# Does Value Added Overstate Productivity Dispersion? Identification and Estimation of the Gross Output Production Function

Amit Gandhi, Salvador Navarro, David Rivers

University of Wisconsin-Madison and University of Western Ontario

May 7, 2011

## 1    Introduction

The measurement of productivity at the plant level is critical to addressing a wide of range of economic policy issues. A plant's productivity is typically defined as the residual from an underlying production function, and thus the measurement of productivity is closely tied to the estimation of the production function itself. There is a large literature that has generated several stylized facts about heterogeneity of productivity at the plant level. Among these is the general understanding that even narrowly defined industries exhibit "massive" unexplained productivity dispersion (Dhrymes (1991), Bartelsman and Doms (2000), Syverson (2004a,b), Collard-Wexler (2010), Fox and Smeets (2011)), and that productivity is closely related to other dimensions of plant level heterogeneity, such as importing (Kasahara and Rodrigue (2008)), exporting (Bernard and Jensen (1995), Bernard and Jensen (1999), Bernard et al. (2003)), wages (Baily et al. (1992)), etc.

In this paper, we show that a fundamentally different understanding of productivity

differences among plants can emerge depending on whether productivity is measured using a *gross output* or *value added* production function. A value added production function subtracts the value of a plant's "flexible" inputs (materials, energy, etc) from the plant's value of gross output to form a value added measure of output. The key obstacle to using gross output production in applied work is the lack of a clear source of identification. As first pointed out by Marschak and Andrews (1944), using the inputs and outputs of profit maximizing firms to estimate production functions gives rise to an endogeneity problem. The endogeneity problem is caused by the transmission of a plant's productivity to the firm's optimal choice of inputs. This "transmission bias" is most severe for the flexible inputs (materials, energy, etc): in a standard competitive environment, there do not exist *any* exclusion restrictions that can identify the elasticities of the flexible inputs. Since these are the same inputs that are subtracted from gross output to form value added, the value added production function can be identified via instruments when the remaining inputs, i.e., capital and labor, are "inflexible" in some way, either because of timing restrictions or adjustment costs (Bond and Söderbom (2005), Ackerberg et al. (2006)). However the value added solution does not come for free. If the "flexible inputs" do not enter the production function in Leontief fashion, then value added construction can cause its own source of bias for measuring productivity (Bruno (1978), Basu and Fernald (1995, 1997)).

A key contribution of the paper is that we present an identification strategy for gross output production functions that is grounded in standard economic assumptions. To solve the endogeneity problem caused by flexible inputs, we show that the firm's first order condition enables the revenue share of a flexible input to non-parametrically identify the input's elasticity of production. We further show the non-parametric estimator of the flexible input elasticities can be combined with standard restrictions on "inflexible" inputs to identify and estimate the full production function.

We apply our estimator to plant level data from Chile and Colombia and find that existing "stylized facts" that have emphasized the large levels of unobserved productivity differences

among observably different groups of firms (exporters versus non-exporters, importers versus non-importers, big versus small firms, etc) become orders of magnitude smaller and sometimes economically insignificant when we analyze the data through the lens of gross output. For example, the standard 90/10 productivity ratio taken among all manufacturing firms in Chile is roughly 9 (meaning that the 90th percentile firms is 9 times more productive as the 10th percentile firm), whereas under our gross output estimates the ratio falls to 2. The 95/5 ratio is even more stark: value added implies a ratio of 20 whereas gross output only entails a ratio of 3. Furthermore, exporters appear 20 percent more productive than non-exporters under value added whereas this productivity difference is only 1 percent with gross output. Similar findings hold for importing, and wages: value added overestimates in an economically significant way the productivity premium of firms who import, bigger firms, and higher wage firms as compared to gross output. Our results suggest that the bias introduced from using value added is at least as important, if not more so, than the transmission bias itself.

The rest of the paper is organized as follows. In Section 2 we describe our data, and we provide preliminary evidence of the large differences in productivity heterogeneity suggested by value-added relative to gross output. In Section 3 we discuss our identification strategy that allows us to estimate parameters on flexible inputs and work with gross output specifications, while controlling for the transmission bias. Section 4 contains a description of our estimation routine. In Section 6 we provide estimates using our approach, and compare them to both results based on value added and results that do not correct for transmission bias. Section 7 concludes.

