

Notes on Bartik Instruments*

Paul Goldsmith-Pinkham

Isaac Sorkin

Henry Swift

This version: October 31, 2016.

Preliminary.

Abstract

Bartik instruments are used to generate plausibly exogenous labor demand shocks. They are constructed by interacting the distribution of industry shares across locations with national industry growth rates. We develop a formal econometric structure for the Bartik instrument. Our structure delivers three main insights. First, the necessary exogeneity condition is with respect to the industry shares in a location. In particular, we relate this condition to identification in continuous difference-in-differences. Second, the Bartik instrument exploits the inner product structure of the endogenous variable to reduce the dimensionality of the first-stage estimation problem. Third, this structure provides guidance about whether to use historical industry shares and how finely to divide industries. With these insights, we develop a checklist of recommendations for how to implement the Bartik instrument and how to test the plausibility of the exclusion restriction. We illustrate this checklist in the context of the canonical case of estimating the inverse elasticity of labor supply. We show that industry shares are correlated with education and other characteristics and that controlling for these characteristics significantly reduces the magnitude of the inverse elasticity of labor supply. We also show evidence of quantitatively important pre-trends.

*Goldsmith-Pinkham: Federal Reserve Bank of New York. Email: paulgp@gmail.com. Sorkin: Department of Economics, Stanford University. Email: sorkin@stanford.edu. Swift: Unaffiliated. Email: henryswift@gmail.com. The views expressed are those of the authors and do not necessarily reflect those of the Federal Reserve Bank of New York or the Federal Reserve Board. All errors are our own.

The Bartik, or shift-share, instrument was initially used in Bartik (1991) and popularized in Blanchard and Katz (1992). Since then, it has been widely used across many fields in economics, including labor, public, macroeconomics, international trade, and finance. (See Table 1.) The Bartik instrument is the prediction of local employment growth that comes from interacting local variation in industry employment shares with national industry employment growth rates. A common intuitive argument in favor of the Bartik instrument is that the national nature of the average industry growth rates avoids correlation with local economic conditions.¹ While there are many applications of Bartik-like instruments, for concreteness our running example is the canonical application of national industry growth rates interacted with industry composition.²

This note develops a formal econometric structure for the Bartik instrument. The key observation is that the endogenous variable has an inner product structure. Formally, $X = \mathbf{Z}'\mathbf{G} = \sum_k Z_k G_k$, where X is the location-specific employment growth rate, \mathbf{Z} is the vector of location-specific industry shares (where k denotes an industry) and \mathbf{G} is the vector of location-specific growth. The Bartik instrument is constructed by replacing the location-specific \mathbf{G} with some national average, $\bar{\mathbf{G}}$, so that $B = \mathbf{Z}'\bar{\mathbf{G}}$, and the researcher uses B as an instrument for X .

We first show that the validity of the Bartik instrument rests on the exogeneity of the location-specific industry shares, and not the growth rates. Indeed, using the Bartik instrument in two stage least squares is numerically equivalent to a generalized methods of moments estimator with the industry shares as instruments, where the weight matrix is constructed from the national industry growth weights. The role of the industry growth rates in the Bartik instrument is about power, not identification.³

We then show that the Bartik instrument should be thought of as an estimator designed to overcome the problem of a high-dimensional first-stage, rather than an instrument. Under the assumption of the exogeneity of industry shares (\mathbf{Z}), the first-stage estimation problem is to compute $\mathbb{E}[X|\mathbf{Z}]$. With the inner product structure of X , this first-stage estimation

¹For example, Bound and Holzer (2000, pg. 31): “The [Bartik] index should capture exogenous shifts in local labor demand that are predicted by the city-specific industry mix, while avoiding the endogeneity associated with local employment growth rates. We use this index as an instrument for the overall local employment growth.”

²For example, Altonji and Card (1991) interact initial immigrant composition with flows from sending countries, Chodorow-Reich (2014) interacts which banks firms borrow from with bank-level shocks stemming from the failure of Lehman Brothers, and Nakamura and Steinsson (2014) interact the geographic composition of defense spending with defense spending shocks.

³There is a sense in which properties of the growth rates are relevant to identification: if part of the industry component of the growth rate enters the location level error term in a way that is proportional to industry composition, then Bartik is no longer valid. Formally, however, Bartik fails because then the error term is correlated with industry composition, and not because of the growth rates.

problem can be rewritten as follows:

$$\mathbb{E}[X|\mathbf{Z}] = \mathbb{E}[\mathbf{Z}'\mathbf{G}|\mathbf{Z}] = \mathbf{Z}'\mathbb{E}[\mathbf{G}|\mathbf{Z}].$$

The dimension of this first-stage, however, is proportional to the number of industries, and thus high-dimensional. The Bartik instrument avoids a high-dimensional first-stage by focusing on the $\mathbb{E}[\mathbf{G}|\mathbf{Z}]$ component and making the following approximation:

$$\mathbb{E}[\mathbf{G}|\mathbf{Z}] \approx \mathbb{E}[\mathbf{G}].$$

This approximation ignores the information in how growth rates vary by industry composition; for example, it might be that the restaurant industry grows quickly in locations with a large entertainment industry presence, or the restaurant industry growth rate depends on the initial restaurant industry share. By ignoring this information, the Bartik instrument avoids the appearance of a high-dimensional first stage.⁴

Nevertheless, Bartik is still leveraging a high-dimensional set of instruments and so thinking about the implicit dimensionality of the first stage provides guidance as to how finely to divide industries and locations. In finite samples, having a high-dimensional first stage can induce bias by overfitting the first stage (i.e., subject to a many instruments problem). In practice, then, researchers should make choices to avoid having “too many” parameters to estimate relative to the amount of data. For example, by using a very fine set of industry shares with very few firms in each industry, in finite samples the Bartik estimator can do a poor job of approximating the first stage. A similar observation explains the desirability of the common practice of using a leave-one-out estimator to construct the national averages.

Given that the identifying assumption is the exogeneity of industry shares, our structure shows that if changes in local industry composition are partially driven by endogenous shocks, then updating industry shares may bias estimates. For example, if the endogeneous component of growth is serially correlated, then the industry shares are mechanically correlated with this endogeneous component, rendering the instrument invalid. Hence, researchers should use the earliest potential version of their industry shares to avoid this correlation.

By clarifying that the identifying assumption is the exogeneity of industry shares, our structure provides guidance to researchers about how to explore the plausibility of their

⁴Given the same identifying assumption as the Bartik instrument, there are potentially more powerful estimators that exploit knowledge of the industry growth rates across different locations. However, we show that in this context, industry-specific growth rates explain almost 25% percent of the variance of industry-by-location growth rates (using PUMAs and 3 digit industries). As a result, the use of industry averages in Bartik is a surprisingly well-supported approximation.

research design. We recommend two tests in particular. First, researchers should present balance tests in terms of industry shares to show that the industry shares are not related to other location characteristics. Because the Bartik instrument is about changes, these tests should be presented in terms of both levels and changes. Second, researchers should examine pre-trends. Testing for pre-trends in the Bartik setting cannot be done directly. If the instrument is valid and correlated through time, then we do not expect parallel trends to hold. By partialling out the expected effect of previous values of the instrument, however, parallel trends ought to hold.

Finally, we illustrate our points about the Bartik instrument in the empirical context of using the Bartik instrument as an instrument to estimate labor supply and in simulations. In the empirical example, we show that the industry shares are correlated with many observable characteristics (including education), controlling for these observable differences attenuates estimates, and that there appear to be pre-trends. Via simulation, we show that endogeneity of growth rates is not necessarily a problem for identification, that a small variance of the industry common component leads to noisy estimates, and that a small number of locations relative to the number of industries also appears to lead to bias.

We suspect that thoughtful users of the Bartik instrument will view some of the points made in this note as obvious or well-known folklore. We found, however, that the process of formalization led us to greater clarity and some new insights. Moreover, we believe that there is value in recording this folklore as it helps codify best practices and understandings that are inconsistently reflected in applied work. We note that while some papers (or literatures) understand these points, this understanding is not reflected in all papers using the Bartik instrument (even in high profile venues). Further details are available upon request.

1 Two simple cases

1.1 Case I: Cross-sectional data

We start with the two industry cross-sectional data case, where it is possible to see that the identifying assumption in Bartik is the exogeneity of the industry shares without the cumbersome notation of the many industry case. Let Y_i and X_i denote wage growth and employment growth in city i . We are interested in estimating β from the following equation:

$$Y_i = \alpha + X_i\beta + \epsilon_i, \quad (1.1)$$

where $i = 1, \dots, n$ and $\text{Cov}(\epsilon_i, X_i) \neq 0$. Hence, the OLS estimator for β is biased.

Suppose that there are two industries in each city, manufacturing and services. Let Z_{i1} and Z_{i2} denote the share of employment for the two industries, and let G_{i1} and G_{i2} denote

the growth in employment in each industry in city i . Observe that city-level growth is the weighted sum of the growth rates in the two industries: $X_i = G_{i1}Z_{i1} + G_{i2}Z_{i2}$. A researcher might be interested in the response of wages to employment growth driven by a labor demand shock.

Practically, a researcher would need exogenous variation in X_i , or instrument for this labor demand shock. In this context, the Bartik instrument is constructed as $B_i = \bar{G}_1 Z_{i1} + \bar{G}_2 Z_{i2}$, where \bar{G}_j is the average growth rate of industry j across all cities.⁵ The usual logic is that while industry-and-location-specific growth rates (G_{i1}, G_{i2}) might be correlated with ϵ_i , the national average ($\bar{G}_{i1}, \bar{G}_{i2}$) is not.⁶

We first define our two estimators.

DEFINITION 1.1. *Define the estimator given by using the Bartik instrument:*

$$\hat{\beta}_{2SLS}(B) = \frac{\sum_i B_i Y_i - n^{-1} \sum_i B_i \sum_j Y_j}{\sum_i B_i X_i - n^{-1} \sum_i B_i \sum_j X_j}, \quad (1.2)$$

and the estimator given by using industry shares as instruments:

$$\hat{\beta}_{2SLS}(Z_2) = \frac{\sum_i Z_{2i} Y_i - n^{-1} \sum_i Z_{2i} \sum_j Y_j}{\sum_i Z_{2i} X_i - n^{-1} \sum_i Z_{2i} \sum_j X_j}. \quad (1.3)$$

We can now show that these estimators are both consistent estimators of the parameter of interest.

PROPOSITION 1.1. *If $\bar{G}_1 - \bar{G}_2 \neq 0$, $\hat{\beta}_{2SLS}(B) = \hat{\beta}_{2SLS}(Z)$. If Z_{i1} is also independent of ϵ_i and the data is independently sampled across i , then both estimators are consistent estimates of β .*

Proof. See Appendix A. □

To understand this equivalence, it is helpful to write out the first stage explicitly and note that the industry growth rates function as weights on the shares.

