-0.17	-0.16	-0.17	-0.14	0.28	0.47	0.29	0.96
	Coverage	0.95	0.88	0.95	0.67	0.95	0.93	0.95	0.84	0.93	0.93	0.93	0.92	0.96	0.97	0.95	1.00
2	Median	0.75	0.74	0.75	0.42	-2.51	-2.46	-2.51	-1.38	-0.30	-0.30	-0.30	-0.30	0.10	0.11	0.10	0.42
	Mean	0.75	0.74	0.75	0.41	-2.51	-2.47	-2.51	-1.39	-0.31	-0.32	-0.32	0.00	0.10	0.11	0.10	-0.12
	P2.5	0.60	0.56	0.59	0.19	-3.51	-3.50	-3.54	-2.48	-0.52	-0.53	-0.54	-1.05	0.04	0.03	0.04	0.17
	P97.5	0.90	0.92	0.91	0.63	-1.52	-1.43	-1.49	-0.30	-0.19	-0.19	-0.19	-0.12	0.19	0.23	0.19	1.86
	Coverage	0.95	0.91	0.95	0.08	0.94	0.94	0.95	0.41	0.93	0.93	0.94	0.92	0.96	0.94	0.96	0.72
3	Median	0.75	0.73	0.75	0.56	-2.51	-2.43	-2.51	-1.86	-0.30	-0.30	-0.30	-0.30	0.10	0.11	0.10	0.24
	Mean	0.75	0.73	0.75	0.56	-2.51	-2.43	-2.51	-1.86	-0.31	-0.31	-0.31	-0.31	0.10	0.12	0.10	0.26
	P2.5	0.62	0.56	0.61	0.39	-3.29	-3.26	-3.31	-2.68	-0.44	-0.45	-0.45	-0.52	0.04	0.04	0.04	0.12
	P97.5	0.88	0.89	0.89	0.72	-1.73	-1.59	-1.69	-1.02	-0.22	-0.21	-0.21	-0.19	0.18	0.23	0.18	0.51
	Coverage	0.95	0.88	0.95	0.29	0.94	0.93	0.94	0.62	0.94	0.94	0.94	0.94	0.95	0.92	0.95	0.55
4	Median	0.75	0.74	0.75	0.51	-2.50	-2.48	-2.51	-1.69	-0.30	-0.30	-0.30	-0.30	0.10	0.10	0.10	0.29
	Mean	0.75	0.74	0.75	0.51	-2.51	-2.48	-2.51	-1.69	-0.31	-0.31	-0.31	-0.31	0.10	0.11	0.10	0.32
	P2.5	0.65	0.62	0.64	0.34	-3.17	-3.16	-3.19	-2.44	-0.42	-0.42	-0.43	-0.51	0.06	0.05	0.06	0.16
	P97.5	0.85	0.86	0.86	0.66	-1.85	-1.79	-1.82	-0.93	-0.22	-0.22	-0.22	-0.20	0.15	0.18	0.16	0.64
	Coverage	0.95	0.91	0.95	0.07	0.94	0.94	0.95	0.36	0.94	0.95	0.95	0.94	0.95	0.92	0.95	0.11
5	Median	0.75	0.74	0.75	0.61	-2.50	-2.45	-2.50	-2.03	-0.30	-0.30	-0.30	-0.30	0.10	0.11	0.10	0.19
	Mean	0.75	0.74	0.75	0.61	-2.51	-2.46	-2.51	-2.03	-0.30	-0.30	-0.30	-0.30	0.10	0.11	0.10	0.20
	P2.5	0.66	0.62	0.66	0.49	-3.03	-3.02	-3.04	-2.59	-0.39	-0.39	-0.39	-0.41	0.06	0.06	0.06	0.12
	P97.5	0.84	0.85	0.84	0.72	-1.98	-1.89	-1.96	-1.46	-0.24	-0.24	-0.24	-0.23	0.15	0.18	0.15	0.32
	Coverage	0.95	0.88	0.95	0.25	0.95	0.93	0.95	0.57	0.95	0.95	0.95	0.95	0.95	0.89	0.95	0.29
6	Median	0.75	0.71	0.75	0.46	-2.49	-2.36	-2.50	-1.53	-0.30	-0.30	-0.30	-0.30	0.10	0.12	0.10	0.35
	Mean	0.75	0.71	0.75	0.46	-2.50	-2.35	-2.50	-1.53	-0.33	-0.34	-0.33	-0.25	0.11	0.14	0.11	0.27
	P2.5	0.55	0.47	0.55	0.21	-3.73	-3.64	-3.76	-2.84	-0.65	-0.71	-0.67	-1.27	0.02	0.03	0.02	0.12
	P97.5	0.95	0.94	0.95	0.71	-1.22	-1.00	-1.21	-0.21	-0.17	-0.16	-0.17	-0.10	0.22	0.34	0.23	1.84
	Coverage	0.95	0.90	0.95	0.35	0.95	0.94	0.95	0.68	0.93	0.93	0.93	0.91	0.96	0.97	0.96	0.97
7	Median	0.75	0.74	0.75	0.46	-2.50	-2.45	-2.50	-1.54	-0.30	-0.30	-0.30	-0.30	0.10	0.11	0.10	0.35
	Mean	0.75	0.74	0.75	0.46	-2.50	-2.45	-2.50	-1.53	-0.31	-0.31	-0.31	-0.31	0.10	0.11	0.10	0.35
	P2.5	0.64	0.61	0.64	0.30	-3.23	-3.21	-3.24	-2.31	-0.45	-0.45	-0.45	-0.61	0.05	0.05	0.05	0.19
	P97.5	0.86	0.87	0.87	0.62	-1.77	-1.70	-1.76	-0.74	-0.22	-0.21	-0.21	-0.18	0.16	0.18	0.16	0.82
	Coverage	0.95	0.91	0.95	0.03	0.95	0.94	0.95	0.29	0.95	0.94	0.95	0.93	0.96	0.93	0.95	0.16

about factor loadings (in the case of M2, which relies only on PCA), accounting for common factors removes most of the bias and leads to much improved coverage estimates. When we know the factor loadings (in case of M1), there is no bias and the coverage estimates are accurate. In addition, combining the PCA estimate and the estimate using the known loadings results in almost the same accuracy as M1. This simulation illustrates the importance of accounting for factors in using GIV when loadings correlate with size.

8 Conclusion

We developed granular instrumental variables (GIVs): we remark that idiosyncratic shocks offer a rich source of instruments, and we lay out econometric procedures to optimally extract them from aggregate shocks.

We provided two empirical applications. We plan to put on our web page a series of GIVs, and their control shocks η_t 's. They might be useful for empirical work.

Many more applications seem within reach — the introduction listed some. We hope that GIVs will help identifications in new settings and help researchers investigate and understand causal relationships in the economy.

References

- Acemoglu, Daron, Vasco M Carvalho, Asuman Ozdaglar, and Alireza Tahbaz-Salehi**, “The network origins of aggregate fluctuations,” *Econometrica*, 2012, *80* (5), 1977–2016.
- Adao, Rodrigo, Michal Kolesár, and Eduardo Morales**, “Shift-Share Designs: Theory and Inference,” *NBER Working Paper No. 24944*, 2018.
- Altavilla, Carlo, Marco Pagano, and Saverio Simonelli**, “Bank exposures and sovereign stress transmission,” *Review of Finance*, 2017, *21* (6), 2103–2139.
- Amiti, Mary and David E Weinstein**, “How much do idiosyncratic bank shocks affect investment? Evidence from matched bank-firm loan data,” *Journal of Political Economy*, 2018, *126* (2), 525–587.
- , **Oleg Itskhoki, and Jozef Konings**, “International shocks, variable markups, and domestic prices,” *Forthcoming at the Review of Economic Studies*, 2019.
- Andrews, Isaiah, James Stock, and Liyang Sun**, “Weak instruments in iv regression: Theory and practice,” *Annual Review of Economics*, forthcoming.
- Arellano, Manuel and Stephen Bond**, “Some tests of specification for panel data: Monte Carlo evidence and an application to employment equations,” *The review of economic studies*, 1991, *58* (2), 277–297.
- Autor, David H, David Dorn, and Gordon H Hanson**, “The China syndrome: Local labor market effects of import competition in the United States,” *American Economic Review*, 2013, *103* (6), 2121–68.
- Bai, Jushan and Serena Ng**, “Determining the number of factors in approximate factor models,” *Econometrica*, 2002, *70* (1), 191–221.
- and —, “Confidence intervals for diffusion index forecasts and inference for factor-augmented regressions,” *Econometrica*, 2006, *74* (4), 1133–1150.
- Baqaei, David Rezza and Emmanuel Farhi**, “The Macroeconomic Impact of Microeconomic Shocks: Beyond Hulten’s Theorem,” *NBER Working Paper No. 23145*, 2018.
- Bartik, Timothy J**, *Who benefits from state and local economic development policies?*, W.E. Upjohn Institute for Employment Research, 1991.
- Baumeister, Christiane and James D. Hamilton**, “Structural Interpretation of Vector Autoregressions with Incomplete Identification: Revisiting the Role of Oil Supply and Demand Shocks,” *American Economic Review* forthcoming, 2018.

- Berry, Steven, James Levinsohn, and Ariel Pakes**, “Automobile prices in market equilibrium,” *Econometrica*, 1995, pp. 841–890.
- Blasques, Francisco, Siem Jan Koopman, Andre Lucas, and Julia Schaumburg**, “Spillover dynamics for systemic risk measurement using spatial financial time series models,” *Journal of Econometrics*, 2016, *195* (2), 211 – 223.
- Blundell, Richard and Stephen Bond**, “Initial conditions and moment restrictions in dynamic panel data models,” *Journal of econometrics*, 1998, *87* (1), 115–143.
- Borusyak, Kirill, Peter Hull, and Xavier Jaravel**, “Quasi-experimental shift-share designs,” *NBER Working Paper No. 24997*, 2018.
- Bramoullé, Yann, Habiba Djebbari, and Bernard Fortin**, “Identification of peer effects through social networks,” *Journal of Econometrics*, 2009, *150* (1), 41–55.
- Caldara, Dario, Michele Cavallo, and Matteo Iacoviello**, “Oil price elasticities and oil price fluctuations,” *Journal of Monetary Economics*, 2018.
- Carvalho, Vasco Miguel and Basile Grassi**, “Large Firm Dynamics and the Business Cycle,” *American Economic Review*, 2019.
- Chodorow-Reich, Gabriel**, “Geographic Cross-Sectional Fiscal Spending Multipliers: What Have We Learned?,” *NBER Working Paper No. 23577*, 2017.
- Cochrane, John H.**, “The Dog That Did Not Bark: A Defense of Return Predictability,” *The Review of Financial Studies*, 2008, *21* (4), 1533–1575.
- di Giovanni, Julian and Andrei A Levchenko**, “Country size, international trade, and aggregate fluctuations in granular economies,” *Journal of Political Economy*, 2012, *120* (6), 1083–1132.
- , —, and **Isabelle Méjean**, “Firms, destinations, and aggregate fluctuations,” *Econometrica*, 2014, *82* (4), 1303–1340.
- Fama, Eugene F and Kenneth R French**, “Common risk factors in the returns on stocks and bonds,” *Journal of Financial Economics*, 1993, *33* (1), 3–56.
- Fan, Jiangqing, Yuan Liao, and Weichen Wang**, “Projected Principal Component Analysis in Factor Models,” *The Annals of Statistics*, 2016, *44* (1), 219–254.
- Farhi, Emmanuel and Jean Tirole**, “Deadly embrace: Sovereign and financial balance sheets doom loops,” *The Review of Economic Studies*, 2017, *85* (3), 1781–1823.
- Gabaix, Xavier**, “Power Laws in Economics and Finance,” *Annual Review of Economics*, May 2009, *1* (1), 255–294.

- , “The granular origins of aggregate fluctuations,” *Econometrica*, 2011, 79 (3), 733–772.
- Gaubert, Cecile and Oleg Itskhoki**, “Granular comparative advantage,” Technical Report, National Bureau of Economic Research 2018.
- Goldsmith-Pinkham, Paul, Isaac Sorkin, and Henry Swift**, “Bartik Instruments: What, When, Why, and How,” *NBER Working Paper No. 24408*, 2018.
- Graham, Bryan S**, “Identifying social interactions through conditional variance restrictions,” *Econometrica*, 2008, 76 (3), 643–660.
- Hamilton, James D.**, “Causes and Consequences of the Oil Shock of 2007-08,” *Brookings Papers on Economic Activity*, 2009, 40 (1 (Spring)), 215–283.
- He, Zhiguo, Bryan Kelly, and Asaf Manela**, “Intermediary asset pricing: New evidence from many asset classes,” *Journal of Financial Economics*, 2017, 126, 1–35.
- Heckman, James**, “Dummy Endogenous Variables in a Simultaneous Equation System,” *Econometrica*, 1978, 46 (4), 931–59.
- Kelejian, Harry H and Ingmar R Prucha**, “A generalized moments estimator for the autoregressive parameter in a spatial model,” *International Economic Review*, 1999, 40 (2), 509–533.
- Kilian, Lutz**, “Not all oil price shocks are alike: Disentangling demand and supply shocks in the crude oil market,” *American Economic Review*, 2009, 99 (3), 1053–69.
- and **Daniel P Murphy**, “The role of inventories and speculative trading in the global market for crude oil,” *Journal of Applied Econometrics*, 2014, 29 (3), 454–478.
- Kline, Brendan and Elie Tamer**, “Econometric Analysis of Models With Social Interactions,” in Bryan Graham and Aureo De Paula, eds., *The Econometric Analysis of Network Data*, Academic Press, Forthcoming.
- Koijen, Ralph SJ and Motohiro Yogo**, “An equilibrium model of institutional demand and asset prices,” *Journal of Political Economy*, forthcoming.
- Kramarz, Francis, Julien Martin, and Isabelle Mejean**, “Volatility in the small and in the large: the lack of diversification in international trade,” 2016.
- Kuersteiner, Guido M and Ingmar R Prucha**, “Dynamic spatial panel models: Networks, common shocks, and sequential exogeneity,” *arXiv preprint arXiv:1802.01755*, 2018.
- Leary, Mark T and Michael R Roberts**, “Do peer firms affect corporate financial policy?,” *The Journal of Finance*, 2014, 69 (1), 139–178.