REMARK 1.1 (Bartik Weighting). *Consider the two-stage system of equations for 2SLS*

$$Y_i = \alpha + X_i \beta + \epsilon_i \quad (1.4)$$

$$X_i = \tau + B_i \gamma + u_i. \quad (1.5)$$

⁵In what follows, we will show that the more appropriate estimator is to use the leave-out mean (excluding the i observation in estimating \bar{G} for B_i) and since Autor and Duggan (2003), this has become standard practice.

⁶E.g., the quote from Bound and Holzer (2000) in the main text, or the following quote from Autor and Duggan (2003, pg. 180): “Provided that national industry growth rates (excluding own state industry employment) are uncorrelated with state-level labor supply shocks, this approach will identify plausibly exogenous variation in state employment.”

Substitute the Bartik instrument into the first stage and use the fact that $Z_{i1} + Z_{i2} = 1$:

$$X_i = \tau + B_i\gamma + u_i \quad (1.6)$$

$$= \tau + (\bar{G}_1 Z_{i1} + \bar{G}_2 Z_{i2}) \gamma + u_i \quad (1.7)$$

$$= \underbrace{(\tau + \bar{G}_1 \gamma)}_{\tilde{\tau}} + Z_{i2} \underbrace{(\bar{G}_2 - \bar{G}_1) \gamma}_{\tilde{\gamma}} + u_i \quad (1.8)$$

$$= \underbrace{(\tau + \bar{G}_1 \gamma)}_{\tilde{\tau}} + Z_{i2} \underbrace{\Delta_G \gamma}_{\tilde{\gamma}} + u_i. \quad (1.9)$$

Here, the difference in growth rates, $\Delta_G = \bar{G}_2 - \bar{G}_1$, is a constant that weights Z_{i2} . Since $Z_{i1} + Z_{i2} = 1$, γ can be transformed exactly into $\tilde{\gamma}$ by division by Δ_G . The Δ_G drops out in the two-stage estimator since the equation is exactly identified.

Note the role of Δ_G in this expression: if $\Delta_G = 0$, then there is no variation on the right hand side of the first stage. This observation makes clear that the variation in the growth rates is about weighting the moment restrictions, rather than being an identifying restriction itself.

In the next section, we generalize this insight to the case with more than two industries and show that the Bartik instrument is numerically equivalent to using industry shares as instruments in a generalized method of moments setup if we use the industry growth rates to construct the weight matrix.

Many researchers are concerned about separating supply and demand shocks using Bartik. The following example shows that separating supply and demand shocks in the growth rates matters to the extent that the supply shocks enter the error term in a way that is proportional to industry composition.

REMARK 1.2 (“Supply” vs. “Demand” shocks). Suppose that the industry-level growth rates consist of two components: $G_1 = G'_1 + \epsilon_1$ and $G_2 = G'_2 + \epsilon_2$. In terms of interpretation, think of the ϵ_1 and ϵ_2 as supply shocks and the G'_1 and G'_2 as demand shocks. Moreover, suppose that these shocks enter the first stage error term in a way that is proportional to industry composition:

$$Y_i = \alpha + X_i\beta + \underbrace{Z_{i1}\epsilon_1 + Z_{i2}\epsilon_2 + \epsilon_i}_{\text{error term}}. \quad (1.10)$$

In this example, if $\epsilon_1 \neq 0$ and $\epsilon_2 \neq 0$ then Bartik is not valid because the instrument (the Z s) enters the error term.

Alternatively, given the same structure of G_1 and G_2 , suppose that the ϵ_1 and ϵ_2 enter the error term in a way that is not proportional to the error term. Let $\text{Corr}(R_{i1}, Z_{i1}) = 0$ and $\text{Corr}(R_{i2}, Z_{i2}) =$

0, and suppose that the shocks enter the first stage error term through the R :

$$Y_i = \alpha + X_i\beta + \underbrace{R_{i1}\epsilon_1 + R_{i2}\epsilon_2 + \epsilon_i}_{\text{error term}}. \quad (1.11)$$

In this example, Bartik is valid because the instrument does not enter the error term.

These examples emphasize that while identification comes from the exogeneity of the Z s, identification can fail if there are components of the growth rates (the G s) that enter the error term in a way that is proportional to the Z s.

1.2 Case II: Panel data

We now show that when Bartik is used in a panel, it is equivalent to allowing for time-variation in the weights on the industry shares. We also show that when there is serial correlation in the endogenous component of X , researchers should use the initial industry shares to estimate Bartik.

Maintaining the assumption of two industries, define the Bartik instrument so that it varies over time:

$$B_{it} = \bar{G}_{1t}Z_{i1t} + \bar{G}_{2t}Z_{i2t},$$

where now there are N locations and T time periods, indexed by i and t respectively.

In many applications, researchers are faced with the choice about whether to allow Z_{it} to vary over time, or fix it in an initial period. Let the initial period shares be denoted by Z_{ij}^0 . To see the need for fixing industry shares, consider the expression for next period's industry shares in the two-industry case, denoted by Z_{ik}^1 and let G_{ik}^0 be period 0 growth rates:

$$Z_{i1}^1 = \frac{G_{i1}^0 Z_{i1}^0}{G_{i1}^0 Z_{i1}^0 + G_{i2}^0 Z_{i2}^0}; \quad Z_{i2}^1 = \frac{G_{i2}^0 Z_{i2}^0}{G_{i1}^0 Z_{i1}^0 + G_{i2}^0 Z_{i2}^0}. \quad (1.12)$$

Note that the Bartik instrument is used due to concerns that either $\text{Cov}(G_{i1}^0, \epsilon_i^0) \neq 0$ or $\text{Cov}(G_{i2}^0, \epsilon_i^0) \neq 0$. If the ϵ_i^t are not independent over time (i.e., $\text{Cov}(\epsilon_i^0, \epsilon_i^1) \neq 0$), then updating the weights can induce a correlation between the instrument and the error term, rendering the instrument invalid.

REMARK 1.3 (General Industry Shares based on Initial Condition). *Generally, it is worth noting that this implies recursively:*

$$Z_{ij}^t = Z_{ij}^0 \prod_{s=0}^{t-1} \pi_{ij}^s \quad (1.13)$$

where $\pi_{ij}^s = \frac{G_{ij}^s}{X_i^s}$ is the relative growth in industry j in location i in period s . Hence, we can write the first-stage estimand as

$$\mathbb{E}(X|\mathbf{Z}^0) = \mathbf{Z}^{0'} \mathbb{E}\left(\prod_{s=0}^{t-1} \pi^s \mathbf{G}_t | \mathbf{Z}^0\right) \quad (1.14)$$

where π^s is a $K \times K$ diagonal matrix with the π_k^s (relative growth of industry k in period s) along the diagonals.

As a result of this potential bias, in what follows, we will use the initial industry share for each location, and hold it fixed across time periods. Hence, the Bartik estimator will be calculated as

$$B_{it} = \bar{G}_{1t} Z_{i1}^0 + \bar{G}_{2t} Z_{i2}^0.$$

To see the relationship between the cross-sectional and the panel estimating equations, it is helpful to write the setup with place and time fixed effects:

$$Y_{it} = \alpha_i + \alpha_t + X_{it}\beta + \epsilon_{it} \quad (1.15)$$

$$X_{it} = \tau_i + \tau_t + B_{it}\gamma + u_{it}. \quad (1.16)$$

Now substitute in the Bartik instrument and rearrange the first stage:

$$X_{it} = \tau_i + \tau_t + (\bar{G}_{1t} Z_{i1}^0 + \bar{G}_{2t} Z_{i2}^0) \gamma + u_{it} \quad (1.17)$$

$$= \tau_i + \underbrace{(\tau_t + \bar{G}_{1t}\gamma)}_{\tilde{\tau}_t} + Z_{i2}^0 \underbrace{(\bar{G}_{2t} - \bar{G}_{1t})\gamma}_{\tilde{\gamma}_t} + u_{it} \quad (1.18)$$

$$= \tau_i + \underbrace{(\tau_t + \bar{G}_{1t}\gamma)}_{\tilde{\tau}_t} + Z_{i2}^0 \underbrace{\Delta_{G,t}\gamma}_{\tilde{\gamma}_t} + u_{it}. \quad (1.19)$$

Here a critical difference emerges between Bartik and using the shares as instruments: using the time-invariant industry shares as instruments necessitates using time-varying coefficients, as otherwise the predicted effect of a given industry on growth is fixed, and would be subsumed by the location fixed effects. Bartik puts the time-variation into the calculation of the instrument. To see this, compare the first-stage using Z_{i2}^0 as an instrument:

$$X_{it} = \tilde{\tau}_i + \tilde{\tau}_t + Z_{i2}^0 \tilde{\gamma} + u_{it}. \quad (1.20)$$

Formally, Bartik and industry shares as instruments are only equivalent if the coefficient on Z_{i2}^0 ($\tilde{\gamma}$) can be transformed into γ : $\gamma = \tilde{\gamma} / \Delta_{G,t}$. This will only hold generally if $\tilde{\gamma}$ is time-varying. Intuitively, Bartik allows a high service share to predict high growth in one period, but low growth in a different period, whereas the simplest form of using industry shares as instruments forces a high service share to always predict high (or low) growth.

To recover the equivalence between Bartik and using shares as instruments in the panel setting, we can interact industry shares with time fixed effects. Anticipating a point that we return to later, this approach means having many more instruments than in Bartik. Bartik is a single instrument whereas interacting T time periods and K industries gives $T \times K$ instruments.

Since the industry shares are time-invariant, the approach is similar to a differences-in-differences estimator. Here, the size of the policy is measured by the dispersion in national industry growth, $\Delta_{G,t} = \bar{G}_{2t} - \bar{G}_{1t}$, and the exposure to the policy is given by Z_{i2}^0 . And, because the industry shares are time invariant, $\tilde{\gamma}_t$ can only be estimated relative to a base period. This gives the familiar estimating equation:

$$Y_{it} = \alpha_i + \alpha_t + X_{it}\beta + \epsilon_{it} \quad (1.21)$$

$$X_{it} = \tau_i + \underbrace{(\tau_t + \bar{G}_{1t}\gamma)}_{\tilde{\tau}_t} + Z_{i2}^0 \underbrace{\Delta_{G,t}\gamma}_{\tilde{\gamma}_t} + u_{it} \quad (1.22)$$

$$= \tau_i + \tilde{\tau}_t + Z_{i2}^0 \underbrace{\Delta_{G,t}\gamma}_{\tilde{\gamma}_t} + u_{it} \quad (1.23)$$

$$= \tau_i + \tilde{\tau}_t + \sum_{s \neq 0} Z_{i2}^0 \tilde{\gamma}_s 1(t = s) + u_{it}. \quad (1.24)$$

This setup is familiar to the now-standard difference-in-difference setup with continuous treatment exposure, and non-parametrically estimated exposure. Cross-sectionally, some groups are exposed more or less to a treatment, which is measured by Z_{i2}^0 . to test for parallel trends.