- Lee, Youngki, Luís A Nunes Amaral, David Canning, Martin Meyer, and H Eugene Stanley**, “Universal features in the growth dynamics of complex organizations,” *Physical Review Letters*, 1998, *81* (15), 3275.
- Lewbel, Arthur**, “Using heteroscedasticity to identify and estimate mismeasured and endogenous regressor models,” *Journal of Business & Economic Statistics*, 2012, *30* (1), 67–80.
- Long, John B and Charles I Plosser**, “Real business cycles,” *Journal of Political Economy*, 1983, *91* (1), 39–69.
- Manski, Charles F**, “Identification of endogenous social effects: The reflection problem,” *The Review of Economic Studies*, 1993, *60* (3), 531–542.
- Merton, Robert C**, “On the pricing of corporate debt: The risk structure of interest rates,” *The Journal of Finance*, 1974, *29* (2), 449–470.
- Nakamura, Emi and Jon Steinsson**, “Identification in Macroeconomics,” *Journal of Economic Perspectives*, August 2018, *32* (3), 59–86.
- Pasten, Ernesto, Raphael Schoenle, and Michael Weber**, “Price Rigidity and the Origins of Aggregate Fluctuations,” *NBER Working Paper No. 23750*, 2017.
- Ramey, Valerie A**, “Macroeconomic shocks and their propagation,” in “Handbook of Macroeconomics,” Vol. 2, Elsevier, 2016, pp. 71–162.
- Rasmussen, Carl and Christopher Williams**, *Gaussian Processes for Machine Learning*, The MIT Press, 2005.
- Rigobon, Roberto**, “Identification through heteroskedasticity,” *Review of Economics and Statistics*, 2003, *85* (4), 777–792.
- Sarto, Andres Pablo**, “Recovering Macro Elasticities from Regional Data,” 2018.
- Shi, Wei and Lung fei Lee**, “Spatial dynamic panel data models with interactive fixed effects,” *Journal of Econometrics*, 2017, *197* (2), 323–347.
- Stock, James H and Mark W Watson**, “Dynamic factor models, factor-augmented vector autoregressions, and structural vector autoregressions in macroeconomics,” in “Handbook of macroeconomics,” Vol. 2, Elsevier, 2016, pp. 415–525.
- **and** —, “Identification and Estimation of Dynamic Causal Effects in Macroeconomics Using External Instruments,” *The Economic Journal*, 2018, *128* (610), 917–948.

- **and Motohiro Yogo**, “Testing for weak instruments in linear IV regression,” in Donald WK Andrews and James H Stock, eds., *Identification and Inference for Econometric Models: Essays in Honor of Thomas Rothenberg*, Cambridge University Press, 2005.
- Zou, Hui, Trevor Hastie, and Robert Tibshirani**, “Sparse principal component analysis,” *Journal of Computational and Graphical Statistics*, 2006, 15 (2), 265–286.

9 Appendix: Proofs omitted in the paper

Variance facts We will repeatedly use the fact that if $(u_i)_{i=1\dots N}$ is a series of uncorrelated random variables with mean 0 and common variance σ_u^2 , then

$$\mathbb{E}[u_\Gamma u_E] = 0, \quad (75)$$

and

$$\mathbb{E}[u_\Gamma^2] = \mathbb{E}[u_S u_\Gamma] = h^2 \sigma_u^2. \quad (76)$$

Hence, the standard deviation of the granular residual $u_{\Gamma t}$ is proportional to the Herfindahl. In the general heteroskedastic case, the quasi-equal weight vector is

$$\tilde{E} = \frac{(V^u)^{-1} \iota}{\iota (V^u)^{-1} \iota}$$

Then, for any Γ such that $\iota' \Gamma = 0$, we have:⁴³

$$\mathbb{E}[u_\Gamma u_{\tilde{E}}] = 0. \quad (77)$$

Proof of Proposition 2 The proof is quite elementary, and uses well-known ingredients. We have

$$\hat{\alpha}_T - \alpha = \frac{\mathbb{E}_T[(\alpha y_{St} + \varepsilon_t) u_{\Gamma t}]}{\mathbb{E}_T[y_{St} u_{\Gamma t}]} - \alpha = \frac{\mathbb{E}_T[\varepsilon_t u_{\Gamma t}]}{\mathbb{E}_T[y_{St} u_{\Gamma t}]} = \frac{A_T}{D_T}.$$

Next, the law of large number gives:

$$D_T = \mathbb{E}_T[y_{St} u_{\Gamma t}] \xrightarrow{a.s.} D,$$

with

$$D = \mathbb{E}[y_{St} u_{\Gamma t}] = \mathbb{E}[(\eta_t + u_{St}) u_{\Gamma t}] = \mathbb{E}[u_{St} u_{\Gamma t}] = \mathbb{E}[(u_{\Gamma t} + u_{Et}) u_{\Gamma t}] = \mathbb{E}[u_{\Gamma t}^2] = \sigma_{u_\Gamma}^2.$$

For the numerator, the central limit theorem gives the convergence in distribution:

$$\sqrt{T} A_T \xrightarrow{d} \mathcal{N}(0, \sigma_A^2),$$

⁴³Here is the proof. We have $\tilde{E} = k (V^u)^{-1} \iota$ for $k = \frac{1}{\iota (V^u)^{-1} \iota}$. So

$$\mathbb{E}[u_\Gamma u_{\tilde{E}}] = \mathbb{E}[(\tilde{E}' u)(u' \Gamma)] = \tilde{E}' \mathbb{E}[u u'] \Gamma = \tilde{E}' V^u \Gamma = k \iota' (V^u)^{-1} V^u \Gamma = k \iota' \Gamma = 0$$

as $\iota' \Gamma = 0$.

where

$$\sigma_A^2 = \mathbb{E} [\varepsilon_t^2 u_{\Gamma t}^2] = \mathbb{E} [\varepsilon_t^2] \mathbb{E} [u_{\Gamma t}^2] = \sigma_\varepsilon^2 \sigma_{u_\Gamma}^2,$$

so that

$$\frac{\sigma_A}{D} = \frac{\sigma_\varepsilon \sigma_{u_\Gamma}}{\sigma_{u_\Gamma}^2} = \frac{\sigma_\varepsilon}{\sigma_{u_\Gamma}} =: \sigma_\alpha.$$

Hence,

$$\sqrt{T} (\hat{\alpha}_T - \alpha) \rightarrow^d N(0, \sigma_\alpha^2).$$

Then the u'_{it} s are i.i.d. across i 's, then $\sigma_{u_\Gamma} = h\sigma_u$, see (76).

Proof of Proposition 3 We have

$$\hat{\alpha}_T - \alpha = \frac{\mathbb{E}_T [(\alpha y_{St} + \varepsilon_t) z_t]}{\mathbb{E}_T [y_{St} z_t]} - \alpha = \frac{\mathbb{E}_T [\varepsilon_t z_t]}{\mathbb{E}_T [y_{St} z_t]},$$

so the same proof as for Proposition 2 yields the asymptotic error

$$\sigma_\alpha(\Gamma) = \frac{\sigma_\varepsilon \sigma_z}{|\mathbb{E} [y_{St} z_t]|} = \frac{\sigma_\varepsilon \sigma_z}{|\mathbb{E} [u_{St} z_t]|} = \frac{\sigma_\varepsilon}{\sigma_{y_S} |\text{corr}(u_{St}, z_t)|}.$$

So, the best estimator $z_t = u_{\Gamma t}$ maximizes the square correlation $C(\Gamma) := \text{corr}(u_{St}, u_{\Gamma t})^2$:

$$\max_{\Gamma} C(\Gamma) \text{ subject to } \iota' \Gamma = 0$$

We next solve this problem.

Call V the variance covariance matrix of the u_i . We have:

$$C^2 \text{var}(u_{St}) = \frac{\mathbb{E} [u_{St} u_{\Gamma t}]^2}{\text{var}(u_{\Gamma t})} = \frac{(S' V \Gamma)^2}{\Gamma' V \Gamma}.$$

The problem is invariant to changing Γ into $\lambda \Gamma$ for a non-zero λ . So, we can fix say $S' V \Gamma$ at some value. Given this, we want the minimum value of $\Gamma' V \Gamma$. So, we minimize over Γ the Lagrangian

$$\mathcal{L} = \frac{1}{2} \Gamma' V \Gamma - b \Gamma' \iota - c \Gamma' V S \tag{78}$$

with some Lagrange multipliers b, c . The first order condition in Γ' is: $0 = V \Gamma - b \iota - c V S$, i.e.

$$\Gamma = c S + b V^{-1} \iota.$$

Now, using $\iota' \Gamma = 0$ gives $0 = c + b \iota' V^{-1} \iota$, i.e., with $\tilde{E} := \frac{V^{-1} \iota}{\iota' V^{-1} \iota}$,

$$\Gamma = c (S - \tilde{E}).$$

The factor c doesn't affect the results, (as Γ and $c'\Gamma$ gives the same estimator $\hat{\alpha}_T$), so we may choose $c = 1$.

Proof of Proposition 4: Sketch The full proof is in the online appendix (Section 13.1). Here we provide a proof sketch. For simplicity, we omit the controls C_t . We use a projection matrix Q (defined in (23) with $W = (V^u)^{-1}$) satisfying:

$$Q\Lambda = 0, \quad (79)$$

$$QV^u(I - Q') = 0. \quad (80)$$

Then, premultiplying (53) by Q , we have $Qy_t = Qu_t$. Next, we define $\Gamma := Q'S$, and the GIV as the scalar:

$$z_t := \Gamma'y_t, \quad (81)$$

i.e. $z_t = S'Qy_t = S'Qu_t = \Gamma'u_t$, i.e.

$$z_t = u_{\Gamma t}. \quad (82)$$

Assumption 1 ensures $\Gamma \neq 0$. Assumption 2 ensures that V^u can be recovered from the knowledge of Qu_t .

Recall that we have (55),

$$y_{St} = Mu_{St} + \varepsilon_t$$

for ε_t a shock correlated with the η_t but not with the u_{it} 's. Finally, we have

$$u_{St} = S'u_t = S'Qu_t + S'(I - Q)u_t = z_t + v_t,$$

with $v_t = S'(I - Q)u_t$. Now, (80) ensures $\mathbb{E}[z_tv_t] = 0$. Then, we can write:

$$y_{St} = Mz_t + \varepsilon_t^{ys},$$

with $\varepsilon_t^{ys} := Mv_t + \varepsilon_t$ orthogonal to z_t . Hence, we can estimate the multiplier M by OLS.

Likewise, we have (via (55) and (57))

$$F_t^f = \alpha^f M z_t + \varepsilon_t^f$$

for some shock ε_t^f orthogonal to z_t . Hence, we can estimate $\alpha^f M$ by regressing F_t^f on z_t .

In both regressions, we can add the estimated common shocks η_t^e as controls, which improves the precision. The full proof shows that η_t^e is orthogonal to z_t , so those controls still lead to a consistent estimators of M and $\alpha^f M$.

10 Appendix: Complements and extensions

10.1 The model with heterogeneous loadings on endogenous and exogenous factors

In the main text, we assume for simplicity a homogeneous or parametric sensitivity on endogenous factors, e.g. on the price in the simple supply and demand example. We show how our framework generalizes easily.

The model of section [3.1](#) implies the representation:

$$y_{it} = \theta_i u_{St} + \lambda_i \eta_t + u_{it}, \quad (83)$$

where $\theta_i = \frac{\sum_f \lambda_i^f \alpha^f}{1 - \sum_f \lambda_S^f \alpha^f}$ is the sensitivity to endogenous factors, η_t is a vector of exogenous factors, and λ_i is a vector of factor loadings, both r -dimensional. The new difficulty is to estimate a heterogeneous sets of θ_i — in our more basic case we considered the case of a common θ_i .

We first motivate the procedure before stating it. Assume first that we somehow could estimate the $v_{it} = u_{it} - u_{Et}$ (we will soon take care of the fact that we can only hope to estimate $u_{it} - u_{Et}$). Then, we can form

$$z_t = v_{St}, \quad z_{it} = z_t - \Gamma_i^u v_{it}, \quad (84)$$

where $\Gamma_i^u := \frac{\mathbb{E}[v_{it} z_t]}{\mathbb{E}[v_{it}^2]}$, which is equal to $\frac{\Gamma_i}{1 - E_i}$ when the u_{it} are uncorrelated. This implies that $\mathbb{E}[v_{it} z_{it}] = 0$. So z_{it} is like the traditional GIV, but is made of the idiosyncratic shocks of the factors other than i : it is uncorrelated with v_{it} .

Then, given

$$y_{it} = \theta_i z_{it} + \lambda_i \eta_t + (1 + \theta_i \Gamma_i^u) u_{it}, \quad (85)$$

we can estimate θ_i by OLS of y_{it} on z_{it} . Then, the residual u_{it} can be estimated.

This suggests the following iterative procedure. We call the set of parameters to be estimated $\omega = (\eta_t, \lambda_i, \theta_i)_{i,t}$. Call $n \geq 1$ the round of iterations. Given a guess for $z_{it}^{(n-1)}$, we define:

$$w_{it}^e(\omega) := y_{it} - \theta_i z_{it}^{(n-1)} - \lambda_i \eta_t,$$

and perform [44](#)

$$\min_{\omega} \sum_t \sum_i w_{it}^e(\omega)^2. \quad (86)$$

This involves both a PCA for λ_i , η_t and a plain OLS for θ_i . [45](#) We call $\omega^{(n)}$ one the minimizers (this

⁴⁴It's actually more efficient to do a weighted least square version, i.e. to do: $\min_{\omega} \sum_t v_t^e(\omega)' V^u v_t^e(\omega)$, where $v_{it}^e(\omega) = \frac{w_{it}^e(\omega)}{1 + \theta_i^{(n-1)} \Gamma_i^u}$

⁴⁵In the PCA-OLS step, there's a danger that an estimate η_t will soak up variation due to z_{it} . So perhaps it's best to first perform the OLS step to estimate θ_i , and only then do the PCA step.

involves the usual choice of rotations and scaling for the PCA).

Then, we define

$$v_{it}^{(n)} = \frac{w_{it}^e(\omega^{(n)})}{1 + \theta_i^{(n)} \Gamma_i^u}, \quad (87)$$

which will be the proxy for our u_{it} , and we form⁴⁶

$$z_t^{(n)} := v_{St}^{(n)}, \quad z_{it}^{(n)} := z_t^{(n)} - \Gamma_i^u v_{it}^{(n)}, \quad (88)$$

and $\Gamma_i^u = \frac{\mathbb{E}[v_{it} z_t]}{\mathbb{E}[v_{it}^2]}$ can be estimated. If it's not estimated, we may take $\Gamma_i^u = \frac{\Gamma_i}{1-E_i}$.⁴⁷

Then, we iterate, until convergence.

Once the model is estimated, we can get z_t and use it to estimate the sensitivity α^f of the endogenous factors via OLS on z_t , like in the GIV with homogeneous sensitivity to endogenous factors $F_t^f = \alpha^f M z_t + \lambda^f \eta_t + \varepsilon_t^f$.

One can check that in the simplest case of supply-and-demand ($y_{it} = \phi^d p_t + \eta_t + u_{it}$), we recover our basic GIV. Indeed, then the estimate η_t is: $\eta_t^e = \eta_t + u_{Et}$, $v_{it} = u_{it} - u_{Et}$, and $\Gamma_i^u = \frac{\Gamma_i}{1-\frac{1}{N}}$.

Remarks This will give $v_t = Q u_t$ with $Q = I - \lambda(\lambda' \lambda)^{-1} \lambda'$. For instance, in the simplest model with $\lambda = \iota$, we have $v_{it} = u_{it} - u_{Et}$, so that $v_{St} = u_{\Gamma t}$. We recover our basic GIV (which we developed with homogeneous exposures to the common factor).

We can do a parametric version of that, e.g. set $\theta_i = \Theta' x_{it}$ for some vector x_{it} of characteristics. Then, we estimate Θ in the PCA-OLS step.

A potential variant is the following. We do the above, but replace y_{it} by

$$y_{it}^{(n-1)} = y_{it} + \theta_i^{(n-1)} v_{it}^{(n-1)}, \quad (89)$$

and form

$$w_{it}^e(\omega) := y_{it}^{(n-1)} - \theta_i z_{it}^{(n-1)} - \lambda_i \eta_t, \quad (90)$$

and do the PCA + OLS estimation of ω , with $y_{it}^{(n-1)}$ on the left-hand side. We just set

$$v_{it}^{(n)} = w_{it}^{(n)}. \quad (91)$$

The advantage of this procedure is that it doesn't involve the division by $1 + \theta_i^{(n)} \Gamma_i^u$ step in (87), which potentially leads to instabilities.

⁴⁶At the first round, before any estimation, we take $v_{it}^{(0)} = y_{it} - y_{Et}$.

⁴⁷In the case with uncorrelated heteroskedastic shocks and $\lambda = \iota$, one can show that $\Gamma_i^u = \frac{\Gamma_i}{1-E_i}$, so that $\Gamma_i^u = \Gamma_i + O(\frac{1}{N})$. So if we don't estimate it, a good approximation is $\Gamma_i^u = \frac{\Gamma_i}{1-E_i}$ or even $\Gamma_i^u = \Gamma_i$.

10.2 When the influence matrix is not proportional to size

10.2.1 Position of the problem

Suppose a model

$$y_{it} = \gamma \sum_j G_{ij} y_{jt} + \lambda_i \eta_t + u_{it}, \quad (92)$$

i.e.

$$y_t = \gamma G y_t + \Lambda \eta_t + u_t, \quad (93)$$

with a given “influence” matrix G . For instance, if we have an “industrial similarity” matrix H_{ij} , we might set

$$G_{ij} = \frac{H_{ij} S_j}{\sum_k H_{ik} S_k}.$$

In our basic “reflection problem”, $G = \iota S'$.

We’d like to identify γ . With $V = \mathbb{E}[u_t u_t']$, we define $\tilde{E} = \frac{V^{-1} \iota}{\iota' V^{-1} \iota}$, and the “generalized size vector”:

$$S := G' \tilde{E}, \quad (94)$$

which is the analogue of “size” in our simpler setup where $G = \iota S'$.

10.2.2 A simple approach, when the loading on common shocks is uniform

In this subsection we assume that

$$G \iota = \iota, \quad (95)$$

which is satisfied in many examples (Section 12.6 has the general case). Consider some vector \mathcal{E} , and define:

$$z_t := \mathcal{E}' (G - I) y_t. \quad (96)$$

Then, we have the key relation:

$$\mathbb{E} [z_t (y_{\tilde{E}t} - \gamma y_{St})] = 0, \quad (97)$$

which allows to identify γ by $\gamma = \frac{\mathbb{E}[z_t y_{\tilde{E}t}]}{\mathbb{E}[z_t y_{St}]}$.

Relation (97) works for any z_t of the type (96). It is sensible to take $\mathcal{E} = \tilde{E}$ (one can show that this is the optimum choice in the sense of minimizing the asymptotic error). Then, the GIV is again: $z_t = y_{St} - y_{\tilde{E}t}$.

Derivation of (97) Indeed,

$$y_{\tilde{E}t} - \gamma y_{St} = \tilde{E}' (I - \gamma G) y_t = \tilde{E}' (u_t + \eta_t \iota).$$

Given $G\iota = \iota$, and $(G - I)$ and $(I - \gamma G)^{-1}$ commute, we have the useful relation:

$$(G - I) (I - \gamma G)^{-1} \iota = 0. \quad (98)$$

As a result,

$$z_t = \mathcal{E}' (G - I) y_t = \mathcal{E}' (G - I) (I - \gamma G)^{-1} (u_t + \eta_t \iota) = \mathcal{E}' (G - I) (I - \gamma G)^{-1} u_t.$$

Hence,

$$\begin{aligned} a &:= \mathbb{E} [z_t (y_{\tilde{E}t} - \gamma y_{St})] = \mathbb{E} \left[\mathcal{E}' (G - I) (I - \gamma G)^{-1} u_t (u_t + \eta_t \iota)' \tilde{E} \right] \\ &= \mathcal{E}' (G - I) (I - \gamma G)^{-1} V \tilde{E} = \mathcal{E}' (G - I) (I - \gamma G)^{-1} \frac{\iota}{\iota' V^{-1} \iota} = 0. \end{aligned}$$

Online Appendix for “Granular Instrumental Variables”

Xavier Gabaix and Ralph S.J. Koijen

April 8, 2019

This online appendix gives complements to the theory, the underlying models, and the empirical examples. It gives also additional proofs.

11 Microfoundations for the model of doom loops

We provide a microfoundation for the empirical model of doom loops in Section [4](#).

11.1 Model setup

We make a number of simplifying assumptions. The safe interest rate is normalized to 0, and pricing is risk neutral. Time is continuous in $[0, T]$. We neglect the $O(dt)$ terms, which are irrelevant for the regression analysis we are interested in, i.e. will write $df(X_t) = f'(X_t)dX_t$.^{[48](#)[49](#)}

Bank i has debt D_i , fundamental value V_{it} , and equity pays only that terminal payoff at date T (which should be thought about as faraway), as in the [Merton](#) ([1974](#)) model. So the equity value is

$$\mathcal{E}_{it} = \mathbb{E}_t [(V_{iT} - D_i)^+] . \quad (99)$$

The net liability that the bank imposes on the government is:

$$G_{it} = -\tau \mathcal{E}_{iT} + q \mathbb{E}_t [(D_i - V_{iT})^+] , \quad (100)$$

and reflects two channels. The first channel is the “healthy revenues” channel (captured by $-\tau \mathcal{E}_{iT}$): a bank with equity \mathcal{E}_{it} will generate at T a revenue $\tau \mathcal{E}_{iT}$, where τ is the tax rate.^{[50](#)} The second channel is the “bank bailout” channel (captured by $q \mathbb{E}_t [(D_i - V_{iT})^+]$) that if bank i defaults at T ,

⁴⁸Formally, we write all the differential expressions $dY_t = a_t dZ_t$ modulo an equivalence by terms $b_t dt$ (or, to be pedantic, we quotient by the ring of expressions of the type $b_t dt$ where b_t is an adapted function). So, $df(X_t) = f'(X_t)dX_t$ modulo dt , where we keep the “modulo dt ” implicit.

⁴⁹We only care, for the regressions, about the “ dZ_t ” terms, that depends on innovation to underlying Brownian shocks dZ_t , as those are the loading detected by the regressions.

⁵⁰This can be a stand-in for richer variants: e.g. a well-capitalized bank might lend more to investing firms, hence increase GDP, and boost government revenues by that channel too. This a bit more complex to model, and the results would be similar.

the government takes over a fraction q of what bank i still owes the bond holders. So, the expected liability government liability for bank i is: $q\mathbb{E}_t[(D_i - V_{iT})^+]$.⁵¹

At time T , the government itself has debt due. The government's debt is B , and value of the debt (per unit of facial value) is thus:

$$Q_t = \mathbb{E}_t \left[\max \left(1, \frac{K_T - G_T}{B} \right) \right] = e^{-(T-t)y_t}. \quad (101)$$

where K_T is a stochastic “repayment capacity”, $G_t = \sum_i G_{it}$ is the total liability, and y_t is the yield spread on government debt, as we normalize the safe interest rate to 0. We set K_T as a martingale, $K_T = K_0 e^{-\frac{\sigma^2}{2}T + \sigma W_T}$ for a standard Wiener process W_t .

11.2 Model solution

We now report the model solution and comment on its economics (the next section the spells out the derivation). Two important quantities are p_{Gt} , the probability of default of government debt, and p_{it} , the probability of default for bank i :

$$p_{Gt} = \mathbb{P}_t \left[\frac{K_T - G_T}{B} < 1 \right], \quad p_{it} = \mathbb{P}_t \left[\frac{V_{iT}}{D_i} < 1 \right]. \quad (102)$$

Both are the risk-neutral probabilities, i.e. they embed a risk premium. The stock return is

$$dr_{it} = \frac{d\mathcal{E}_{it}}{\mathcal{E}_{it}}. \quad (103)$$

Given the model, a change dr_{it} of bank i affects the bond yield (really, the yield spread over the safe bond) as:

$$dy_t = -a_t^y \frac{p_{Gt}}{B} \sum_i Y_{it} dr_{it} + d\eta_t^y, \quad (104)$$

where $a_t^y = \frac{1}{T-t}$ is a slowly-varying parameter (as T is far from the interval of times t under study), Y_{it} is the following measure of “size” for the bank risk:

$$Y_{it} = \left(\tau + \frac{qp_{it}}{1 - p_{it}} \right) \mathcal{E}_{it}. \quad (105)$$

It features our two channels: the “healthy revenues” channel (τ) and the “bank bailout” channel (proportional to $\frac{p_{it}}{1 - p_{it}}$) the equity value times the probability of default. The last term in [104](#), $d\eta_t^y$,

⁵¹This too may be is a stand-in for richer variants. Likewise, if if bank i is very leveraged and close to a regulatory constraint, it might cut down on lending. The reduced form version will still be, we conjecture, that the “size” factor will be equity value times default probability, $\frac{p_{it}}{1 - p_{it}} \mathcal{E}_{it}$. Likewise, adding a term $-\tau' V_{iT}$ in [\(100\)](#) would change almost nothing, as changes in equity and assets are very close (equation [\(109\)](#)). Finally, the predictions are quite similar if the government pays a fraction of total bank debt in case of default, $q\mathbb{E}_t[1_{D_i - V_{iT} > 0}]$.

is just a sensitivity to common shocks.

Hence, the impact of a return dr_{it} on the yield is proportional to equity \mathcal{E}_{it} times (i) p_{Gt} for the “general revenue” channel and (ii) $p_{Gt}p_{it}$ for the “bank bailout” channel — the product of the probability of default by both the government (p_{Gt}) and by the bank (p_{it}); indeed, if either is 0, then the impact of a bank return on the bond yield should be 0, as no default will occur. So, the size factor should be (in this model): $S_{it} = \frac{Y_{it}}{Y_t^B}$ with $Y_t^B = \sum_i Y_{it}$.