REMARK 1.4. *In a panel setting, to maintain the analogy to difference-in-difference, the regression should include place and time fixed effects. A key difference compared to the standard setting is that a typical DinD will have a well-defined “shock” that treats one group vs. the other in one time period going forward, while the Bartik instrument generates a continuously varying treatment that can change over the sample.⁷ For example, there is a shock in 1970-1980, 1980-1990, and etc., of varying size across groups. This continuous variation in the treatment makes testing for parallel trends less direct. We return to this issue in section 4.4.*

⁷For example, a simple case would be a change in trade tariffs in period $t = 0$: manufacturing is more affected than services. Hence, for $t > 0$, $\tilde{\gamma}_s$ would be positive, and would be a valid instrument so long as industry shares affected Y_{it} only through X_{it} .

1.3 Examples

We now discuss how several papers map into our notation. This exercise also emphasizes the flexibility that Bartik provides researchers in constructing $\hat{\mathbb{E}}[\mathbf{G}|\mathbf{Z}]$. Nonetheless, it remains the case that identification comes from the exogeneity of \mathbf{Z} .

Autor and Duggan (2003) use Bartik in the canonical way, to construct a labor demand shock by interacting industry growth rates and a measure of national industry performance. Specifically, their X is state-specific employment growth, \mathbf{Z} is state-specific industry composition, and \mathbf{G} is state-specific industry employment growth. In place of \mathbf{G} , for $\hat{\mathbb{E}}[\mathbf{G}|\mathbf{Z}]$ they use the leave-one-out change in national industry *shares*.

Autor, Dorn, and Hanson (2013) use Bartik to construct regional variation in exposure to trade with China. Specifically, their X is commuting zone-specific increase in imports from China, \mathbf{Z} is commuting-zone specific industry composition (lagged by several years to avoid anticipation effects), and \mathbf{G} is commuting zone specific industry growth in imports from China. In place of \mathbf{G} , for $\hat{\mathbb{E}}[\mathbf{G}|\mathbf{Z}]$ they use the growth in imports to *other* high income countries from China. Note that this is an extreme version of a leave-one-out estimator in that they use no U.S. information to construct it.

Greenstone, Mas, and Nguyen (2015) use Bartik to construct regional variation in shocks to the supply of credit during the Great Recession. Specifically, their X is credit growth in a county, \mathbf{Z} is the county specific composition of bank lending, and \mathbf{G} is the county-specific growth of lending of a bank. In place of \mathbf{G} , for $\hat{\mathbb{E}}[\mathbf{G}|\mathbf{Z}]$ they use a residualized measure of bank lending growth, which partials out county fixed effects.

Further afield, we note that the Currie and Gruber (1996b) and Currie and Gruber (1996a) simulated instrument is also encompassed by our framework. To make this link, note that rather than k indexing *industries*, think about k as indexing one of K discrete *types*, where the types are defined by eligibility criterion of the policy. Then X is the change in the share of the population eligible for Medicaid. Slightly different than the Bartik setting, \mathbf{Z} is the *change* in Medicaid eligibility rules in each state for each type. Finally, \mathbf{G} is the share of the population in each state that is of each type. In this case, for $\hat{\mathbb{E}}[\mathbf{G}|\mathbf{Z}]$ they use the national shares of each type of person. What this analogy highlights is the the key identifying assumption is that the changes in Medicaid policy are exogenous.

2 Many industries

We now present the case with multiple industries and time periods. As in the two industry case, our goal is to show the relationship between Bartik and using industry shares as instruments. Let there be K industries, T time periods and N locations, with k, t , and i denoting a particular industry, time or location. Extending the two industry case, we now

have that employment growth in location i at time t is given by:

$$X_{it} = \sum_{k=1}^K G_{ikt} Z_{ikt} \quad (2.1)$$

where G_{ikt} is the location-industry-time growth rate, and Z_{ikt} are the industry shares such that $\sum_{k \in \mathcal{I}} Z_{ikt} = 1, \forall i, t$.

We first derive an expression of the estimator using the Bartik instrument. Define

$$\hat{B}_{it} = \sum_{k \in \Sigma} \bar{G}_{kt} Z_{ikt}, \quad (2.2)$$

where $\bar{G}_{kt} = N^{-1} \sum_i G_{ikt}$ is the average industry growth rate for industry k in time period t .⁸ The two stage least squares system of equations is:

$$Y_{it} = W_{it}\alpha + X_{it}\beta + \epsilon_{it} \quad (2.3)$$

$$X_{it} = W_{it}\tau + \hat{B}_{it}\gamma + u_{it}, \quad (2.4)$$

where W_{it} is a $1 \times L$ vector of controls. Typically in a panel context, W_{it} will include location and year fixed effects, while in the cross-sectional regression, this will simply include a constant. It may also include a variety of other variables. Let $n = N \times T$, the number of location-years. For simplicity, let \mathbf{Y}_n denote the $n \times 1$ stacked vector of Y_{it} , \mathbf{W}_n denote the $n \times L$ stacked vector of W_{it} controls and \mathbf{X}_n denote the $n \times 1$ stacked vector of X_{it} and $\hat{\mathbf{B}}_n$ denote the stacked vector of \hat{B}_{it} . Denote $\mathbf{P}_W = \mathbf{W}_n(\mathbf{W}_n' \mathbf{W}_n)^{-1} \mathbf{W}_n'$ as the $n \times n$ projection matrix of \mathbf{W}_n , and $\mathbf{M}_W = \mathbf{I}_n - \mathbf{P}_W$ as the annihilator matrix. Then, because this is an exactly identified instrumental variable our estimator is

$$\hat{\beta}_{\text{bartik}} = \frac{(\mathbf{M}_W \mathbf{B}_n)' \mathbf{Y}_n}{(\mathbf{M}_W \mathbf{B}_n)' \mathbf{X}_n}. \quad (2.5)$$

We now consider the alternative approach of using industry shares as instruments. The two-equation system is:

$$Y_{it} = W_{it}\alpha + X_{it}\beta + \epsilon_{it} \quad (2.6)$$

$$X_{it} = W_{it}\tau + Z_{it}\gamma_t + u_{it}, \quad (2.7)$$

where Z_{it} is a $1 \times K$ row vector of industry shares, and γ_t is a $K \times 1$ vector, and, reflecting the lessons of previous section, the t subscript allows the effect of a given industry share to

⁸We return to leave-one-out below.

be time-varying. In matrix notation, we write

$$\mathbf{Y}_n = \mathbf{W}_n \alpha + \mathbf{X}_n \beta + \epsilon_n \quad (2.8)$$

$$\mathbf{X}_n = \mathbf{W}_n \tau + \tilde{\mathbf{Z}}_n \gamma_T + u_n, \quad (2.9)$$

where γ_T is a stacked $1 \times (T \times K)$ row vector such that

$$\gamma_T = [\gamma_1 \cdots \gamma_T], \quad (2.10)$$

and $\tilde{\mathbf{Z}}_n$ is a stacked $n \times (T \times K)$ matrix such that

$$\tilde{\mathbf{Z}}_n = \begin{bmatrix} \mathbf{Z}_n \odot \mathbf{1}_{t=1} & \cdots & \mathbf{Z}_n \odot \mathbf{1}_{t=T} \end{bmatrix}, \quad (2.11)$$

where $\mathbf{1}_{t=T}$ is an $n \times K$ indicator matrix equal to one if the n th observation is in period t , and zero otherwise. \odot indicates the Hadamard product, or pointwise product of the two matrices. Let $\mathbf{Z}_n^\perp = \mathbf{M}_W \tilde{\mathbf{Z}}_n$ and $\mathbf{P}_{Z^\perp} = \tilde{\mathbf{Z}}_n^\perp (\tilde{\mathbf{Z}}_n^{\perp'} \tilde{\mathbf{Z}}_n^\perp)^{-1} \tilde{\mathbf{Z}}_n^{\perp'}$. Then, the 2SLS estimator is

$$\hat{\beta}_{2SLS} = \frac{\mathbf{X}_n' \mathbf{P}_{Z^\perp} \mathbf{Y}_n}{\mathbf{X}_n' \mathbf{P}_{Z^\perp} \mathbf{X}_n}. \quad (2.12)$$

Alternatively, using the \mathbf{Z}_n as instruments, the GMM estimator is:

$$\hat{\beta}_{GMM} = \frac{\mathbf{X}_n' \mathbf{M}_W \tilde{\mathbf{Z}}_n \Omega^{-1} \tilde{\mathbf{Z}}_n' \mathbf{M}_W \mathbf{Y}_n}{\mathbf{X}_n' \mathbf{M}_W \tilde{\mathbf{Z}}_n \Omega^{-1} \tilde{\mathbf{Z}}_n' \mathbf{M}_W \mathbf{X}_n}, \quad (2.13)$$

where Ω is a $(K \times T) \times (K \times T)$ weight matrix.

We now turn to the relationship between Bartik and using industry shares as instruments. To show the relationship between GMM and Bartik, it is helpful to write \mathbf{B}_n in terms of the underlying industry shares and growth rates. Define a row vector $\tilde{\mathbf{G}}$ which has an analogous structure to γ_T . We do this in two steps. First, let \bar{G}_t be a $1 \times K$ row vector where the k^{th} entry is \bar{G}_{kt} . Second, stack the \bar{G}_t to have the following row vector

$$\tilde{\mathbf{G}} = [\bar{G}_1 \cdots \bar{G}_T], \quad (2.14)$$

where this vector is $1 \times (K \times T)$. Then

$$\mathbf{B}_n = \tilde{\mathbf{Z}}_n \tilde{\mathbf{G}}'. \quad (2.15)$$

PROPOSITION 2.1. *If $\Omega^{-1} = \tilde{\mathbf{G}}' \tilde{\mathbf{G}}$, then $\hat{\beta}_{GMM} = \hat{\beta}_{bartik}$.*

Proof. See Appendix A. □

This proposition shows that Bartik is equivalent to doing GMM with industry shares as instruments when the inverse of the weight matrix is given by $\tilde{\mathbf{G}}' \tilde{\mathbf{G}}$.

Our framework also shows why—starting with Autor and Duggan (2003)—the literature has often used a leave-one-out estimator of the industry growth rates to construct the Bartik instrument. In our notation,

$$\mathbf{B}_n^{lo} = \frac{N}{N-1} \tilde{\mathbf{Z}}_n \tilde{\mathbf{G}}' - \frac{1}{N-1} \mathbf{X}_n, \quad (2.16)$$

where this expression is the same as \mathbf{B}_n except that the own-location growth rate is subtracted off. Specifically, it is possible to derive an expression for the finite sample bias that emerges from including own-location in the computation of the growth rates.⁹

PROPOSITION 2.2. *In finite samples, the difference between the Bartik estimator of β and the true β is given by:*

$$\hat{\beta}_{bartik} - \beta = \frac{\tilde{\mathbf{G}} \tilde{\mathbf{Z}}_n' \mathbf{M}_W \epsilon_n}{\tilde{\mathbf{G}} \tilde{\mathbf{Z}}_n' \mathbf{M}_W \mathbf{X}_n} \quad (2.17)$$

where

$$\tilde{\mathbf{G}} \tilde{\mathbf{Z}}_n' \mathbf{M}_W \epsilon_n = N^{-1} \sum_i Z_i G_i' M_{ii} \epsilon_i. \quad (2.18)$$

Proof. See Appendix A. □

Since G_i and ϵ_i are correlated, there will be bias in estimation when N is not sufficiently large. The leave-one-out mean estimator solves this problem directly, and is analogous to the jackknife estimator in JIVE (Angrist, Imbens, and Krueger (1999)).