So, the sensitivity of the sovereign yield is

$$dy_t = \alpha_t^y dr_{St} + d\eta_t^y$$

with

$$\alpha_t^y = -a_t^y \frac{Y_t^B}{B} p_{Gt}, \quad (106)$$

The return of a bank i has a sensitivity to the yield equal to:

$$dr_{it} = \lambda_i^y dy_t + du_{it} + \lambda_i^\eta d\eta_t, \quad (107)$$

where

$$\lambda_i^y = L_{it}\theta_{it}, \quad L_{it} = (1 - p_{it}) \frac{V_{it}}{\mathcal{E}_{it}}, \quad (108)$$

where L_{it} is the leverage of the firm, adjusted by a probability of default, and θ_{it} is the exposure of the firm’s total value (debt + equity) to the yield. The term du_{it} is the idiosyncratic shock, and $\lambda_i^\eta d\eta_t$ is a sensitivity to the common shock that we need not theorize. Hence, the multiplier is the usual $M_t = \frac{1}{1 - \alpha_t^y \lambda_{St}^y}$.

The model could be enriched in a number of ways.⁵² Still, the main lesson is likely to be robust to those variants: The main channel is that the “weight” of a bank should be proportional to equity values, for the “healthy revenues” channel, and equity values times probability of default for the “bank bailout” channel.

11.3 Derivation of the model

As we assume that pricing is risk-neutral, the fundamental bank value V_{it} is a martingale – this simplifies the calculations. Equation (99) implies that $d\mathcal{E}_{it} = \mathbb{E}_t[1_{V_{iT} > D_i} dV_{it}]$, i.e.

$$d\mathcal{E}_{it} = (1 - p_{it}) dV_{it} = \mathcal{E}_{it} dr_{it}, \quad (109)$$

⁵²For instance, we could endogenize the risk-free rate. The risk-free rate (e.g., German bond yield) might fall when the disaster probability increases, which would generate a “reverse” doom loop where a negative return on German banks create an decrease in the German yield. This would take us too far into modelling for the purposes of this paper, though. The “yield” is best understood as the “yield spread”, i.e. the spread over the yield of the safest bond, the German yield.

where $p_{it} = \mathbb{E}_t [1_{V_{iT} > D_i}]$ is the probability of default.⁵³ Hence, (100) gives:

$$dG_{it} = -\tau d\mathcal{E}_{it} - qp_{it}dV_{it} = -\left(\tau + q\frac{p_{it}}{1-p_{it}}\right)d\mathcal{E}_{it} = -\left(\tau + q\frac{p_{it}}{1-p_{it}}\right)\mathcal{E}_{it}dr_{it} = -qY_{it}dr_{it}, \quad (110)$$

so the right “size” is

$$Y_{it} = \left(\tau + q\frac{p_{it}}{1-p_{it}}\right)\mathcal{E}_{it}. \quad (111)$$

Next (101) gives:

$$dQ_t = -\mathbb{E}_T [1_{B > K_T - G_T}] (dG_T - dK_T) = -\frac{p_{Gt}}{B} (dG_T - dK_T),$$

so (omitting the dK_t as they will be absorbed in the common shock $d\eta_t^y$)

$$dy_t = -\frac{dQ_t}{T-t} = \frac{p_{Gt}}{(T-t)B}dG_t = -\frac{1}{(T-t)B}p_{Gt} \sum_i Y_{it}dr_{it} = \frac{1}{(T-t)B}p_{Gt} \sum_i Y_{it}dr_{it}, \quad (112)$$

i.e. the expression announced in (104), with

$$a_t^y = \frac{1}{T-t}. \quad (113)$$

Finally, bank i ’s underlying value V_{it} evolves as:

$$\frac{dV_{it}}{V_{it}} = \theta_{it}dy_t + \Theta_{it}d\eta_t + du_{it}, \quad (114)$$

This is, it has an exposure θ_{it} to government bond yields, exposure Θ_{it} to another factor η_t , and some idiosyncratic return du_{it} . This gives, via (109)

$$dr_{it} = \frac{d\mathcal{E}_{it}}{\mathcal{E}_{it}} = (1-p_{it})\frac{dV_{it}}{\mathcal{E}_{it}} = (1-p_{it})\frac{V_{it}}{\mathcal{E}_{it}}\frac{dV_{it}}{V_{it}} = (1-p_{it})\frac{V_{it}}{\mathcal{E}_{it}}\theta_{it}dy_t,$$

as announced in (108).

11.4 Discussion

Another variant with a “propagation via GDP” channel Now, consider the variant with just one size channel (to simplify), but with a new term μr_{St} :

$$\begin{aligned} y_t &= \alpha r_{St} + \eta_t^y, \\ r_{it} &= \lambda y_t + \mu r_{St} + u_{it} + \eta_t^r. \end{aligned}$$

⁵³It depends on leverage, so that one could write is as $p_{it} = f(V_{it}/D_i)$.

The new μr_{St} term means that other banks can affect bank i via other channels than through y_t , e.g. as they lend to other firms, that increases output and bank i . The full multiplier is now

$$M = \frac{1}{1 - \alpha\lambda - \mu} \quad (115)$$

and we still have (omitting all common shocks η terms):

$$y_t = \alpha M u_{St}, \quad r_{St} = M u_{St}. \quad (116)$$

We conclude that another general GDP channel changes the structural interpretation of M , but doesn't change our interpretation of α and M .

An objection and answer A potential objection is the following. Suppose that bank i has a disproportionate exposure to say a large sector s of the economy. The causality might be just from sector s to GDP and the government yield, rather than via bank i . Then, a negative return for bank i would just reflect a negative shock to sector s . Bank i would just be a symptom, not a causal factor.

There are a few ways to address that. First (and this is probably the simplest solution), we can control for sector-specific stock market returns. Even though the aggregate stock market might have extraneous variations (e.g. due to very long run forecast, or variations in risk premia), the sector-specific stock market return (controlling for the aggregate) will control for “fundamental news about the economy”. Second, when we interact $\Delta y_t = \gamma (\sum_i p_{it} \mathcal{E}_{it} r_{it})$, where p_{it} is bank i 's CDS, this captures the “pure bank bailout” channel. Third, we can examine what caused (narratively) the main return shocks to the big banks.

A potential extension to estimate the model in the presence of two channels In the augmented model with the two channels, a return r_{it} changes in the bond yield as in (104), where the augmented size is (105)

$$Y_{it} = \left(\tau + q \frac{p_{it}}{1 - p_{it}} \right) \mathcal{E}_{it} \quad (117)$$

Now, we have a new richness that we haven't encountered so far: we have *two* size vectors,

$$Y_{it}^1 = \mathcal{E}_{it}, \quad Y_{it}^2 = \frac{p_{it}}{1 - p_{it}} \mathcal{E}_{it}. \quad (118)$$

Call S^k ($k = 1, 2$) the associated normalized sizes (summing to 1, $S_i^k = \frac{Y_i^k}{\sum_j Y_j^k}$). Then, the model is (omitting the common factor η_t for notational simplicity), in summary and calling y_t the change in

the yield,

$$\begin{aligned} y_t &= \sum_{k=1}^2 \alpha^k r_{S^k t}, \\ r_{it} &= \lambda y_t + u_{it}, \end{aligned}$$

where $(\alpha^1, \alpha^2) = \frac{q}{(T-t)B} p_{Gt}(q, \tau)$ will be estimated. We have $y_t = \sum_{k=1}^2 \alpha^k \lambda y_t + \sum_k \alpha^k u_{S^k t}$, so with $M = \frac{1}{1 - \sum_k \alpha^k \lambda}$, i.e.

$$y_t = \sum_k M \alpha^k u_{S^k t} + b' \eta_t.$$

Hence, we can run the following OLS, which gives the values of $\beta^k = M \alpha^k$:

$$y_t = \sum_k M \alpha^k z_t^k + b' \eta_t, \quad (119)$$

where $z_t^k = u_{S^k - E, t}$. We have two GIVs, the z_t^k with $k = 1, 2$.

The total impact of a shock is then

$$\beta = \sum_{k=1}^2 \beta^k, \quad (120)$$

and this is what counts for the total impact. Then, it gives a “first stage”

$$y_t^1 = \sum_k \beta^k z_t^k.$$

Then, to estimate λ , we run e.g.

$$r_{Et} = \lambda y_t^1 = \lambda \left(\sum_k \beta^k z_t^k \right),$$

and this gives an estimate of λ .

Pragmatically, it's probably a good idea to first explore each channel separately. Perhaps one is much bigger than the other, and that's the one we'll insist on.

12 Complements

12.1 Multi-dimensional actions

Suppose now that the action y_{it} is q -dimensional, for some $q \geq 1$. For instance, y_{it} 's components might be the growth rate, and the labor share of firms of firm i , and then $q = 2$. Then, the general GIV procedure extends well, as we shall now see.

We call $a \in \{1, \dots, q\}$ (as in action) a component of y . The model is:

$$\begin{aligned} y_{S^a t} &= \sum_f \lambda_{S^a, f}^a F^f + u_{S^a t}^a, \\ F_t^f &= \eta_t^f + \sum_a \alpha_a^f y_{S^a, t}^a, \end{aligned}$$

Here u_{it} is q dimensional, α is $q \times r$ dimensional matrix, and λ is $r \times q$ dimensional matrix.

We can also estimate M (hence $\sum_f \alpha^f \lambda^f$), the α^f . For ε_t a composite of aggregate shocks,

$$y_{St} = Hy_{St} + u_{St} + \varepsilon_t,$$

where

$$H = \Lambda A = \sum_f \alpha^f \lambda^f,$$

with $\Lambda_{af} = \lambda_{S^a, f}^a$ and $A_{fa} = \alpha_a^f$ matrices with dimensions $q \times r$ and $r \times q$ respectively, so that H is $q \times q$, and

$$u_{St} = (u_{S^a t}^a)_{a=1 \dots q}.$$

This implies

$$y_{St} = M(u_{St} + \varepsilon_t), \tag{121}$$

there the multiplier M is now a $q \times q$ matrix:

$$M = (I - H)^{-1}.$$

We will form a GIV:

$$z_t = u_{\Gamma t},$$

which is q -dimensional:

$$u_{\Gamma} = (u_{\Gamma^a}^a)_{a=1 \dots q}.$$

We want, with $E^a = S^a - \Gamma^a$,

$$\mathbb{E}[u_{Et} u'_{\Gamma t}] = 0$$

i.e., for all $Q^{ab} = 0$, where

$$Q^{ab} := \mathbb{E}[u_{E^a t}^a u_{\Gamma^b t}^b].$$

Let us focus on the case where u_{it}, u_{jt} are uncorrelated for $i \neq j$, but for a given i , u_{it}^a, u_{it}^b can be correlated (if a firm have a investment boom, it will likely hire more labor, so that the components of its idiosyncratic shock in $y_{it} \in \mathbb{R}^q$ will be correlated).

We have:

$$Q^{ab} = \sum_i E_i^a \Gamma_i^b v_i^{ab}, \quad v_i^{ab} := \mathbb{E} [u_{it}^a u_{it}^b]. \quad (122)$$

For simplicity, we will suppose that there are v^{ab} and σ_i^2 such that

$$v_i^{ab} = \sigma_i^2 v^{ab}. \quad (123)$$

Hence, we can simply take $E_i = \frac{k}{\sigma_i^2}$ with $k = \frac{1}{\sum_j 1/\sigma_j^2}$ and set, for all a , $E_i^a = E_i$ and $\Gamma^a = S^a - E^a$. Then,

$$Q^{ab} = \sum_i \frac{k}{\sigma_i^2} \Gamma_i^b \sigma_i^2 v^{ab} = k v^{ab} \sum_i \Gamma_i^b = 0,$$

so that we have achieved our goal that $\mathbb{E} [u_{Et} u'_{\Gamma t}] = 0$. In the more general case, other Γ_i^a can probably be found.

Given (121), we have

$$y_{St} = M (u_{St} + \varepsilon_t) = M (u_{\Gamma t} + u_{Et} + \varepsilon_t),$$

so

$$\mathbb{E} [y_{St} z'_t] = M \mathbb{E} [z_t z'_t],$$

hence our estimator is

$$M = \mathbb{E} [y_{St} z'_t] \mathbb{E} [z_t z'_t]^{-1}. \quad (124)$$

Finally, we can also estimate $\alpha^f M$ by regressing on z_t :

$$F_t^f = \eta_t^f + \sum_a \alpha_a^f y_{Sa,t}^a = \eta_t^f + \alpha^f y_{St} = \eta_t^f + \alpha^f M (u_{\Gamma t} + u_{Et} + \varepsilon_t) m,$$

so $\beta^f = \alpha^f M$ (a row vector) obtains by simply regressing

$$F_t^f = \beta^f z_t + \varepsilon_t^f,$$

and get $\beta^f = \alpha^f M$, $\beta^f = \mathbb{E} [F_t^f z'_t] \mathbb{E} [z_t z'_t]^{-1}$.

Extension: causal estimation of the actor-specific multiplier The following is a refinement.

We can also identify causally $\mu_i := \lambda_i \alpha = \sum_f \lambda_i^f \alpha^f$. Indeed, use

$$u_{\Gamma t, -i} := u_{\Gamma t} - S_i^u u_{it}, \quad (125)$$

which is the granular shock purged of a correlation with u_{it} . Then, a shock u_{st} creates an impact $\frac{dF_t}{du_{st}} = M\alpha$, hence an impact

$$\frac{dy_{it}}{du_{st}} = M\lambda_i\alpha.$$

Hence, we can identify μ_i , by regression

$$y_{it} = \mu_i M u_{\Gamma t, -i} + \phi^i \mathcal{C}_t + \varepsilon_{it}^y, \quad (126)$$

with some noise ε_{it}^y . This is the average impact of a causal impact of idiosyncratic shocks of the other entities on entity i .

12.2 Full recovery when different factors have different “size” weights

In the basic model, we can identify α^f , $M = \frac{1}{1 - \sum_f \lambda^f \alpha^f}$, but not λ^f .

We give some conditions under which we can actually also identify the λ^f (in addition to α^f and M). We show here that this is the case if assume that the size S^f differ across all factors f , and this knowledge is given to us (from a model).

Here we take the basic set up as in Section 3.1, in the simplified case where $\lambda_i^f = \lambda^f$ for all “endogenous” factors, i.e. for the factors f such that $\alpha^f \neq 0$, the other exogenous factors η all have an impact of 1:

$$y_{it} = u_{it} + \sum_f \lambda^f F_t^f + \eta_t^y, \quad (127)$$

$$F_t^f = \alpha^f y_{S^f, t} + \eta_t^f. \quad (128)$$

This implies

$$y_t = u_t + \iota \sum_f \lambda^f F_t^f + \iota \eta_t^y = u_t + \iota \sum_f \lambda^f \left(\eta_t^f + \alpha^f S^{f'} y_{it} \right) + \iota \eta_t^y.$$

Noting “ ε^k ” some combination of the various η ’s, and as usual $M = \frac{1}{1 - \sum_f \alpha^f \lambda^f}$,

$$\begin{aligned} y_t &= \left(I - \iota \sum_f \lambda^f \alpha^f S^{f'} \right)^{-1} (u_t + \iota \varepsilon_t^1) \\ &= \left(I + M \iota \sum_f \lambda^f \alpha^f S^{f'} \right) (u_t + \iota \varepsilon_t^1) \\ y_t &= u_t + M \iota \sum_f \lambda^f \alpha^f u_{S^f, t} + \iota \varepsilon_t^y, \end{aligned} \quad (129)$$

i.e. so that $F_t^f = \eta_t^f + \alpha^f y_{S^f,t}$ gives:

$$F_t^f = \alpha^f \left(u_{S^f,t} + M \sum_g \lambda^g \alpha^g u_{S^g,t} \right) + \varepsilon_t^f. \quad (130)$$

Hence, suppose that we extracted the $\check{u}_{it} = u_{it} - u_{Et}$ (following our usual procedure). Then, we form

$$z_{\Gamma^f t} := S^{f'} \check{u}_t = u_{S^f t} - u_{Et}. \quad (131)$$

Then, regressing F_t^f on the various $z_{\Gamma^g t}$

$$F_t^f = \sum_g b_g^f z_{\Gamma^g t} + \varepsilon_t^{f,1} \quad (132)$$

(for $\varepsilon^{f,1}$ some residual noise) yields a regression coefficient:

$$b_g^f = \alpha^f (1_{f=g} + M \lambda^g \alpha^g). \quad (133)$$

This allows to recover everything, and with several overidentifying restrictions. Indeed,

$$b^f := \sum_g b_g^f = \alpha^f \left(1 + M \sum_g \lambda^g \alpha^g \right) = \alpha^f M,$$

which identifies $\alpha^f M$. Next, for $f \neq g$,

$$\frac{b_g^f}{b^f} = \lambda^g \alpha^g,$$

which gives $\lambda^g \alpha^g$ (and should be equal for all f), hence M . Hence, we obtained $\alpha^f M$, M and $\lambda^g \alpha^g$ — hence all quantities: α^f, λ^f, M .

12.3 Complements to the general procedure

The procedure can be simplified in some cases. *When we have a long time-series.* Recall that

$$y_{St} = \sum_f \lambda_{St}^f F_t^f + u_{St}. \quad (134)$$

Hence, if all factors with λ_{St}^f possibly non-zero are observables and exogenous, we can measure the λ_{St}^f by OLS with the regression above, and get u_{St} to be the residual. This is useful when we have high-frequency data (e.g. daily financial returns), which can give an acceptably small error.⁵⁴

⁵⁴Indeed, this time-series regressions gives an $O\left(\frac{1}{\sqrt{T}}\right)$ error, which is good enough for large T . Using the cross section, as in the basic procedure, gives an $O\left(\frac{1}{\sqrt{TN}}\right)$ error.

We can aggregate entities into categories. For this discussion, replaced “entity” by “firm”. We could aggregate the firms into $K > 1$ sub-categories (e.g. industries – or even an arbitrary categorization like “blue firms” and “red firms”)— then the above works, but interpreting firm i as “aggregate firm category i ”. Indeed, (48) aggregates without problem: if aggregate k is made of firm $i \in I_k$, we just define the aggregate size of category k as $S_{[k]} := \sum_{i \in I_k} S_i$, the relative weight of firm i in category k as $\omega_{[k]i} = \frac{S_i 1_{i \in I_k}}{S_{[k]}}$, and the action factor loading as value-weighted averages ($y_{[k],t} = \sum_i \omega_{[k]i,t} r_{it}$, $\alpha_{[k]}^f = \sum_i \omega_{[k]i} \alpha_i^f$). Then, the model works, using those aggregated categories. What we do need is that categories has non-trivial idiosyncratic shock (so that a “very small firms” category would not be valid, as it would have $\text{var}(u_{it}) \simeq 0$).

12.4 Typical size of Herfindahls

The GIV instrument is valid as long as $h > 0$, i.e. as long as there is heterogeneity. However, for it to be strong, we need high Herfindahls. In estimates for firms, we typically have $h \in [0.02, 0.5]$. One can have an a priori estimate of its size (for theory purposes). In practice, many size distributions follow a power law with fat tails, $\mathbb{P}(S > x) \sim kx^{-\zeta}$ for large x , with $\zeta \in (1, 2]$ — something also explained via random growth behavior. In that case one can show that (as in Gabaix (2011), Proposition 2)

$$h \sim k' N^{-\psi}, \quad \psi = 1 - \frac{1}{\zeta} \in (0, \frac{1}{2}] \quad (135)$$

for k' a non-zero random variable independent of N . The traditional variance case gives corresponds to $\zeta = 2$, which confirms $h \sim k' N^{-1/2}$ (and then k' is a constant), a very weak instrument. But when $\zeta \in (1, 2)$, we have a decay in $N^{-\psi}$ with $\psi \in (0, \frac{1}{2})$. A fatter tail in the distribution of large firms (lower ζ) creates scaling in N that decays more slowly (ψ is lower) as N grows large. In the limit of Zipf’s law (i.e., $\zeta \rightarrow 1$), we find $\psi \rightarrow 0$ (indeed, one can show that we have $h \sim \frac{k'}{\log N}$), a stronger instrument.

To simulate sizes from a power law distribution with exponent ζ , we can take $V_i = i^{-1/\zeta}$, and $S_i = \frac{V_i}{\sum_j V_j}$.⁵⁵ In the case of Zipf’s law, that yields $h \sim \frac{\alpha}{\log n}$ with $\alpha = \frac{\pi}{\sqrt{6}} \simeq 1.3$.

12.5 When we have disaggregated data for both the demand and the supply side

When we have disaggregated data for both the demand and the supply side, we can refine the “exclusion restriction”. So far we assumed that $\mathbb{E}[u_{it}\varepsilon_t] = 0$, i.e. no covariance between idiosyncratic demand and supply shock. If that’s not the case, we can also decompose each supply with a factor model:

$$s_{it} = \mu_i^p p_t + \sum_f \lambda_i^{s,p} \eta_t^{s,f} + u_{it}^s. \quad (136)$$

⁵⁵More refined, we can simulate n i.i.d. uniform variables U_i , order them $U_{(1)} \leq \dots \leq U_{(n)}$, and take $V_i = U_{(i)}^{-1/\zeta}$.

Then, if the US has a “fracking shock” that affects both supply and demand, it will be captured by both u_{it}^s and u_{it}^d for $i = \text{USA}$, but we still maintain $\mathbb{E}[u_{it}^x \eta_t^{x'}] = 0$ for x, x' in s, d . But now we keep $\mathbb{E}[u_{it}^s u_{it}^d]$ potentially nonzero.

Then, we have two instruments: the old one, i.e. the demand-based GIV, $z_t^d = u_{\Gamma t}^d$ (replacing putting a d on the demand side), the new supply-based GIV:

$$z_t^s := \Gamma^{s'} y_t^s = u_{\Gamma^{s'} t}^s, \quad (137)$$

where $\Gamma^s = S^s - E^s$ is the relative size on the supply side. We can also form the difference:

$$z_t^{d-s} = z_t^d - z_t^s. \quad (138)$$

Then, the “demand minus supply” GIV is always valid, e.g. it satisfies $\mathbb{E}[(y_{Et} - \lambda p_t) z_t^{d-s}] = 0$. If supply and demand shocks are uncorrelated across countries, then both z_t^d and z_t^s are valid instruments.

12.6 When the influence matrix is not proportional to size: When the loading on common shocks is not necessarily uniform

Here we complete our discussion in Section [10.2](#). We now study the more general case where:

$$y_t = \gamma G y_t + \Lambda \eta_t + u_t, \quad (139)$$

where the factor loading Λ (an $N \times r$ matrix) is not necessarily equal to ι (but we keep imposing that the Λ spans ι , i.e. there is a q such that $\iota = \Lambda q$). As before, η_t is a low-dimensional vector of factors. However, we do not assume anymore that $G\iota = \iota$.

First, we suppose that we have a first estimate of γ , which we call if γ^e . We will later iterate on it. Then, we form:

$$\tilde{y}_t(\gamma^e) = (I - \gamma^e G) y_t. \quad (140)$$

If $\gamma^e = \gamma$, then

$$\tilde{y}_t(\gamma) = \Lambda \eta_t + u_t. \quad (141)$$

Hence, we run a factor analysis on $\tilde{y}_t(\gamma^e)$, which recovers Λ and $W = (V^u)^{-1}$. We introduce Q as in [\(159\)](#) so that $Q\Lambda = 0$ and set

$$\check{u}^e = Q \tilde{y}_t(\gamma^e),$$

so that at $\gamma^e = \gamma$,

$$\check{u}_t^e = Q u_t. \quad (142)$$

Then, we define (with $E = \frac{W_\iota}{\iota' W_\iota}$), with $S := G'E$ ⁵⁶⁵⁷

$$z_t = S'Q(1 - \gamma^e G)^{-1} \check{u}_t^e \quad (143)$$

$$= S'Q(1 - \gamma^e G)^{-1} Q(I - \gamma^e G) y_t. \quad (144)$$

Our key moment is as before, equation (97)⁵⁸

$$\mathbb{E}[z_t(y_{\tilde{E}t} - \gamma y_{St})] = 0. \quad (146)$$

This yields an estimate of γ . Hence, we simply replace the definition of the GIV (96) by (144). In fact, we can show that when $\Lambda = \iota$, the estimator (144) is equal to the estimator (96). In that sense it is its natural generalization⁵⁹

We note that we have a fixed point: an initial γ^e gives an estimate of γ ; that's then the new estimate γ^e , and we re-iterate the process, until convergence.

⁵⁶A variant is:

$$z_t = S'(1 - \gamma^e G)^{-1} \check{u}_t^e = S'(1 - \gamma^e G)^{-1} Q(I - \gamma^e G) y_t.$$

The advantage of formulation (97) is that in the simple case of the previous subsection (with $\Lambda = \iota$), then it recovers the estimator of that subsection.

⁵⁷The proof shows that this choice works. This particular choice is heuristically motivated by the analogy with (166) and (94), and the fact that when $\eta_t = 0$, $y_t = (I - \gamma G)^{-1} u_t$, so that $y_{St} = S'y_t = \gamma S'(1 - \gamma^e G)^{-1} u_t$. Hence in some loose sense z_t is a good idiosyncratic-based approximation of y_{St} .

⁵⁸Here is the proof. At the right estimator $\gamma = \gamma^e$,

$$z_t = S'Q(1 - \gamma^e G)^{-1} Qu_t = c'Qu_t, \quad (145)$$

with $c' := S'Q(1 - \gamma^e G)^{-1}$. We also have

$$y_{\tilde{E}t} - \gamma y_{St} = E'(I - \gamma G) y_t = E'(\Lambda \eta_t + u_t).$$

This implies that

$$\mathbb{E}[(y_{\tilde{E}t} - \gamma y_{St}) z_t] = \mathbb{E}[E' u_t u_t' Q' c] = E' V Q' c = \frac{1}{\iota' W_\iota} \iota' Q' c = 0,$$

as $Q\iota = 0$.

⁵⁹In the case $\Lambda = \iota$, then $Q = I - \iota E'$, so

$$E' G Q = E' G (I - \iota E') = E' G - E' = E' (G - I),$$

so that the estimator in (144) can be written:

$$z_t = E' (G - I) y_t = E' G Q y_t.$$

On the other hand,

$$\begin{aligned} A &:= GQ(1 - \gamma^e G)^{-1} Q(I - \gamma^e G) = GQ(1 - \gamma^e G)^{-1} (I - \iota E') (I - \gamma^e G) \\ &= GQ(1 - \gamma^e G)^{-1} (I - \gamma^e G) - GQ(1 - \gamma^e G)^{-1} \iota E' (I - \gamma^e G) \\ &= GQ - 0 \text{ as } G\iota = \iota \text{ and } Q\iota = 0 \\ &= GQ, \end{aligned}$$

Suppose that instead we use⁶⁰

$$Z_t = S' (1 - \gamma^e G)^{-1} \check{u}_t^e. \quad (147)$$

Suppose that $\gamma^e = \gamma$. Then we can write:⁶¹

$$y_{St} = Z_t + \varepsilon_t, \quad \mathbb{E}[\varepsilon_t Z_t] = 0, \quad (148)$$

for $\varepsilon_t = S' (I - \gamma^e G)^{-1} ((I - Q) u_t + \Lambda \eta_t)$. Also, we have:

$$y_{Et} = \gamma y_{St} + E' \Lambda \eta_t + u_{Et}.$$

Hence, we can estimate γ by OLS:

$$y_{Et} = \gamma Z_t + \beta' \eta_t + \varepsilon_t^{yE}. \quad (149)$$

This consistently estimates γ .