3 Estimation

So far we have established that using the Bartik *instrument* is equivalent to using industry shares as an instrument. Hence, what is distinctive about the Bartik instrument is not as an *instrument* per se, but as an *estimation approach* to dealing with a high-dimensional first stage when the endogenous variable has an inner product structure. Put differently, to call it the Bartik instrument is a slight misnomer, and it in fact could be called the Bartik approach.

We first explain the sense in which using industry shares as instruments generates a high-dimensional first stage. In the traditional two-stage estimator set-up, the generic first-stage estimation problem is to estimate $\mathbb{E}[X|\mathbf{Z}]$. If a researcher estimates this conditional

⁹Note that in the presence of leave-one-out, the equivalence between Bartik and GMM no longer holds.

expectation using least squares with a vector of industry shares in the first-stage, the estimation would be very noisy if \mathbf{Z} is high-dimensional and might lead to bias in the two-stage estimator. For example, a state-level analysis might have 50 (or 51) states and there are 312 4-digit NAICS industries. That is, in a single cross-section there would be more instruments than observations. In this case, even once the number of instruments is reduced to match the number of states, the high degree of over-identification instruments would generate substantial bias. As we emphasized above, moving to a panel setting does nothing to alleviate this problem because the industry shares are interacted with period dummies and so the number of instruments scales with the number of periods.

We now explain how Bartik uses the inner product structure of the endogenous variable to do dimension reduction. Given the inner product structure of X in the Bartik setting, the first-stage estimation problem can be rewritten as follows:

$$\mathbb{E}[X|\mathbf{Z}] = \mathbb{E}[\mathbf{Z}'\mathbf{G}|\mathbf{Z}] = \mathbf{Z}'\mathbb{E}[\mathbf{G}|\mathbf{Z}].$$

In this case, instead of estimating $\mathbb{E}[X|\mathbf{Z}]$, a researcher would need to estimate $\mathbb{E}[\mathbf{G}|\mathbf{Z}]$ and then pre-multiply it by \mathbf{Z} . This estimation problem is still very high-dimensional. The Bartik instrument avoids the high-dimensional first-stage by making the following approximation:

$$\mathbb{E}[\mathbf{G}|\mathbf{Z}] \approx \mathbb{E}[\mathbf{G}].$$

Then, and to return to the 4-digit industry example, rather than having 312 instruments, the researcher constructs a single instrument:

$$B = \mathbf{Z}'\hat{\mathbb{E}}[\mathbf{G}].$$

In this case, $\hat{\mathbb{E}}[\mathbf{G}]$ is a valid estimator for $\mathbb{E}[\mathbf{G}|\mathbf{Z}]$, but just not necessarily the most efficient.

Substantively, this Bartik approximation means that the expectation of the growth rate of a particular industry does not depend on a location's industry composition. It would be consistent with the identifying assumption of Bartik to allow for inter-industry spillovers through input-output linkages, i.e., the restaurant industry might have a different growth rate in locations with a large and small entertainment industry presence because restaurant industry demand comes primarily from the entertainment industry. Similarly, it is consistent with the Bartik assumption to allow for the effects of a shock to depend on levels, i.e., it might be that there is curvature in the location-industry production function and a given shock has a larger (or smaller) effect in locations where the industry is more prominent.

There are potentially many ways to efficiently estimate $\mathbb{E}[\mathbf{G}|\mathbf{Z}]$. Surprisingly, however, the Bartik estimator does extremely well in our application compared to other potential alternatives, in large part due to the fact that industry means explain almost 50% of the

overall location-industry variance (see below for more details). Hence, simple means are a very reasonable approximation. In unreported results, we attempted other approximations to $\mathbb{E}[\mathbf{G}|\mathbf{Z}]$ but none worked any better than Bartik.

Since this process constructs a generated instrument to be used in a just-identified two-stage estimation procedure, there is no impact on the asymptotic distribution of the estimator. See Section 6.1.2 in Wooldridge (2002, pg. 117). This implies that under standard asymptotics, any improvement in estimating $\mathbb{E}[X|\mathbf{Z}]$ by estimating $\mathbb{E}[\mathbf{G}|\mathbf{Z}]$ more efficiently will only improve our results in finite samples.

3.1 Industry and location bins: theoretical considerations

An important issue that arises in practice is how finely to divide industry and location bins. So far we have assumed a sampling frame in which it is not possible to discuss this question. The reason is that we have assumed that each location contains information on all industries. Hence, implicitly, as we divide industries and locations more finely, we get additional data on the additional industries and locations (and they are all independent), so there is no reason to not divide them arbitrarily finely. Here we consider an alternative sampling frame where it is possible to discuss this issue.

To be able to discuss the choice of how finely to divide locations (the set-up for industries would be analogous), we consider a sampling frame where the number of locations with a given industry is fixed as the sample size grows. We can motivate this frame by imagining that there are a fixed number of firms in a particular industry, so as we divide locations more finely we get no more information about the industry growth rates. Formally, let there be N locations, T time periods, and K industries. For each industry k , there are N_k locations where growth rates are observed. A simple way to envision this is that there are a number of individuals who are employed, and they are dispersed across locations in finite quantity. Hence, as N grows with K fixed, there are more and more cities with individuals in industry k , and so $N_k \rightarrow \infty$.

We can use this setup to understand the problems that emerge with picking excessively fine location bins. In this setup, we define our Bartik instrument as $B_i = \mathbf{Z}_i' \bar{\mathbf{G}}$, with $\bar{\mathbf{G}}_k = N_k^{-1} \sum_j G_{jk}$. Define $G_{ik} = \mu_k + \epsilon_{ik}$, where $\mu_k = E(G_{ik})$, and let $\bar{\mathbf{G}}_k = \mu_k + N_k^{-1} \sum_j \epsilon_{jk}$. Let $u_k = N_k^{-1} \sum_j \epsilon_{jk}$ and note that for fixed N_k , u_k is mean zero with non-zero variance. Then,

$$B_i = \sum_k \mathbf{Z}_{ik} \mu_k + \sum_k \mathbf{Z}_{ik} u_k. \quad (3.1)$$

Let μ_G be the vector of μ_k and let \mathbf{u} be the vector of u_k .¹⁰ Note that since u_k is independent

¹⁰We are slightly abusing notation here, as we will assume that $\bar{\mathbf{G}}_k$ uses the leave-one-out mean and hence it should be indexed by i as well. We do this for clarity's sake, and the results should still hold under these

of Z_{ik} , we can write the variance as

$$\text{Var}(B_i) = \mu'_G \Sigma_Z \mu_G + \text{Var}(Z'_i \mathbf{u}). \quad (3.2)$$

Then:¹¹

$$\text{Var}(Z'_i \mathbf{u}) = \iota'_K (\Sigma_Z \odot \Sigma_u) \iota_K + \mu'_Z \Sigma_u \mu_Z. \quad (3.8)$$

The notable feature about this is that while the second term is finite, the first term grows as $K \rightarrow \infty$ if Σ_u stays fixed. Hence, if $N_k \not\rightarrow \infty$, then $\text{Var}(Z'_i \mathbf{u}) \rightarrow \infty$. This will destroy the power in the first stage, since the first stage is effectively trading off between $\sum_k Z'_{ik} \mu_k$ as a fraction of $\text{Var}(B_i)$.

We explore this further in Section 5 when we perform simulations.

4 Testing for confounds

So far we have emphasized that the key assumption of Bartik is the exogeneity of industry shares. We now present some simple descriptive results that illustrate how one might go about probing this assumption. The key challenges are that industry shares are high-dimensional, and that Bartik allows for a new shock in every period so that testing for pre-trends is nonstandard. While Bartik is used for many questions (and the generic approach is used in many setting besides industry-location), we focus on variables related to labor supply, since this is the canonical application.

4.1 Dataset

We use the 5% sample of IPUMS (Ruggles et al. (2015)) for 1980, 1990 and 2000 and we pool the 2009-2011 ACSs for 2010. We look at PUMAs and 3-digit IND1990 industries. In Appendix C we show all our results with states and 2 digit industries and results are very

assumptions.

11

$$\text{Var}(Z'_i \mathbf{u}) = E(\mathbf{u}' Z_i Z'_i \mathbf{u}) - E(\mathbf{u}' Z_i) E(Z'_i \mathbf{u}) \quad (3.3)$$

$$= E(\mathbf{u}' Z_i Z'_i \mathbf{u}) - E(\mathbf{u})' E(Z_i) E(Z_i)' E(\mathbf{u}) \quad (3.4)$$

$$= E(\mathbf{u}' Z_i Z'_i \mathbf{u}) \quad (3.5)$$

$$= \sum_{k,l} E(u_k u_l Z_k Z_l) \quad (3.6)$$

$$= \iota'_K (\Sigma_Z \odot \Sigma_u) \iota_K + \mu'_Z \Sigma_u \mu_Z. \quad (3.7)$$

similar.¹² To construct the remaining aspects of our dataset, we follow Autor and Duggan (2003). In the notation given above, our y variable is earnings growth, and X is employment growth. We use people aged 18 and older who report usually working at least 30 hours per week in the previous year. We fix industry shares at the 1980 values, and then construct the Bartik instrument using 1980 to 1990, 1990 to 2000 and 2000 to 2010 leave-one-out growth rates.

4.2 Why might Bartik work?

Our theoretical results emphasize that the Bartik instrument uses cross-industry variation in national growth rates. And, if there is no cross-industry variation, then the instrument has no power. A simple variance decomposition of the industry-location growth rates shows that indeed there is a national component of the industry growth rates.

Consider the following expression for the industry-location growth rates in location i and industry k in a particular time period:

$$g_{ik} = g_k + g_i + \epsilon_{ik}. \quad (4.1)$$

We can operationalize this equation as a regression where we include location and industry fixed effects. By computing the covariance of each estimated component with the overall growth rate, we can then ask how much of the variance of the industry-location growth rates reflect the common industry component, how much is a location component, and how much is in the interaction.¹³

Table 2 provides evidence that there is indeed a common industry component to growth rates. At the 3 digit level and using PUMAs we find that the industry component explains 15-20% of the variance of growth rates. Interestingly, the location component is quite small—explaining less than 5% of the variance of the industry location growth rates. Table A2 shows analogous statistics at the 2-digit level and using locations, and finds a similarly small role for location in explaining the industry-location growth rates, but a substantially larger for industry (the industry component is about 45% of the variance).