Calling z_t the GIV in the basic models (with $G = \iota' S$), $Z_t = \frac{z_t}{1 - \gamma^e}$.

12.7 Identification of social interactions and the reflection problem

There seems to be a contradiction between Section 6.3's finding that we do achieve identification, and Manski (1993)'s Proposition 2 and Bramoullé et al. (2009)'s Proposition 1, which seem also to state the impossibility of identification. Bramoullé et al. (2009) analyze social interactions of the type:

$$\mathbf{y}_t = \beta \mathbf{G} \mathbf{y}_t + \gamma \mathbf{x}_t + \mathbf{y} \mathbf{G} \mathbf{x}_t + \boldsymbol{\varepsilon}_t \quad (150)$$

so that the estimator in (144) can be written:

$$z_t = E' G Q (1 - \gamma^e G)^{-1} Q (I - \gamma^e G) y_t = E' A y_t = E' G Q y_t = E' (G - I) y_t.$$

Hence, when $\Lambda = \iota$, the estimators in (96) and (144) are identical.

⁶⁰Note that the γ^e in the definition of Z_t need not be the same γ^e used above to construct the \check{u}_t^e ; i.e., we could have $Z_t = S' (1 - \gamma^{e,2} G)^{-1} \check{u}_t^e$ for some other $\gamma^{e,2}$. There is still a fixed point though, and in the limit the estimated γ in (149) should also be equal to the γ^e and $\gamma^{e,2}$.

⁶¹Here is the proof. We saw that $\check{u}_t^e = Q u_t$, so,

$$Z_t = S' (I - \gamma^e G)^{-1} Q u_t,$$

while, $a' := S' (I - \gamma^e G)^{-1}$,

$$y_{St} = S' (I - \gamma^e G)^{-1} (Q u_t + (I - Q) u_t + \Lambda \eta_t) = Z_t + a' (I - Q) u_t + a' \Lambda \eta_t.$$

From (160), we have $\mathbb{E}[a' (I - Q) u_t z_t] = 0$. Hence we have $\mathbb{E}[\varepsilon_t Z_t] = 0$.

with $\mathbb{E}[\varepsilon_t|\mathbf{x}_t] = 0$. In their main result, they conclude that if the matrices I, G, G^2 are not linearly independent, then the system is not identified. However, in our setup $G = \iota S'$ (where ι is a vector of 1's) so that $G^2 = G$ and we satisfy [Bramoullé et al. \(2009\)](#)'s condition that seems to guarantee the impossibility of identification. However, we can identify the parameters, as we saw in [Section 6.3](#). How do we solve that seeming contradiction?

The short answer is that [Manski \(1993\)](#) and [Bramoullé et al. \(2009\)](#) do not consider anything like a GIV, as they immediately reason on averages based on observables, eschewing any exploration of the noise. In contrast, GIVs are all about exploring some structure in the noise — the idiosyncratic shocks of large entities. For instance Manski considers something akin to:

$$\mathbb{E}[\mathbf{y}_t|\mathbf{x}_t] = \beta \mathbf{G} \mathbb{E}[\mathbf{y}_t|\mathbf{x}_t] + \gamma \mathbf{x}_t + \mathbf{y} \mathbf{G} \mathbf{x}_t, \quad (151)$$

where all the noise has been averaged out.

Indeed, we do impose some structure, namely:

$$\varepsilon_{it} = \eta_t + u_{it}, \quad u_{it} \text{ i.i.d., orthogonal to } \eta_t, \quad (152)$$

and that was helpful to derive the GIV estimator ([Section 6.3](#)).

It would be interesting to show weaker conditions, or even necessary and sufficient conditions. We leave a full treatment that to future research. Still, we offer a few remarks with more general sufficient conditions for identification via GIV.

We can generalize the noise condition [152](#) (while staying with our setup $G = \iota S'$) to the more general condition:

$$\varepsilon_{it} = \lambda_i \eta_t + \sigma_i v_{it}, \quad (153)$$

where λ_i are scalar and η_t, v_{it} all uncorrelated (including across i 's). More generally, a “low rank” representation where $\eta_t \in \mathbb{R}^k$ with a low k is admissible too.^{[62](#)}

Second, we can generalize to the case where $\mathbf{G}^2 = \mathbf{G}$ (the case where G^2 is a linear combination of G and I is similar^{[63](#)}), which seems to leads to the impossibility of identification in [Bramoullé et al. \(2009\)](#). This is formalized here (and proved in the appendix, with a constructive identification procedure).

Proposition 5 (Identification achieved in the [Bramoullé et al. \(2009\)](#) setup). *Suppose that $G^2 = G$, which is satisfied in our basic setup, but leads to the impossibility of identification in the [Bramoullé et al. \(2009\)](#) setup without further assumptions. Suppose also the “simple noise structure” assumption [152](#). Suppose also the existence of two n -dimensional vectors S and Γ satisfying*

$$G' S = S, \quad G' \Gamma = 0, \quad \iota' S \neq 0, \quad \Gamma' S \neq 0. \quad (154)$$

⁶²Informally, this generates $2n$ unknowns (λ_i, σ_{it}), while the variance-covariance matrix has dimensions $\frac{n(n-1)}{2}$.

⁶³It can be reduced to that case by rescaling $H = b_0 + b_1 G$ with the right coefficient, with $H^2 = H$.

Then GIV is possible in that setup, i.e. with the GIV $z_t = \Gamma'y_t$, we can identify the coefficients (β, γ, y) (and indeed β it was assumed that $\gamma = y = 0$).

In our basic setup, we had S_i the relative sizes, and $G = \iota S'$, $\Gamma = S - \frac{\iota}{N}$. Hence (154) is an abstract generalization of our concrete conditions.

Hence, in many situations of interest we can be quite confident that condition (154) is satisfied.⁶⁴

In conclusion: our GIV approach gives some renewed hope for identification in the context of social influence and reflection problems. Indeed, it provides a way to achieve identification where it seemed impossible. Informally, this is by exploiting the idiosyncratic noise of “large players”. Formally, and less intuitively, it is by exploiting a little bit of structure in the noise (so that there is a low-dimensional common noise). Future research might profitably firm up the exact necessary and sufficient conditions for this.

12.8 When only some shocks are kept in the GIV

If we truncate the residuals, i.e. use

$$z_t = \sum_i \tau(S_i(u_{it} - u_{Et}))$$

for the hard thresholding function

$$\tau(x) = x1_{|x| \geq b}$$

for some $b > 0$, then everything works too. Indeed, we have $\check{u}_{it} := u_{it} - u_{Et}$ is orthogonal to u_{Et} . Let us assume that it is independent. In our basic example of Section 2.2, we still have $\mathbb{E}[(p_t - \alpha y_{St})z_t] = 0$, so that the IV procedure (16) still works. Likewise, in the more complex supply and demand case, the IV relations (34) and (36) still hold.

Furthermore, the OLS estimates still hold. The key is that we can write:

$$u_{\Gamma t} = z_t + z_t^<,$$

where $\tau^<(x) = x1_{|x| < b}$, and $z_t^< = \sum_i \tau^<(S_i \check{u}_{it})$, so that $z_t \perp z_t^<$. Hence, regression $u_{\Gamma t}$ on this truncated z_t gives a coefficient of 1, and all the analysis goes through.

12.9 When the researcher assumes too much homogeneity

Take the supply and demand example, and imagine that the econometrician assumes a homogeneous elasticity of demand ϕ^d , even though they're heterogeneous, they're ϕ_i^d . What happens then?

⁶⁴As $G^2 = G$, one can always find vectors Γ, S satisfying the first 3 conditions (provided n is big enough and G is not the identity nor 0), and the last one is rather “generically” easy to satisfy.

The model (26)-(27) becomes, for the demand:

$$y_{it} = \phi_i^d p_t + \lambda_i \eta_t + u_{it},$$

and for the supply

$$s_t = \phi^s p_t + \varepsilon_t.$$

As supply equals demand, $y_{St} = s_t$, which gives the price

$$p_t = \frac{u_{St} + \eta_t - \varepsilon_t}{\phi^s - \phi_S^d}. \quad (155)$$

In this thought experiment, the econometrician assumes $\phi_i^d = \phi_i$. He runs a panel model for $y_{it} - y_{Et}$, and we assume that it's large enough that he can extract η_t 's, successfully.⁶⁵ The GIV (we use the notation Z_t rather than z_t to denote the GIV before controls by η_t) is then

$$Z_t := y_{\Gamma t} = \phi_{\Gamma}^d p_t + \lambda_{\Gamma} \eta_t + u_{\Gamma t} = \left(1 + \frac{\phi_{\Gamma}^d}{\phi^s - \phi_S^d}\right) u_{\Gamma t} + \lambda^Z \tilde{\eta}_t = \frac{1}{\psi} u_{\Gamma t} + \lambda^Z \tilde{\eta}_t$$

so

$$Z_t = \frac{1}{\psi} u_{\Gamma t} + \lambda^Z \tilde{\eta}_t, \quad \frac{1}{\psi} = \frac{\phi^d - \phi_E^d}{\phi^s - \phi_S^d}, \quad (156)$$

where $\frac{1}{\psi} = 1$ in the common-elasticity case, $\tilde{\eta}_t = (\eta_t, \varepsilon_t, u_{Et})$ gathers the common shocks, and λ^Z is a vector of loadings.

Hence, when we run the first stage

$$p_t = b^p Z_t + \beta^p \eta_t + \varepsilon_t^p,$$

we will gather

$$b^p = \frac{1}{\phi^s - \phi_E^d}.$$

If we run

$$s_t = b^s Z_t + \beta^s \eta_t + \varepsilon_t^s,$$

we will estimate

$$b^s = \frac{\phi^s}{\phi^s - \phi_E^d}$$

The ratio of the two coefficients still gives ϕ^s . Likewise, the IV on the elasticity of demand will give ϕ_E^d .

In the polar opposite where η_t cannot be estimated or controlled for, then the simple procedure becomes biased, however, as (156) shows. To fix it, one can estimate the model with non-parametric

⁶⁵One of the factors, formally, will be p_t . We assume that it is not included in the vector of factors η_t .

coefficients (Section [10.1](#)).

12.10 Link between our initial examples and the general framework

Let us make the link between our initial examples and the general setup of Section [3.1](#).

The basic example of Section [2.2](#) was:

$$y_{it} = \eta_t + u_{it}, \quad p_t = \alpha y_{St} + \varepsilon_t.$$

The factors are p_t and η_t :

$$F_t^1 = p_t, \quad \alpha^1 = \alpha, \quad \lambda^1 = 0, \quad \eta_t^1 = \varepsilon_t,$$

$$F_t^2 = \eta_t, \quad \alpha^2 = 0, \quad \lambda^2 = 1, \quad \eta_t^2 = \eta_t.$$

Factor F_t^1 is endogenous, factor F_t^2 is exogenous. They are both observable (via $y_{Et} = \eta_t + O\left(\frac{1}{\sqrt{N}}\right)$). So, $\mathcal{C}_t = \eta_t$. The multiplier is $M = 1$.

Let us next make the link with supply-and-demand example of Section [2.3](#), which was:

$$y_{it} = \phi^d p_t + \eta_t + u_{it}, \quad p_t = \frac{y_{St}}{\phi^s} - \frac{\varepsilon_t}{\phi^s}.$$

The factors are also p_t and η_t :

$$F_t^1 = p_t, \quad \alpha^1 = \frac{1}{\phi^s}, \quad \lambda^1 = \phi^d, \quad \eta_t^1 = -\frac{\varepsilon_t}{\phi^s},$$

$$F_t^2 = \eta_t, \quad \alpha^2 = 0, \quad \lambda^2 = 1, \quad \eta_t^2 = \eta_t.$$

Here F_t^1 is exogenous and hidden, while F_t^2 is endogenous and observable. So, \mathcal{C}_t is empty.^{[66](#)} The multiplier is $M = \frac{1}{1 - \alpha^1 \lambda^1} = \frac{\phi^s}{\phi^s - \phi^d}$, as was estimated in [\(39\)](#).

13 Proofs omitted in the paper

13.1 Proof of Proposition [4](#)

13.1.1 Parametric Identification

We start with the parametric case, deferring the semi-parametric case.

⁶⁶This is why we could do only plain OLS regression in [\(39\)](#), without any controls like y_{Et} , unlike in the very simple initial example of Section [2.2](#)

The solution is, with $\lambda_S = S'\Lambda$ a $1 \times r$ vector, $M = \frac{1}{1-\lambda_S\alpha}$,

$$y_{St} = M(u_{St} + \lambda_S\eta_t + C_{St}m), \quad (157)$$

$$y_t = u_t + \Lambda[\alpha M(u_{St} + \lambda_S\eta_t + C_{St}m) + \eta_t] + C_tm. \quad (158)$$

We take the parametric case (the semi-parametric case will then be an easy corollary). This is, we have some characteristics x_{it} of actors (e.g. countries or firms), and a priori knowledge that $\lambda_{it} = X_{it}R$ for some r -dimensional vector $X_{it} = (1, x_{it})$ with x_{it} is a $r-1$ dimensional vector, and R is a $r \times r$ matrix. By rotation-invariance of the η_t (which is an r dimensional vector), we can take the case where $R = I$. Hence, in that sense we know the loadings $\lambda_{it} = x_{it}$ – but don't know the variance-covariance matrix V^η of the η_t .