The table also emphasizes a point that we return to in our simulations: there are many zeros in the industry-location growth rates—at the 3 digit level and with PUMAs, about a quarter of the industry-location level observations are zeros. In contrast, Table A2 shows that at the state and 2 digit industry level less than 5% of the state-industry observations

¹²We have also looked at the other 2 possible combinations of location and industry, and results are again quite similar. There are 244 3-digit IND1990 industries and 91 such 2 digit industries. There are 543 PUMAs.

¹³Formally, $\frac{Cov(g_{ik}, \hat{g}_k)}{Var(\hat{g}_k)}$ (the industry share); $\frac{Cov(g_{ik}, \hat{g}_i)}{Var(\hat{g}_i)}$ (the location share); and $\frac{Cov(g_{ik}, \hat{\epsilon}_{ik})}{Var(\hat{\epsilon}_{ik})}$ (the residual share). We have also considered versions where we compute the industry share as the national leave-one-out mean and compute the g_k component directly, and get similar answers.

are zeros.

4.3 Balance

While the previous section showed that there is an important national component to industry growth rates, our theoretical results emphasized that identification in Bartik comes from the exogeneity of industry shares. A natural way to begin to probe the plausibility of this assumption is to examine the relationship between industry composition and observable characteristics of a place.

Table 3 shows that observables in both levels and changes are closely related to the Bartik projection of industry composition. The left-hand side of the table shows the relationship between these three different values of Bartik and *levels* of observable characteristics of the PUMAs. The right-hand side of the table shows these relationships for *changes* in observable characteristics of places. Observable characteristics explain a large share of the variance in the Bartik instrument. In the first period, the R^2 is over 0.6. The characteristic that is most consistently related to the Bartik instrument is education. Thus, any trend that is correlated with these characteristics—for example, skill biased technological change, immigration, or rising female labor force participation—will be correlated with Bartik.

The relationship between the Bartik instrument and the observable characteristics is because industry composition is related to observable characteristics, and not because of properties of the growth rates. In Table 4, we show the results of an analogous exercise where we do dimension reduction on the 1980 industry shares using principal component analysis, which does not use any information in the subsequent growth rates. We then relate the first principal component of 1980 industry shares to time-varying observable characteristics of a location. We find, if anything, a tighter relationship between the first principal component of the 1980 industry shares and characteristics. Notably, the relationship between the levels and the characteristics does not fade over time.

There are two broad classes of reactions to this observation. First, if we think that we live in a selection on observables world, then we have measured the observables and so we can control for them. Second, we might think that we live in a selection on unobservables world, and then worry that the extent of selection on observables suggests how much we should worry about selection on unobservables.

Table 5 pursues the selection on observables logic and reports the results of the IV estimates of the inverse labor supply elasticity with and without controlling for observables. The main take-away is that the inverse elasticity estimates are sensitive to the inclusion of controls for observable differences of locations. The table shows the results of pooling data for three ten-year changes (to 1990, 2000 and 2010). The benchmark IV estimate of the inverse elasticity of labor supply is shown in column (6) and is 1.08. Columns (7) through (9)

report results when we control for observable characteristics. Controlling for observables attenuates estimates by over 20% (from 1.08 to 0.82). Even though the regression includes location fixed effects, it is when we control for *levels* of observable characteristics rather than *changes* that the attenuation occurs.

4.4 Pre-trends

Besides balance, another natural way to explore the validity of an instrument is to consider pre-trends. While looking at balance forced us to focus on observables, examining pre-trends allows us to say something the relationship between unobservables and Bartik. A notable feature of Bartik is that it allows for a new shock in every period so that it implies that we do not expect parallel trends to hold. Here we develop a simple procedure to test for parallel trends even in the face of a time-varying instrument and show that there is evidence of pre-trends using the Bartik instrument.

4.4.1 Why parallel trends might not hold, even if Bartik is a valid instrument

It is consistent with the validity of the Bartik instrument to find evidence of pre-trends. To see this, suppose we have two periods of data, $t = \{1, 2\}$ and the same set-up we have been using:

$$y_{t,i} = \alpha_0 + \beta X_{t,i} + \epsilon_{t,i} \quad (4.2)$$

$$X_{t,i} = \alpha_1 + \gamma B_{t,i} + \nu_{t,i}, \quad (4.3)$$

where $B_{t,i}$ is a valid instrument. Note that y is already in changes. I.e., y might be wage growth. Testing for pre-trends then amounts to asking whether $\text{Corr}(y_1, B_2) = 0$. That is, do places with a higher period 2 instrument have faster (or slower) wage growth in period 1?

To see why there might be pre-trends even if the Bartik instrument is valid, note that we can write:

$$\text{Cov}(y_1, B_2) = \text{Cov}(\alpha_0 + \beta\alpha_1 + \beta\gamma B_1 + \beta\nu_1 + \epsilon_1, B_2). \quad (4.4)$$

Hence, $\text{Cov}(y_1, B_2)$ can be nonzero if $\text{Cov}(B_1, B_2)$ is nonzero. That is, if the Bartik instrument is correlated through time (because, for example, industry growth rates are correlated through time), then we will find evidence of pre-trends.

4.4.2 Adjusting for mechanical pre-trends

Mechanically, we ask whether the residuals from the second stage in the current period can be predicted by the values of Bartik in a future period. That is, we remove the part of wage growth that we would predict from the Bartik instrument. Formally, we compute

$$\tilde{y}_{1,i} = y_{1,i} - \beta\gamma B_1. \quad (4.5)$$

Then this adjustment purges the most mechanical reason for pre-trends, and we have:

$$Cov(\tilde{y}_1, B_2) = Cov(\alpha_0 + \beta\alpha_1 + \beta\nu_1 + \epsilon_1, B_2). \quad (4.6)$$

4.4.3 Evidence from Bartik

Column (1) of Table 6 shows that there is evidence of pre-trends. Namely, we pool wage growth from 1980 to 1990 and 1990 to 2000, and regress it on Bartik constructed one period forward, so the Bartik constructed using 1990 to 2000, and 2000 to 2010. Column (1) shows that we can predict past wage growth using future values of the instrument.

There is reason to think that some of this relationship might be mechanical. Namely, the correlation between adjacent period values of the Bartik instrument is 0.4272. Hence, it might be that the future value of the instrument is simply correlated with the past good shocks that led to the wage growth.

Columns (2) through (5) of Table 6 show that after addressing the mechanical reason for correlation that values of the Bartik instrument are correlated over time, we still find evidence of pre-trends. The columns correspond to the residualization in columns (6)-(9) of table 5 and show that for various ways of residualizing earnings growth for that predicted by the Bartik instrument we can still predict past values of wage growth using future values of the Bartik instrument. Two aspects of the table are quantitatively notable. First, controlling for observable characteristics of the location does reduce the magnitude of the implied pre-trends. Second, however, the size of the coefficient is still large. For example, the reduced-form for the effect of Bartik on wage growth for the specification in column (5) is 0.2624.¹⁴ Hence, the coefficient in column (5) of 0.0403 of future values of Bartik on past values of wage growth is large.

5 Simulations

We now show a number of simulations that illustrate a few of the points we have made. Appendix B provides details on the simulations and Table A1 shows the parameters. The

¹⁴This multiplies the main effect in in column (5) and column (9) of Table 5.

baseline parameters are chosen such that Bartik “works.”

Table 7 summarizes the interesting simulations. The first point that emerges from considering the various simulations is the role of the relative number of industries and locations. In the baseline simulation, we have 300 locations and 10 industries, that is, the ratio of locations to industries is 30. In this case, IV appears to be median and mean unbiased. When we drop the number of locations to 20 (so that the ratio of locations to industries is 2), we find that IV appears to be median and mean biased. Figure 1 shows the continuous version of this point. A similar way of showing this point is to hold the number of locations constant and raise the number of industries. Figure 2 shows that as we increase the number of industries the mean estimate of θ drifts down, while the variance across simulation runs increases. Table 7 shows that when we reach 200 industries (so that the ratio of locations to industries is 1.5), that IV is median and mean biased.

In these simulations, it appears that when the ratio of locations to industries approaches 2 or 3 that there is bias, while in our empirical work above we considered 543 locations and 244 industries (and in the appendix we considered 51 states and 91 industries). While we have not yet designed a simulation that we feel fully captures the empirical setting, this suggests that the fine-ness of industry and location bins can play a large role.

The other interesting point that emerges from the simulations is that reducing the variance of common industry component has large effects. The table shows the effect of dropping the variance by a factor of 7 (from $\sigma_k^2 = 7$ to $\sigma_k^2 = 1$). These simulations leave classic traces of a weak instrument, in the sense that the IV estimates are incredibly unstable (the 2.5th to 97.5th percentile of $\hat{\theta}$ across the simulations is -7.9 to 9.0). Figure 3 provides the continuous version of these results.

We have also explored how varying other dimensions of our baseline simulation affect results, but there are no notable or interpretable results. (See figures A1 to A3).

To understand the identification result, we consider three alternative simulations. First, having the industry common component in the error term does not necessarily constitute a problem. We draw a set of random vectors R_l in a way analogous to the Z , interact these with G_k , and enter them in the error term. The row titled “ G_k in ϵ (R)” shows that having the industry growth rates in the first stage error does not by itself constitute a problem. Second, having some function of the industry shares in the error does constitute a problem for identification. We can show this in a couple ways. By analogy to the previous simulation, we take the inner product of the local industry shares and the national industry components and enter them in the error term. The next row titled “ G_k in ϵ (Z)” shows that this leads to enormous bias in IV. We also generate a random vector, ϵ_K , which is unrelated to the G_k , interact this vector with the industry shares and enter it in the error term. This

also leads to substantial bias.¹⁵

6 Summary: recommendations for practice

This note develops a formal econometric structure to study the Bartik instrument. Developing such a structure is valuable because the Bartik instrument is widely used but not fully understood. We summarize this paper by the implications for practice that our structure delivers:

- The argument about exogeneity is in terms of the industry shares, and not growth rates.
- Use the initial period shares – don’t update.
- Make sure to use leave-one-out means.
- Check to see how “filled” industries are – don’t go too fine-grained in industry cuts if there aren’t many cities covering the means. Check this. You want $N_k \propto N$.
- Balance tests– check for confounders using initial composition.
- Try looking for pre-trends after partialling out the direct effects of previous values of Bartik.

We wish to emphasize that while some of these recommendations may seem obvious parts of the applied microeconomics toolkit, we are struck by how rarely, if at all, they are used in the context of Bartik instruments.

¹⁵Admittedly, to get this to work, we need to make the element-by-element variance quite large: 70.