Given a symmetric matrix W of size $N \times N$ (which, later, will optimally be $W = (V^u)^{-1}$, but we don't use that here yet) we define another $N \times N$ matrix:⁶⁷

$$Q^{\Lambda, W} = I - \Lambda(\Lambda'W\Lambda)^{-1}\Lambda'W, \quad (159)$$

so that $Q = Q^{\Lambda, W}$ satisfies:

$$Q\Lambda = 0, \quad Q'W\Lambda = 0, \quad (I - Q)W^{-1}Q' = 0, \quad Q^2 = Q. \quad (160)$$

Roughly, Q is the projection on the space orthogonal to the Λ , but with a scalar product that depends on W . Hence, (158) implies:

$$Qy_t = Qu_t + QC_tm. \quad (161)$$

Defining, for a vector Y_t ,

$$\check{Y}_t := QY_t, \quad (162)$$

we have

$$\check{y}_t = \check{u}_t + \check{C}_tm. \quad (163)$$

The controls C_t^k are all assumed to have non-zero cross-sectional variation: this is what allows to identify their m . A variable that's an “aggregate control” without cross-sectional variation (e.g. a time fixed effect, or maybe the world price of oil if we study the macroeconomics of a small country not affecting it) will be classified as an F_t^f – it's in $\mathcal{F}^{\text{Exo}, O}$ the set of observable, exogenous factors.

⁶⁷For instance, in our basic example with uniform loading $\Lambda = \iota$, $Q = I - \iota E'$, where $E = \frac{W\iota}{\iota'W\iota}$.

13.1.2 Estimating multipliers α^f, M by GIV

We assume that we have identified W^u (up to a multiplicative factor), either because we know for instance that $W^u = \sigma_u^2 I$, or because of the material in Section 13.1.4.

We treat now the more GIV-specific topic of how to estimate the α^f and M . We set $Q = Q^{\Lambda, W^u}$ as in (159). Then, (161) gave

$$Qy_t = Qu_t + QC_tm. \quad (164)$$

Let us define

$$\Gamma := Q'S, \quad (165)$$

and define the GIV to be $z_{\Gamma t} := \Gamma'(y_t - C_tm)$, which gives:

$$z_{\Gamma t} := \Gamma'(y_t - C_tm) = \Gamma'u_t = u_{\Gamma t}. \quad (166)$$

*This relation means that we identify $u_{\Gamma t}$ exactly, even though we estimate the η_t with errors.*⁶⁸

The GIV is possible if and only if

$$\Gamma \neq 0. \quad (167)$$

This is exactly what motivated Assumption 1 mentioned in Proposition 4.⁶⁹

We define $E = S - \Gamma = (I - Q')S$, so that

$$u_{St} = u_{\Gamma t} + u_{Et},$$

then (160) implies that $\mathbb{E}[u_{Et}u_{\Gamma t}] = S'(I - Q)V^uQ'S = 0$, so that we have the relation:^{70,71}

$$\mathbb{E}[u_{Et}u_{\Gamma t}] = 0. \quad (170)$$

⁶⁸This may be surprising, but consider the following simple case to see how this is true: if $y_{it} = \lambda p_t + \eta_t + u_{it}$, then as $\Gamma'v = 0$, $y_{\Gamma t} = u_{\Gamma t}$. So we perfectly measure the $u_{\Gamma t}$. This relation with the Q generalizes that simple example with more complex factors.

⁶⁹If $V^u S$ was spanned by the Λ , we could write $S = W\Lambda b$ for some vector b , and we'd have $Q'S = 0$, by 160. Conversely, if S is not spanned by the $V^u S$, then it is easy to check that $\Gamma \neq 0$.

⁷⁰In addition, the value

$$z_{Ct} := C_S'tm \quad (168)$$

is also a valid instrument. Hence, the following is an instrument:

$$z_t = z_{\Gamma t} + z_{Ct}. \quad (169)$$

In this paper we mostly use $z_{\Gamma t}$ as an instrument though, to insist on what is GIV-specific.

⁷¹As always in this paper, this relation leads to clean relations with finite N , but it relies on doing “Generalize least squares” with the proper weight matrix $W = (V^u)^{-1}$. In the general case, $\mathbb{E}[u_{Et}u_{\Gamma t}] = O\left(\frac{1}{\sqrt{N}}\right)$, so that this relation is likely to be approximately true in most cases of interest.

Hence, (157) reads:

$$y_{St} = M(u_{St} + \lambda_S \eta_t + C_{St}m) = (u_{\Gamma t} + u_{Et} + \lambda_S \eta_t + C_{St}m) = Mz_t + \varepsilon_t^{y_S} \quad (171)$$

for $\varepsilon_t^{y_S} = M(u_{Et} + \lambda_S \eta_t)$ uncorrelated with z_t . Hence, M will be consistently estimated by regressing y_{St} :

$$\mathbb{E}[(y_{St} - Mz_t)z_t] = 0. \quad (172)$$

Likewise, for f an observable, endogenous control (i.e., one such that α^f need not be 0 a priori), we can regress:

$$F_t^f = \alpha^f Mz_t + \varepsilon_t^{F^f}$$

with $\varepsilon_t^{F^f} = \eta_t + \alpha M \varepsilon_t^{y_S}$ so that we can estimate αM consistently by the OLS regression of F_t on z_t .

This shows identification even when we do not control for estimated factors η_t^e . To gain statistical power, it is useful to control for estimated η_t^e . We go on to this topic now.

13.1.3 Controlling for other estimated exogenous factors η_t^e

If we have a cross-sectionally important factor η_t^f , we may want to control for it to gain statistical precision. To do so, we first want to extract the factor. For notational simplicity, we assume that we removed all the controls $C_t m$ (i.e., we replace y_t by $y_t - C_t m$).

We define the $r \times N$ matrix⁷²

$$L_t := (\Lambda_t' W \Lambda_t)^{-1} \Lambda_t' W, \quad (173)$$

so that

$$L_t \Lambda_t = I_r. \quad (174)$$

Next, using $Q \Lambda_t = 0$ in the factor structure (53) gives:

$$(I - Q) y_t = \Lambda_t F_t + (I - Q) u_t.$$

Premultiplying this by L_t gives:

$$L_t (I - Q) y_t = F_t + L_t (I - Q) u_t.$$

Hence, an estimate of the F_t is

$$F_t^e := L_t (I - Q) y_t. \quad (175)$$

Indeed, we will have

$$F_t^e = F_t + \varepsilon_t^{F^e},$$

⁷²We might also take $L_t := (\Lambda_t' \Lambda_t)^{-1} \Lambda_t'$.

where $\varepsilon_t^{F^e} = -L_t(I - Q)u_t$ is a small error. In addition, this error is orthogonal to $z_{\Gamma t}$ ⁷³

$$\mathbb{E}[F_t^e z_{\Gamma t}] = 0, \quad (176)$$

so that the measurement error in the factors does not introduce a bias when estimating M .

Given our assumptions with $\Lambda_{it} = (1, x_{it})$, with $y_t = \Lambda_t F_t + u_t + C_t^y m$, we can write $F_t = (F_t^1, F_t^x)$, and where F_t^1 is the factor multiplying the “1” and F_t^x is the factor multiplying the x_{it} , so that $\Lambda_{it} F_t = F_t^1 + x_{it} F_t^x$. Given this, decompose $F_t = (F_t^{\text{endo}}, F_t^{\text{exo}}) = (F_t^1, F_t^x)$ with endogenous factors (i.e., affected by u_{it}) and exogenous factors. So here, $F_t^x = F_t^{\text{endo}}$.

We keep $\eta_t^e := F_t^{\text{exo}, e}$ as use them as controls, as it satisfies (176). Note that the standard errors returned by OLS will be trustworthy, because of (176) again.

This is all a bit abstract, so to get a concrete sense of the situation, let us take the main example, where $\Lambda_{it} = (1, x_{it})$ with $x_{Et} = 0$ and $y_{it} = \gamma p_t + \varepsilon_t + x_i \eta_t + u_{it}$, and $W = I \sigma_u^{-2}$. So factors are: $F_t = (F_t^{\text{endo}}, F_t^{\text{exo}}) = (\gamma p_t + \varepsilon_t, \eta_t)$, we have $I - Q = \Lambda L$, so $L(I - Q) = L \Lambda L = L$, hence $Ly_t = \left(y_{Et}, \frac{x_t' y_t}{\|x_t\|^2}\right)$, so that the error is

$$F_t^e - F_t = \varepsilon_t^{F^e} = -L(I - Q)u_t = \left(u_{Et}, \frac{x_t' u_t}{\|x_t\|^2}\right), \quad (177)$$

and the standard deviation of its components are $\frac{\sigma_u}{\sqrt{N}} \left(1, \frac{1}{\sigma_x}\right)$. The factor analysis recovers (up to that error) $F_t = (F_t^{\text{endo}}, F_t^{\text{exo}}) = (\gamma p_t + \varepsilon_t, \eta_t)$, so it recovers

$$\eta_t^e = \eta_t + \frac{x_t' u_t}{\|x_t\|^2} \quad (178)$$

but not ε_t . We can use that η_t^e as a control in the regression.

In conclusion, with a factor model (with known factor loadings Λ , but unknown factor covariance matrix V^η), we have identified V^u , and gotten a GIV, which gave $M, \alpha M$.

Even though all worked with finite N (but as always, $T \rightarrow \infty$), and we don't consistently estimate η_t , we still have a consistent estimator for the GIV.

⁷³Indeed, (166) gives

$$-\mathbb{E}[F_t^e z_t] = L(I - Q)V^u Q' S = L(I - Q)W^{-1}Q' S = 0$$

using (160) $(I - Q)W^{-1}Q' = 0$.

13.1.4 Estimating the variance-covariance matrix of the residuals, V^u

First, we estimate m , using (163) Basically, we can estimate m by OLS. It's pretty easy, as we have $(N - r) \times T$ effective values to use (where r is the number of factors):

$$m^e = \mathbb{E}_T [\check{C}'_t W \check{C}_t]^{-1} \mathbb{E}_T [\check{C}'_t W \check{y}_t]. \quad (179)$$

Next, we estimate V^u . We have $\check{u}_t := Q(y_t - C_t m^e) = Q u_t$, so that:

$$V^{\check{u}} = \mathbb{E} [\check{u}_t \check{u}'_t] = Q \mathbb{E} [u_t u'_t] Q. \quad (180)$$

We consider the case where we have a priori knowledge that V^u diagonal. Let us call $D^u = (V^u_{ii})_{i=1\dots N}$ and similarly for $D^{\check{u}} = (V^{\check{u}}_{ii})_{i=1\dots N}$ (so they are vectors of dimension N) and a new matrix $R_{ij} := (Q^W_{ij})^2$. Then:

$$D^{\check{u}}_i := V^{\check{u}}_{ii} = \sum_j Q^W_{ij} V^u_{jj} Q^W_{ij} = \sum_j R_{ij} D_j,$$

i.e., $D^{\check{u}} = R D^u$ so we recover

$$D^{u,e} = R^{-1} D^{\check{u}}. \quad (181)$$

Parametric variant We can do a parametric variant. We parametrize $\ln \sigma^2_{u_i} = \beta^\sigma x_i^\sigma$ for some vector of characteristics x_i^σ , e.g. log size or log volatility (and x_i^σ has 1 as the first component) — σ is just a superscript here, not an exponent. So, we estimate β^σ by regressing the log estimated variance from (181) on the characteristics:

$$\ln D^{u,e}_i = \beta^\sigma x_i^\sigma + \varepsilon_i,$$

and take the fitted values for the diagonal covariance matrix of the u_i 's:

$$D^u_i = e^{\beta^\sigma x_i^\sigma}. \quad (182)$$

Loop over W This gives a consistent estimate of D^u , for any W . Now, Bayesian considerations indicate that the optimum W is

$$W = (V^u)^{-1}. \quad (183)$$

So, we can loop: a good initial W is probably $1/\text{var}(y_i)$. This gives an estimate of D^u , and a new, better estimate of $W = \text{diag}(1/D^u_i)$. We keep looping until convergence. We have consistently estimated the variance matrix of V^u .

13.1.5 Semi-parametric case

Suppose now the semi-parametric case

$$\lambda_{it}^f = \lambda_0^f + \lambda_1^f x_{it}^f + \zeta_i^f,$$

and then we apply the above parametric procedure. For notational simplicity, we assume that all the control and constants are 0, as they are inessential. So, with $X_{it} = (1, x_{it})$ and $\lambda_X^f = (\lambda_0^f, \lambda_1^f)$ we have:

$$\Lambda_t = X_t \lambda_X + \zeta \quad (184)$$

and

$$y_t = (X_t \lambda_X + \zeta) F_t + u_t,$$

Recall also that we use $Q = Q^{X,W}$, so $QX = 0$. Then, the GIV will be as in (166)

$$z_{\Gamma t} := \Gamma' y_t = S' Q [(X_t \lambda_X + \zeta) F_t + u_t] = S' Q [\zeta F_t + u_t]$$

$$z_{\Gamma t} = u_{\Gamma t} + \zeta_{\Gamma} \eta_t \quad (185)$$

Hence, our GIV is partially polluted by a small $\zeta_{\Gamma} \eta_t$. However, as we will control for η_t^e in the regression, this part $\zeta_{\Gamma} \eta_t$ will be largely controlled for, and will not impact the results.