References

- Altonji, Joseph G. and David Card. 1991. "The Effects of Immigration on the Labor market Outcomes of Less-skilled Natives." In *Immigration, Trade and the Labor Market*, edited by John M. Abowd and Richard B. Freeman. University of Chicago Press, 201–234.
- Angrist, Joshua D., Guido W. Imbens, and Alan B. Krueger. 1999. "Jackknife Instrumental Variables Estimation." *Journal of Applied Econometrics* 14:57–67.
- Autor, David, David Dorn, and Gordon Hanson. 2013. "The China Syndrome: Local Labor Market Effects of Import Competition in the United States." *American Economic Review* 103 (6):2121–2168.
- Autor, David and Mark Duggan. 2003. "The Rise in the Disability Rolls and the Decline in Unemployment." *Quarterly Journal of Economics* 118 (1):157–205.
- Bartik, Timothy. 1991. *Who Benefits from State and Local Economic Development Policies?* W.E. Upjohn Institute.
- Blanchard, Olivier and Lawrence Katz. 1992. "Regional Evolutions." *Brookings Papers on Economic Activity* 1992 (1):1–75.
- Bound, John and Harry J. Holzer. 2000. "Demand Shifts, Population Adjustments, and Labor Market Outcomes during the 1980s." *Journal of Labor Economics* 18 (1):20–54.
- Chodorow-Reich, Gabriel. 2014. "The Employment Effects of Credit Market Disruptions: Firm-Level Evidence From the 2008-9 Financial Crisis." *Quarterly Journal of Economics* 129 (1):1–59.
- Currie, Janet and Jonathan Gruber. 1996a. "Health Insurance Eligibility, Utilization, Medical Care and Child Health." *Quarterly Journal of Economics* 111 (2):431–466.
- . 1996b. "Saving Babies: The Efficacy and Cost of Recent Changes in the Medicaid Eligibility of Pregnant Women." *Journal of Political Economy* 104 (6):1263–1296.
- Greenstone, Michael, Alexandre Mas, and Hoai-Luu Nguyen. 2015. "Do Credit Market Shocks affect the Real Economy? Quasi-Experimental Evidence from the Great Recession and 'Normal' Economic Times." Working paper.
- Nakamura, Emi and Jon Steinsson. 2014. "Fiscal Stimulus in a Monetary Union: Evidence from US Regions." *American Economic Review* 104 (3):753–792.

Ruggles, Steven, Katie Genadek, Ronald Goeken, Josiah Grover, and Matthew Sobek. 2015. *Integrated Public Use Microdata Series: Version 6.0 [Machine-readable database]*. Minneapolis: University of Minnesota.

Wooldridge, Jeffrey M. 2002. *Econometric Analysis of Cross Section and Panel Data*. MIT Press.

Table 1: Literature

Aizer (2010, AER)	Domestic Violence
Allcott and Kenniston (2014, WP R& R at REStud)	
Altonji and Card (1991, NBER Chapter)	Immigration
Autor and Duggan (2003, QJE)	Public, Disability insurance
Autor, Dorn and Hanson (2013, AER)	Trade, local labor markets
Autor, Dorn, Hanson and Song (2014, QJE)	Trade, local labor markets
Baum-Snow and Ferrreira (2015, Handbook of Urban and Regional)	Urban
Beaudry, Green and Sand (2012, ECMA)	Macro- Labor
Bertrand, Kamenica and Pan (2015, QJE)	
Blanchard and Katz (1992, BPEA)	Macro - Labor
Bloom, Draca and Van Reenen (2016, REStud)	Trade and Productivity
Bound and Holzer (2000, JoLE)	Labor
Brunner, Ross and Washington (2011, Restat)	Political economy
Cadena and Kovak (2016, AEJ:Applied)	Labor
Card (2001, JoLE)	Immigration
Card (2009, AER)	Immigration
Chodorow-Reich and Wieland (2016, WP)	Macro-labor, Reallocation
Davis and Haltiwanger (2014, WP– Jackson Hole)	
Diamond (2016, AER)	Urban/Public
Dinerstein, Hoxby, Meer and Villanueva (2014, NBER Chapter)	Education
Gould, Weinberg and Mustard (2002, REStat)	
Greenstone, Mas and Nguyen (2015, WP–R&R at AEJ:Policy)	Finance, Macro
Guerrieri, Hartley and Hurst (2013, JPubE)	Urban, Public
Hagedorn, Karahan, Manovskii and Mitman (2016, WP)	Public, Macro-Labor
Juhn and Kim (1999, JoLE)	
Kovak (2013, AER)	Trade
Lewis (2011, QJE)	
Lin (2011, REStat)	
Luttmer (2005, QJE)	Public
Moretti (2013, AEJ: Applied)	
Nakamura and Steinsson (2014, AER)	Macro, Fiscal Multipliers
Nekarda and Ramey (2011, AEJ: Macro)	Macro,
Notowidigdo (2013, WP)	
Oberfield and Raval (2014, WP – R&R at ECMA)	Macro and IO
Ottaviano, Peri and Wright (2013, AER)	Trade
Saiz (2010, QJE)	Urban
Saks and Wozniak (2011, JoLE)	
Suarez Serrato and Zidar (2016, AER))	Public Economics

Table 2: Variance decomposition of industry-location growth rates (Puma 3-Digit)

	Industry	Location	Residual	Share of zeros
1990	.1896	.0356	.7748	.2675
2000	.2015	.0255	.773	.2561
2010	.1352	.0183	.8465	.2911

Table 3: (Puma) 1980 3-Digit Industry Share Bartik Instrument

	Levels			Changes		
	1990	2000	2010	1990	2000	2010
Male	-0.31***	-0.24***	0.02	0.04	0.10*	-0.12**
White	-0.08	0.04	-0.04	-0.07	-0.13***	0.01
Native Born	-0.32***	-0.14**	-0.27***	-0.13***	0.10*	-0.09*
12th Grade Only	0.07	0.52***	0.20	-0.36***	-0.59***	-0.22**
Some College	0.69***	0.93***	0.95***	-0.53***	-1.14***	-1.28***
Veteran	0.29***	-0.14	0.37**	0.15***	-0.19***	0.10**
# of Children	0.05	0.03	0.26***	-0.04	0.18***	-0.22***
Total Income	-0.63*	0.11	-0.80***	0.56***	0.47***	0.78***
Social Security Income	0.17	-0.04	0.35***	-0.03	-0.01	-0.07
5-Year Same State	-1.14***	-0.03	-0.29	0.71**	0.23	
R^2	0.66	0.42	0.27	0.62	0.30	0.43
F	104.45	42.08	20.80	85.78	17.32	35.31
p	0.00	0.00	0.00	0.00	0.00	0.00

Standard errors in parentheses

* $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$

Table 4: (Puma) 1980 3-Digit Industry Share Principal Component 1

	Levels				Changes		
	1980	1990	2000	2010	1990	2000	2010
Male	-0.39***	-0.38***	-0.38***	0.02	0.24***	0.09*	
White	-0.18***	-0.11***	-0.10***	-0.22***	-0.22***	-0.08	
Native Born	-0.61***	-0.44***	-0.42***	-0.15**	-0.20***	-0.20***	
12th Grade Only	0.08*	0.29***	0.20***	-0.34***	-0.18	-0.90***	
Some College	0.26***	0.50***	0.46***	-0.11	-0.02	-1.05***	
Veteran	0.17*	-0.11	-0.09	-0.39***	-0.31***	-0.18***	
# of Children	-0.14***	-0.12***	-0.17***	-0.07	0.28***	0.06	
Total Income	1.91***	0.93***	0.54***	0.63***	0.26*	-0.13	
Social Security Income	-0.74***	-0.26***	-0.22***	0.01	0.10**	0.07*	
5-Year Same State	-0.21	-0.10	-0.05	0.36	0.76		
R^2	0.82	0.87	0.85	0.73	0.57	0.52	
F	213.22	259.03	229.78	121.45	66.23	59.52	
p	0.00	0.00	0.00	0.00	0.00	0.00	

Standard errors in parentheses

* $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$

Table 5: Across-Time Puma Regressions with 3-Digit Industry bartik_1980

	OLS		First Stages			Second Stages			
	Δ Wage (1)	Δ Emp (2)	Δ Emp (3)	Δ Emp (4)	Δ Emp (5)	Δ Wage (6)	Δ Wage (7)	Δ Wage (8)	Δ Wage (9)
Δ Emp	0.46***					1.08	0.82***	1.17***	0.79***
Bartik (1980)		0.38***	0.39***	0.27***	0.33***				
Male			-0.07		0.20*		-0.24***		-0.30***
White			0.22		0.16		0.22*		0.13
Native Born			-0.04		-0.20		0.15		0.13
12th Grade Only			0.47***		0.14		-0.36***		-0.11
Some College			0.63***		0.66***		-0.71***		-0.58***
Veteran			0.04		-0.20		0.15*		0.20**
# of Children			0.76***		0.23*		-0.42***		-0.34***
Δ Male				0.05	0.08*			-0.01	-0.05*
Δ White				0.01	0.01			-0.07	-0.03
Δ Native Born				-0.05	-0.04			0.20***	0.12***
Δ 12th Grade Only				-0.19***	-0.23***			0.18***	0.18***
Δ Some College				-0.18*	-0.06			0.24***	0.08
Δ Veteran				-0.13**	-0.15**			0.21***	0.17***
Δ # of Children				-0.26***	-0.27***			0.11*	0.00
Year FE	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes
Puma FE	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes
Observations	1,629	1,629	1,629	1,629	1,629	1,629	1,629	1,629	1,629
R-squared	0.90	0.73	0.77	0.79	0.80	0.78	0.88	0.78	0.90
F
p-value	0.00	0.00	0.00

Standard errors in parentheses

* $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$

Table 6: Regression of 1-Lagged Wage Growth Residuals on Bartik, Puma 3-Digit

	Lag Wage Growth (1)	Residualized			
		No Controls (2)	Levels (3)	Changes (4)	Levels+Changes (5)
bartik_1980	0.112** (0.0360)	0.0988*** (0.0125)	0.0558*** (0.0104)	0.0567*** (0.0106)	0.0403*** (0.0101)
<i>N</i>	1086	1086	1086	1086	1086

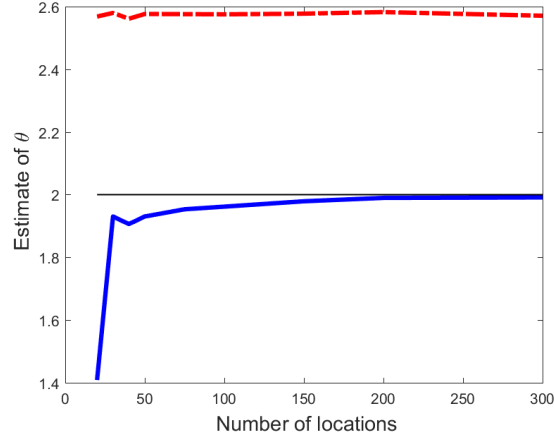
Standard errors in parentheses

* $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$

Table 7: Simulation results

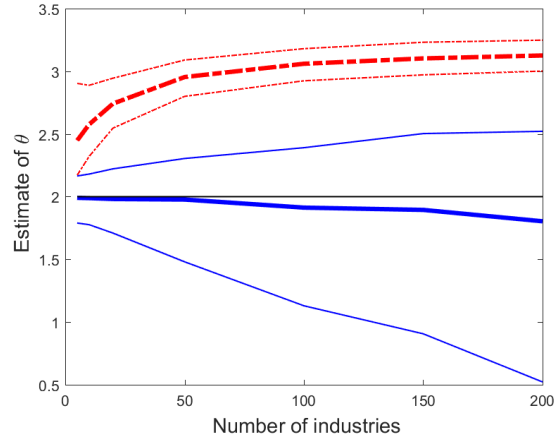
	OLS				IV			
	Mean($\hat{\theta}$)	Median ($\hat{\theta}$)	$\hat{\theta}_{2.5}$	$\hat{\theta}_{97.5}$	Mean($\hat{\theta}$)	Median ($\hat{\theta}$)	$\hat{\theta}_{2.5}$	$\hat{\theta}_{97.5}$
Benchmark	2.5760	2.5606	2.3321	2.8806	1.9927	2.0016	1.7686	2.1764
$L = 20$	2.5684	2.5472	2.0478	3.1621	1.4084	1.9365	-0.1614	2.5931
$K = 200$	3.1263	3.1262	3.0017	3.2485	1.8036	1.9053	0.5218	2.5212
$\sigma_k^2 = 1$	3.1030	3.1038	2.9735	3.2267	0.7719	1.8535	-7.8606	8.9783
G_k in ϵ (R)	2.5736	2.5638	2.2961	2.8884	1.9930	1.9992	1.7434	2.1963
G_k in ϵ (Z)	3.0628	3.0613	2.9821	3.1590	3.0008	3.0018	2.8482	3.1512
ϵ_k in ϵ (Z)	2.6803	2.6695	-0.8200	6.1420	2.2531	2.1527	-5.1831	10.6310

Figure 1: Bartik simulation: changing number of locations



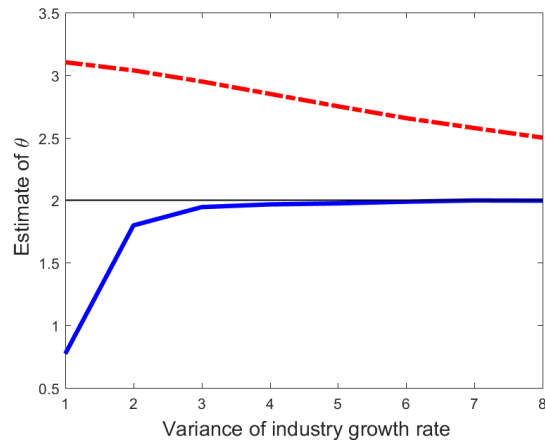
Notes: The blue line shows the TSLS estimates using Bartik. The red dashed line shows OLS. The thin black line shows the truth ($\theta = 2$). The thin lines show the 2.5 to 97.5th percentiles across the 1000 simulations at each point. The parameter values are as in Table A1, except that L varies. The benchmark number of locations is 300.

Figure 2: Bartik simulation: changing number of industries



Notes: The blue line shows the TSLS estimates using Bartik. The red dashed line shows OLS. The thin black line shows the truth ($\theta = 2$). The thin lines show the 2.5 to 97.5th percentiles across the 1000 simulations at each point. The parameter values are as in Table A1, except that K varies. The benchmark number of industries is 10.

Figure 3: Bartik simulation: changing variance of industry shocks



Notes: The blue solid line shows the TSLS estimates using Bartik. The red dashed line shows OLS. The thin black line shows the truth ($\theta = 2$). There are 1000 simulations at each point. The parameter values are as in Table A1, except that σ_k^2 varies. The benchmark value of σ_k^2 is 7.

A Omitted proofs

Proposition 1.1

Proof. Note that

$$\sum_i B_i Y_i = \sum_i (\bar{G}_1 Z_{i1} + \bar{G}_2 Z_{i2}) Y_i \quad (A1)$$

$$= \sum_i \bar{G}_1 Y_i + \sum_i (\bar{G}_2 - \bar{G}_1) Z_{i2} Y_i \quad (A2)$$

$$= \bar{G}_1 \sum_i Y_i + (\bar{G}_2 - \bar{G}_1) \sum_i Z_{i2} Y_i \quad (A3)$$

$$n^{-1} \sum_i B_i \sum_j Y_j = n^{-1} \left[\sum_i \bar{G}_1 \sum_j Y_j + \sum_i (\bar{G}_2 - \bar{G}_1) Z_{i2} \sum_j Y_j \right] \quad (A4)$$

$$= \bar{G}_1 \sum_j Y_j + n^{-1} (\bar{G}_2 - \bar{G}_1) \sum_i Z_{i2} \sum_j Y_j \quad (A5)$$

where the first and fourth line holds by definition, the second because $Z_{i1} + Z_{i2} = 1$, and the third and fifth due to the fact that \bar{G}_1 and \bar{G}_2 are constant across i . Hence,

$$\sum_i B_i Y_i - n^{-1} \sum_i B_i \sum_j Y_j = (\bar{G}_2 - \bar{G}_1) \sum_i Z_{i2} Y_i - n^{-1} (\bar{G}_2 - \bar{G}_1) \sum_i Z_{i2} \sum_j Y_j. \quad (A6)$$

It is easy to show the same argument for the denominator of $\hat{\beta}_{2SLS}(B)$ such that

$$\sum_i B_i X_i - n^{-1} \sum_i B_i \sum_j X_j = (\bar{G}_2 - \bar{G}_1) \sum_i Z_{i2} X_i - n^{-1} (\bar{G}_2 - \bar{G}_1) \sum_i Z_{i2} \sum_j X_j. \quad (A7)$$

As a result,

$$\hat{\beta}_{2SLS}(B) = \frac{(\bar{G}_2 - \bar{G}_1) \sum_i Z_{i2} Y_i - n^{-1} (\bar{G}_2 - \bar{G}_1) \sum_i Z_{i2} \sum_j Y_j}{(\bar{G}_2 - \bar{G}_1) \sum_i Z_{i2} X_i - n^{-1} (\bar{G}_2 - \bar{G}_1) \sum_i Z_{i2} \sum_j X_j} \quad (A8)$$

$$= \frac{\sum_i Z_{i2} Y_i - n^{-1} \sum_i Z_{i2} \sum_j Y_j}{\sum_i Z_{i2} X_i - n^{-1} \sum_i Z_{i2} \sum_j X_j} \quad (A9)$$

$$= \hat{\beta}_{2SLS}(Z_2). \quad (A10)$$

Hence, estimation in the first stage using B_i is identical to using Z_{i2} as an instrument, regardless of the values of Δ_G (assuming $\Delta_G \neq 0$). Also, note that if Z_{i2} is a valid instrument, then $\hat{\beta}_{2SLS}(B)$ and $\hat{\beta}_{2SLS}(Z_2)$ are consistent estimators of β . \square

Proposition 2.1

Proof. Start with the Bartik estimator,

$$\hat{\beta}_{\text{bartik}} = \frac{(\mathbf{M}_W \mathbf{B}_n)' \mathbf{Y}_n}{(\mathbf{M}_W \mathbf{B}_n)' \mathbf{X}_n} \quad (\text{A11})$$

$$= \frac{\mathbf{B}_n' \mathbf{M}_W \mathbf{Y}_n}{\mathbf{B}_n' \mathbf{M}_W \mathbf{X}_n} \quad (\text{A12})$$

$$= \frac{\tilde{\tilde{\mathbf{G}}} \tilde{\tilde{\mathbf{Z}}}_n' \mathbf{M}_W \mathbf{Y}_n}{\tilde{\tilde{\mathbf{G}}} \tilde{\tilde{\mathbf{Z}}}_n' \mathbf{M}_W \mathbf{X}_n} \quad (\text{A13})$$

$$= \frac{\mathbf{X}_n' \mathbf{M}_W \tilde{\tilde{\mathbf{Z}}}_n \tilde{\tilde{\mathbf{G}}} \tilde{\tilde{\mathbf{G}}} \tilde{\tilde{\mathbf{Z}}}_n' \mathbf{M}_W \mathbf{Y}_n}{\mathbf{X}_n' \mathbf{M}_W \tilde{\tilde{\mathbf{Z}}}_n \tilde{\tilde{\mathbf{G}}} \tilde{\tilde{\mathbf{G}}} \tilde{\tilde{\mathbf{Z}}}_n' \mathbf{M}_W \mathbf{X}_n}, \quad (\text{A14})$$

where the second equality is algebra, the third equality follows from the definition of \mathbf{B}_n , and the fourth equality follows because $\mathbf{X}_n' \mathbf{M}_W \tilde{\tilde{\mathbf{Z}}}_n \tilde{\tilde{\mathbf{G}}}'$ is a scalar. By inspection, if $\Omega^{-1} = \tilde{\tilde{\mathbf{G}}} \tilde{\tilde{\mathbf{G}}}'$, then $\hat{\beta}_{GMM} = \hat{\beta}_{\text{bartik}}$. \square

Proposition 2.2

Proof. Note that if \mathbf{B}_n is defined using a leave-one-out estimator for each location, this equivalence does not hold exactly. Instead, \mathbf{B}_n can be written as

$$\mathbf{B}_n = \frac{N}{N-1} \tilde{\tilde{\mathbf{Z}}}_n \tilde{\tilde{\mathbf{G}}} - \frac{1}{N-1} \mathbf{X}_n. \quad (\text{A15})$$

To see this, note that in this case,

$$B_{it}^{lo} = \sum_k \sum_i Z_{ikt} (N-1)^{-1} \sum_{j \neq i} G_{jkt} \quad (\text{A16})$$

$$= \sum_k \sum_i Z_{ikt} (N-1)^{-1} (N \bar{G}_{kt} - G_{ikt}) \quad (\text{A17})$$

$$= (N-1)^{-1} \left[N \sum_k \sum_i Z_{ikt} \bar{G}_{kt} - \sum_k \sum_i Z_{ikt} G_{ikt} \right] \quad (\text{A18})$$

$$= (N-1)^{-1} N \sum_i \sum_k Z_{ikt} \bar{G}_{kt} - (N-1)^{-1} \sum_k \sum_i Z_{ikt} G_{ikt} \quad (\text{A19})$$

$$= (N-1)^{-1} N \sum_i \sum_k Z_{ikt} \bar{G}_{kt} - (N-1)^{-1} \sum_i X_i. \quad (\text{A20})$$

Hence, this implies that

$$\hat{\beta}_{\text{bartik}} = \frac{\left[\tilde{\mathbf{G}}\tilde{\mathbf{Z}}'_n - N^{-1}\mathbf{X}_n' \right] \mathbf{M}_W \mathbf{Y}_n}{\left[\tilde{\mathbf{G}}\tilde{\mathbf{Z}}'_n - N^{-1}\mathbf{X}_n' \right] \mathbf{M}_W \mathbf{X}_n}. \quad (\text{A21})$$