We sometimes use the close cousin

$$Z_t := y_{\Gamma t} = u_{\Gamma t} + \lambda_{\Gamma} \eta_t \quad (186)$$

Then again, it will be controlled for in the regression, as we control for η_t^e .

Note that there are two “ Γ ” here. The plain one is $\Gamma^0 = S - E$ (with $E_i = \frac{1}{N}$, if we stay in the homoskedastic case). The other, elaborated in this section, is $\Gamma^Q = Q'S$, where $Q = Q^{X,W}$ (with $W = I$ in the homoskedastic case) is given by (23). In principle, it is better to use Γ^Q than Γ^0 , as it “fully purges” the parametric factor in the purely parametric case.. By extension it should also be a bit better in the semi-parametric case (as we need to purge a “small” ζ_{Γ} — which is 0 in the parametric case — rather than a “potentially big” λ_{Γ}). We advocate the Γ^Q , but in practice using Γ^0 gives similar results.

If we use $\Gamma = \Gamma^Q$, then $\lambda_{\Gamma} = \zeta_{\Gamma^Q}$ so that

$$Z_t = u_{\Gamma t} + \zeta_{\Gamma^Q} \eta_t = z_t. \quad (187)$$

13.2 Other proofs

Proof of Proposition 5 The identification goes as follows. By rescaling S , we impose $\iota'S = 1$. Define $E := S - \Gamma$ (which is $\frac{1}{N}\iota$ in our framework), and form

$$y_{Et} = E'y_t, \quad y_{St} := S'y_t$$

which are our generalized “equal weighted” and “value weighted” averages – for more abstract setting. Then, premultiplying (150) by Γ' and S' gives:

$$z_t := \Gamma'y_t = \gamma x_{\Gamma t} + u_{\Gamma t}.$$

Hence, estimating this by OLS we can obtain γ , and $\text{var}(u_{\Gamma t})$, so that we obtain also σ_u^2 . Next,

$$y_E = \beta y_S + \gamma x_E + y x_S + \eta + u_E,$$

so that

$$\mathbb{E}[(y_{Et} - \beta y_{St} - \gamma x_{Et} - y x_{St})'(z_t, x_{St})] = (\mathbb{E}u_{Et}u_{\Gamma t}, 0). \quad (188)$$

The right-hand side is known, as $\mathbb{E}u_{Et}u_{\Gamma t} = E'\Gamma\sigma_u^2$, which is known. So, we have two unknowns β , y and 2 equations: we can solve the system. The condition $\Gamma'S = 0$ ensures that $\mathbb{E}[y_{St}z_t] \neq 0$.

14 A Bayesian perspective on GIVs

We will see that, under conditions of Gaussianity, our estimators are basically the MLE. As variables may not be Gaussian, we keep the general exposition (showing identification) free of distributional assumptions. If we assume that variables (u_{it}, η_t) are Gaussian, then a Bayesian analysis can be performed. We detail it here.

14.1 The general model: Bayesian version

Here we treat the general model of Section 3.1, in the case where the λ_{it} are the same, and equal to λ_i , and all factors are observed (except the η_t^y , as in $y_{it} = \sum_f \lambda^f F_t^f + u_{it} + \eta_t^y$). The very general case would involve the same ideas, too many layers of notations.

The data D is $D = \left(y_t, F_t^f\right)_{f=1\dots d_F, t=1\dots T}$, made of i.i.d. draws from a fixed distribution. To simplify the notations, we'll just denote by f the collection of all variables corresponding to factors (without explicitly mentioning that $f = 1 \dots d_F$)

The solution of the system features:

$$\begin{aligned} y_{St} - My_{\Gamma t} &= b^y \varepsilon_t, \\ F_t^f - \alpha^f My_{\Gamma t} &= b^f \varepsilon_t, \end{aligned}$$

for some vector b^y, b^f , and $\varepsilon_t := (u_{Et} + \eta_t^y, \eta_t^f)$.

Hence, we form: $\theta = (M, \alpha^f M)$; W a parametrization of the relevant variance matrices; $E(W) = \frac{V^u(W)^{-1}_t}{\iota' V^u(W)^{-1}_t}$ the corresponding quasi-equal weights vector, $\omega = (\theta, W)$, and form the key quantities:

$$Y_t(\omega) = (y_{St} - My_{\Gamma(W),t}, F_t^f - \alpha^f My_{\Gamma(W),t}). \quad (189)$$

We also keep track of

$$\check{y}_{it}(\omega) = y_{it} - y_{E(W),t} \quad (190)$$

and stack those two vectors together as $X_t(\theta)$, which contains all our information:

$$X_t(\theta) = (Y_t(\theta), \check{y}_t(\omega)). \quad (191)$$

The “trick to tractability” is to transform the data into that X_t .

There is an invertible matrix $A(\theta)$ such that $D_t = A(\theta) X_t$. Hence, there is no loss of information in using X_t as “conveniently processed” data, rather than the “unprocessed” data D_t . Hence, instead of $\ln \mathbb{P}(D_t|\omega)$, we’ll consider

$$\ln \mathbb{P}(X_t|\omega) = \ln \mathbb{P}(D_t|\omega) + \ln |A(\theta)|. \quad (192)$$

The Jacobian $|A(\theta)| := \det A$ is independent of all parameters ω .⁷⁴ Hence it can be discarded as a constant in the calculations.

The key simplifying observation is that (under the correct model), $\mathbb{E}\check{y}_t y_{Et} = 0$, so that $Y_t(\theta)$ and $\check{y}_t(\omega)$ have zero covariance. Hence, the log likelihood decouples, and we have

$$-2 \ln \mathbb{P}(D_t|\omega) = Y_t'(\omega) V^Y(W)^{-1} Y_t(\omega) + \check{y}_t(\omega)' (V^{\check{y}}(\omega))^{-1} \check{y}_t(\omega) + \ln |V^Y(W)| + \ln |V^{\check{y}}|. \quad (193)$$

As \check{y}_t lives in a space of dimension $N - 1$ (as $E'\check{y}_t = 0$), the value of $V^{\check{y}}$ is understood as being of the corresponding dimensions, $(N - 1) \times (N - 1)$.

Now, imagine that W has already been estimated, and do only the optimization w.r.t. θ . That gives:

$$\min_{\theta} \mathbb{E}_T Y_t'(\theta, W) V^Y(W)^{-1} Y_t(\theta, W) \quad (194)$$

⁷⁴First, go from X_t to $\tilde{D}_t = (F_t, y_{Et}, \check{y}_t)$ is upper triangular with 1 on the diagonals, so has determinant 1; second, go from \tilde{D} to D , which is independent of the ω .

The first order conditions are:

$$\mathbb{E}_T [(y_{\Gamma t}, 0) (V^Y)^{-1} Y_t] = 0, \quad \mathbb{E}_T [(0, y_{\Gamma t}) (V^Y)^{-1} Y_t] = 0,$$

i.e. (as $0 = \mathbb{E}_T [y_{\Gamma t} (V^Y)^{-1} Y_t] = (V^Y)^{-1} \mathbb{E}_T [y_{\Gamma t} Y_t]$) $\mathbb{E}_T [y_{\Gamma t} Y_t] = 0$ i.e.

$$\mathbb{E}_T \left[y_{\Gamma t} \left(y_{St} - M y_{\Gamma t}, F_t^f - \alpha^f M y_{\Gamma t} \right) \right] = 0.$$

Those are precisely the first order condition of the OLS estimation:

$$\min_M \mathbb{E}_T [(y_{St} - M y_{\Gamma t})^2], \quad \min_{\alpha^f M} \mathbb{E}_T \left[\left(F_t^f - \alpha^f M y_{\Gamma t} \right)^2 \right]. \quad (195)$$

Hence, our GIV is also the MLE estimator of $M, M\alpha^f$, when we have Gaussian distributions.

We can also go beyond MLE, and calculate full Bayesian posteriors. Then, the GIV gives an easy way to do finite-sample Bayesian updating. Assuming again for simplicity that we know the variance matrices, we have

$$\ln \mathbb{P}(\theta | D) = \ln \mathbb{P}(\theta) - \frac{1}{2} \sum_t Y_t(\theta) (V^Y)^{-1} Y_t(\theta) + K(D), \quad (196)$$

where $K(D)$ ensures that the probability sums to 1.

The rest of this subsection examines instantiations and variants of the general idea we just saw.

14.2 The supply and demand model of Section 2.3

This model which corresponds exactly to the general case, with a factor $F_t^f = p_t$, $p_t = \alpha^f y_{St} + \eta_t^f$ with $\alpha^f = \frac{1}{\mu}$ and $\eta^f = -\frac{\varepsilon}{\mu}$. Then, everything goes through.

14.3 The basic example with self-loop of Section 6.3.

We give a Bayesian treatment of this model of Section 6.3:

$$y_{it} = \gamma y_{St} + \eta_t + u_{it}.$$

We are given $D_t = y_t$. We wish to estimate $M = \frac{1}{1-\gamma}$. The vector of parameters of interest is $\theta = M$

$$Y_t(\omega) = y_{St} - M y_{\Gamma(W),t}.$$

As in the general procedure, we set:

$$\check{y}_{it}(\omega) = y_{it} - y_{E(W),t} \quad (197)$$

and

$$X_t(\theta) = (Y_t(\theta), \check{y}_t(\omega)). \quad (198)$$

In the true model, we have $Y_t = M(u_{Et} + \eta_t)$, so

$$-2 \ln \mathbb{P}(D_t | \omega) = \frac{Y_t^2(\omega)}{\sigma_Y^2} + \check{y}_t(\omega)' V^{\check{y}} \check{y}_t(\omega) + \ln \sigma_Y^2 + \ln |V^{\check{y}}|.$$

Suppose first that we have know the variance terms. Then, the MLE is simply to do

$$\max_M \mathbb{E}_T [Y_t(\omega)^2],$$

which is the identification condition we used, and it corresponds to running the OLS $\min_M \mathbb{E}_T (y_{St} - M y_{\Gamma(W)_t})^2$

Next, for the estimation of the variance terms, we optimize on σ_Y^2 , $V^{\check{y}}$. Asymptotically, that gives the true values.

14.4 The basic example without loop of Section 2.2.

We now detail the Bayesian version of our example in Section 2.2.

$$y_{it} = \eta_t + u_{it}, \quad p_t = \alpha y_{St} + \varepsilon_t,$$

and we'd like to estimate α especially (or, in a Bayesian context, update our prior on α). This example is actually a bit non-generic, as it endows the economist with a knowledge that $\lambda^f = 0$, which creates some subtle changes: it features the “recovered” factor y_{Et} , used as a regressor.

The data D is a set of $D = (y_t, p_t)_{t=1 \dots T}$, assumed to be i.i.d. draws from a fixed distribution.

We call $\theta = (\alpha, \beta)$ the set of “key” model parameters, and W , the variance-covariance matrix $V^{(u+\eta, \varepsilon)}$ (or, it could be some parametrization of it, e.g. if we assume that u is diagonal), the auxiliary parameter, and $\omega = (\theta, W)$ the full set of parameters. The correct value is ω^* .

Given y_t, p_t , we form

$$Y_t(\theta) = p_t - (\alpha y_{\Gamma t} + \beta y_{Et})$$

and $X_t(\theta) = (Y_t(\theta), y_t)$. As the correct parameter ω^* ,

$$Y_t(\theta^*) = \varepsilon_t^\perp$$

which is defined in the analysis is the “enriched OLS estimator” (Section 13). Hence, at the correct value, Y_t and y_t are uncorrelated. Call $V^X(\omega)$ the variance-covariance matrix of X_t .

We can start the Bayesian analysis:

$$\mathbb{P}(\omega | D) = \mathbb{P}(D | \omega) \mathbb{P}(\omega)$$

and

$$\ln \mathbb{P}(D|\omega) = \sum \ln \mathbb{P}(D_t|\omega)$$

with

$$-2 \ln \mathbb{P}(D_t|\omega) = \frac{Y_t(\theta)^2}{\sigma_{\varepsilon^\perp}^2} + y_t' V^y(W) y_t + \ln \sigma_{\varepsilon^\perp}^2 + \ln |V^y(W)|, \quad (199)$$

where here $|A|$ is the determinant of a matrix A .

Hence, the MLE estimator maximizes $\sum_t \ln \mathbb{P}(D_t|\omega)$ over $\omega = (\theta, W)$. The problem for θ separates as:

$$\min_{\alpha, \beta} \sum_t Y_t(\theta)^2,$$

i.e.

$$\min_{\alpha, \beta} \sum_t (p_t - (\alpha y_{\Gamma t} + \beta y_{Et}))^2,$$

which is the “enriched GIV-OLS estimator”. This shows that, with Gaussian distribution, the MLE is just our enriched GIV-OLS estimator.

Maximizing over the other parameters W will allow to recover the variance matrix (including that of ε_t, η_t).

If we have a small sample, we can just update rather than do MLE. The above shows that the “simplifying trick” is to form that statistics $Y_t(\theta)$, which allows for an interpretable updating of the parameters. For simplicity, suppose that we know the value of $V^y(W)$, and $\sigma_{\varepsilon^\perp}^2$.⁷⁵ However, we have a prior on $\theta = (\alpha, \beta)$, perhaps Gaussian distributed. Then, our posterior after observing the data D is:

$$\ln \mathbb{P}(\theta|D) = \ln \mathbb{P}(\theta) - \sum_t \frac{Y_t(\theta)^2}{2\sigma_{\varepsilon^\perp}^2} + K(D),$$

where $K(D)$ ensures that the probability sums to 1.

⁷⁵Otherwise, we can update our knowledge of those, which is standard though tedious to lay out.