Note that for sufficiently large N , $\hat{\beta}_{\text{bartik}} \approx \hat{\beta}_{\text{GMM}}$. This also highlights where the finite sample bias from failing to use the leave-one-out mean in \bar{G} comes from. Note that

$$\hat{\beta}_{\text{bartik}} - \beta = \frac{\tilde{\mathbf{G}}\tilde{\mathbf{Z}}'_n \mathbf{M}_W \epsilon_n}{\tilde{\mathbf{G}}\tilde{\mathbf{Z}}'_n \mathbf{M}_W \mathbf{X}_n} \quad (\text{A22})$$

and

$$\tilde{\mathbf{G}}\tilde{\mathbf{Z}}'_n \mathbf{M}_W \epsilon_n = \sum_i \sum_j Z_i \bar{G}' M_{ij} \epsilon_j. \quad (\text{A23})$$

For $i \neq j$, we have assumed independence, so bias would arise when $i = j$:

$$\tilde{\mathbf{G}}\tilde{\mathbf{Z}}'_n \mathbf{M}_W \epsilon_n = \sum_i Z_i \bar{G}' M_{ii} \epsilon_i \quad (\text{A24})$$

$$= \sum_i Z_i N^{-1} \sum_j G_j' M_{ii} \epsilon_i \quad (\text{A25})$$

$$= N^{-1} \sum_i Z_i G_i' M_{ii} \epsilon_i. \quad (\text{A26})$$

□

B Simulation details

B.1 Overview

We first set-up a simulation where Bartik works. Here is the notation:

- K : number of industries
- L : number of locations
- $g_k \sim N(0, \sigma_k^2)$: industry growth rates
- $g_{kl} \sim N(0, \sigma_{kl}^2)$: industry-location growth rates
- $g_l \sim N(0, \sigma_l^2)$: location growth rate
- $g_{lk}^{tot} = g_k + g_{kl} + g_l$: total growth observed in the industry-location
- G_l : the vector version of g_{lk}^{tot}
- To construct Z_l : draw independent standard normal variables, take the absolute values, and then normalize such that the Z 's sum to one in each location
- $X_l = Z_l' G_l$: growth rate of “employment” in the location
- $y_l = 0.5 + \theta X_l + \epsilon_l$: per capita earnings growth in the location, where $\epsilon_l \sim N(0, \sigma_\epsilon^2)$ and ϵ_l is correlated with X_l through $\text{Corr}(g_l, \epsilon_l) = \rho$

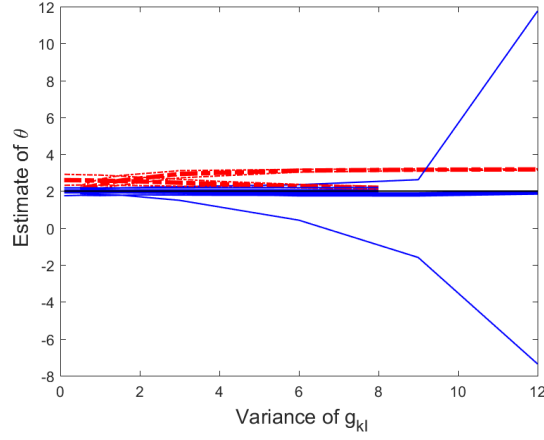
Table A1: Base simulation

Parameter	Value
K	10
L	300
σ_k^2	7
σ_{kl}^2	1
σ_l^2	1.5
θ	2
σ_ϵ^2	2
σ_z^2	1
ρ	0.6

We show simulations where we do the following:

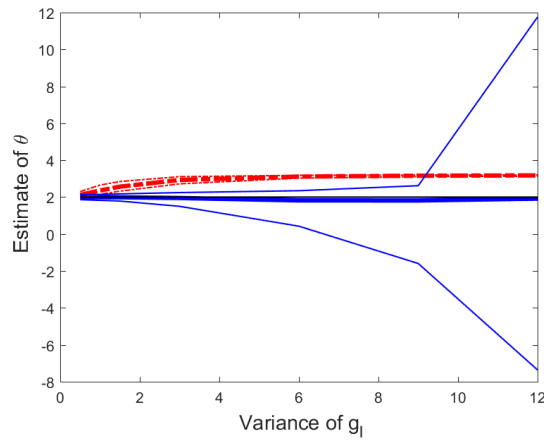
- Vary K , L , σ_k^2 , σ_{kl}^2 , σ_l^2 and σ_ϵ^2 .
- Define: $y_l = 0.5 + \theta X_l + \epsilon_l + Z_l' G_k$, where G_k is the vectorized version of g_k

Figure A1: Bartik simulation: changing variance of industry-location shocks



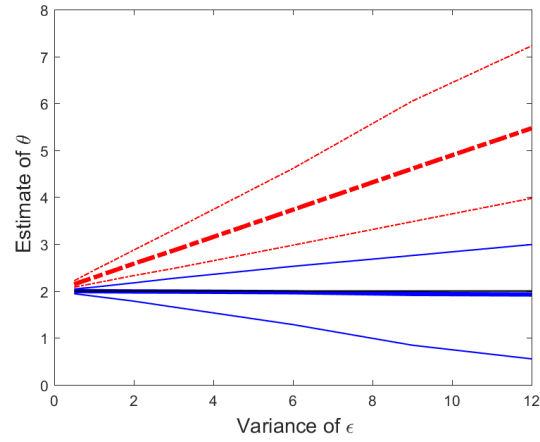
Notes: The blue solid line shows the TSLS estimates using Bartik. The red dashed line shows OLS. The thin black line shows the truth ($\theta = 2$). There are 1000 simulations at each point. The parameter values are as in Table A1, except that g_k varies. The benchmark value of σ_{kl}^2 is 1.

Figure A2: Bartik simulation: changing variance of location shocks



Notes: The blue solid line shows the TSLS estimates using Bartik. The red dashed line shows OLS. The thin black line shows the truth ($\theta = 2$). The thin lines show the 2.5 to 97.5th percentiles across the 1000 simulations at each point. The parameter values are as in Table A1, except that σ_l^2 varies. The benchmark value of σ_l^2 is 1.5.

Figure A3: Bartik simulation: changing variance of error term



Notes: The blue solid line shows the TSLS estimates using Bartik. The red dashed line shows OLS. The thin black line shows the truth ($\theta = 2$). The thin lines show the 2.5 to 97.5th percentiles across the 1000 simulations at each point. The parameter values are as in Table A1, except that σ_ϵ^2 varies. The benchmark value of σ_ϵ^2 is 2.

C Additional tables

Table A2: Shares State 2-Digit

	Industry	Location	Residual	Share of zeros
1990	.4562	.0543	.4896	.0509
2000	.4756	.0662	.4582	.0269
2010	.4006	.0323	.5672	.0304

Table A3: (State) 1980 2-Digit Industry Share Bartik Instrument

	Levels			Changes		
	1990	2000	2010	1990	2000	2010
Male	-0.35*	-0.46**	-0.08	0.18	0.04	0.06
White	-0.14	0.06	-0.07	-0.03	0.08	0.24
Native Born	-0.46***	-0.34***	-0.25	-0.12	-0.12	-0.10
12th Grade Only	-0.02	0.23	0.73	-0.29*	-1.05***	0.36
Some College	0.18	0.43*	1.43**	-0.82***	-1.50***	-1.78**
Veteran	0.66*	0.40	-0.12	0.14	-0.12	0.22
# of Children	0.13	0.07	0.27	-0.07	0.11	-0.27
Total Income	0.43	0.20	-0.96*	0.04	0.20	0.39
Social Security Income	1.47*	0.35	0.52	0.12	-0.08	-0.23
5-Year Same State	-1.14	-0.14	-1.70***	-0.78	1.77	
R^2	0.74	0.63	0.47	0.63	0.67	0.65
F	22.83	9.15	8.51	29.30	12.55	20.11
p	0.00	0.00	0.00	0.00	0.00	0.00

* $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$

Table A4: (State) 1980 2-Digit Industry Share Principal Component 1

	Levels			Changes		
	1990	2000	2010	1990	2000	2010
Male	-0.27	-0.52	-0.00	-0.15	0.18	0.11
White	0.46***	0.51*	0.35*	-0.12	-0.06	-0.38*
Native Born	0.17	0.14	0.44**	0.59**	1.11***	0.65***
12th Grade Only	0.31*	0.14	-0.17	0.02	0.54	0.18
Some College	0.09	-0.28	-0.61	1.05*	0.53	1.87**
Veteran	-0.10	-0.07	-0.86	-0.56***	-0.63**	-0.72***
# of Children	-0.01	-0.03	-0.23	0.28	0.45*	0.48*
Total Income	-2.06	0.77	0.83*	0.31	1.43	0.39
Social Security Income	-2.19*	-0.29	0.04	-0.21	0.13	0.11
5-Year Same State	4.00***	1.91***	2.15***	-2.44**	-2.17	
R^2	0.85	0.75	0.79	0.77	0.58	0.63
F	59.34	16.38	19.15	23.51	7.52	8.63
p	0.00	0.00	0.00	0.00	0.00	0.00

* $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$

Table A5: Across-Time State Regressions with 2-Digit Industry bartik_1980

	OLS			First Stages			Second Stages		
	Δ Wage	Δ Emp	Δ Emp	Δ Emp	Δ Emp	Δ Emp	Δ Wage	Δ Wage	Δ Wage
Δ Emp	0.46***						0.95***	0.76***	2.52
Bartik (1980)		0.70***	0.70**	0.19	0.36				1.18**
Male			-0.75		0.03			0.20	-0.25
White			2.56*		1.34			-0.87	-1.03
Native Born			-0.08		-0.58			0.32	0.46
12th Grade Only			0.65		0.29			-0.46*	-0.30
Some College			1.47*		1.10			-1.21***	-1.21*
Veteran			0.20		-0.07			0.36	0.33
# of Children			1.58***		0.30			-0.58*	-0.57
Δ Male				0.30	0.37			-0.82	-0.52*
Δ White				0.24*	0.23			-0.61	-0.24
Δ Native Born				0.10	-0.03			-0.28	-0.08
Δ 12th Grade Only				-0.43*	-0.32			1.04	0.43*
Δ Some College				-1.10**	-0.80			2.71	0.88
Δ Veteran				-0.26	-0.25			0.82	0.46*
Δ # of Children				-0.42***	-0.41**			0.73	0.15
Year FE	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes
State FE	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes
Observations	153	153	153	153	153	153	153	153	153
R-squared	0.93	0.67	0.80	0.85	0.87	0.84	0.92	0.20	0.87
F
p-value	0.00	0.00	0.00	0.00

* $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$

Table A6: Regression of 1-Lagged Wage Growth Residuals on Bartik, State -Digit

	OLS	Residualized			
	Lag Wage Growth	No Controls	Levels	Changes	Levels+Changes
bartik_1980	0.459*** (0.129)	0.140** (0.0486)	0.0306 (0.0331)	0.0225 (0.0330)	0.0133 (0.0315)
<i>N</i>	102	102	102	102	102

Standard errors in parentheses

* $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$