# 2   Data and Evidence

To motivate the general interest in identification and estimation of the gross output production function, we first present descriptive evidence that suggests gross output and value added have very different empirical implications in terms of the overall picture of productivity

differences among plants. In order to do so, we make use of two commonly used plant-level manufacturing datasets. The first comes from a Colombian manufacturing census covering all manufacturing plants with more than 10 employees, from 1981-1991. This dataset has been used in several studies, including Roberts and Tybout (1997), Clerides et al. (1998), and Das, Roberts, and Tybout (2007). The second dataset comes from a census of Chilean manufacturing plants conducted by Chile's Instituto Nacional de Estadistica (INE). It covers all firms from 1979-1996 with more than 10 employees. This dataset has also been used extensively in previous studies, both in the production function estimation literature (Levinsohn and Petrin (2003)), and in the international literature (Pavcnik (2002) and Alvarez and López (2005)).

To examine the effects in the raw data, we will presently abstract from the endogeneity problem and recover productivity using simple OLS.[1] Specifically, we estimate a flexible second-order parametric approximation, i.e. translog, to both the value-added

$$
\begin{aligned}
va_{j,t} \;=\; & \beta_l l_{j,t} + \beta_k k_{j,t} + \\
& + \beta_{ll} l_{j,t}^2 + \beta_{kk} k_{j,t}^2 \\
& + \beta_{lk} l_{j,t} k_{j,t} + \nu_{jt},
\end{aligned}
\tag{1}
$$

and gross output

$$
\begin{aligned}
go_{j,t} \;=\; & \alpha_l l_{j,t} + \alpha_k k_{j,t} + \alpha_m m_{j,t} \\
& + \alpha_{ll} l_{j,t}^2 + \alpha_{kk} k_{j,t}^2 + \alpha_{mm} m_{j,t}^2 \\
& + \alpha_{lk} l_{j,t} k_{j,t} + \alpha_{lm} l_{j,t} m_{j,t} \\
& + \alpha_{km} k_{j,t} m_{j,t} + \omega_{jt},
\end{aligned}
\tag{2}
$$

---

[1]The remainder of the paper will then be largely concerned with whether these raw effects indeed remain robust after we properly correct for the endogeneity problem.

production functions. In both cases, productivity is estimated as the residual from these regressions: $\nu_{jt}$ and $\omega_{jt}$.

We estimate separate production functions for five of the largest 3-digit manufacturing industries in both Chile and Colombia, which are Food Products (311), Textiles (321), Apparel (322), Wood Products (331), and Fabricated Metal Products (381). We also estimate an aggregate specification grouping all manufacturing together.

In Table 1 we report estimates of the average input elasticities for both the value-added and gross output models. In the last row of the table we report the sum of the input elasticities. Since our second-order approximation does not impose homotheticity of the production function, this is not strictly-speaking an estimate of returns to scale, but it has a similar interpretation. In every case the value-added model substantially overestimates the sum of elasticities relative to gross output, with an average difference of 15% in Chile and 7% in Colombia.

Value added also recovers dramatically different patterns of productivity as compared to gross output. In Tables 2A and 2B, for Colombia and Chile respectively, we report estimates of several frequently analyzed statistics of the resulting productivity distributions. In the first three rows of each table we report ratios of percentiles of the productivity distribution, which measure the overall level of productivity dispersion in each industry. For Colombia, the interquartile range is between 1.65 and 2.17, and for Chile it is between 2.47 and 3.07 using the value added measure of productivity. Using the gross output measure it falls to between 1.16 and 1.23 for Colombia and between 1.30 and 1.46 for Chile. The 90/10 ratio ranges from 2.8 to 5.2 for Colombia and from 6.6 to 10.2 for Chile under value added. However, under gross output, the ratios shrink substantially to 1.4 to 1.6 and 1.7 to 2.1. The change is even more dramatic for the 95/5 split. The value added estimates indicate a range from 4.2 to 10.9 in Colombia and from 12.2 to 25.3 in Chile, whereas the gross output estimates indicate a range from 1.7 to 2.0 and from 2.1 to 2.8. There are two important implications of these results. First, value added hsuggests a much larger amount of heterogeneity in productivity

across plants *within* an industry, as the various percentile ratios are much smaller under gross output. Second, value added also implies much more heterogeneity *across* industries, which is captured by the finding that the range of the percentile ratios across industries are much tighter using the gross output measure of productivity.

In addition to having much larger overall amounts of productivity dispersion, results based on value added also suggest a much different relationship between productivity and other dimensions of plant-level heterogeneity. In the last five rows of Tables 2A and 2B we report percentage differences in productivity between plants that export some percentage of their output, import intermediate inputs, have positive advertising expenditures, have above the median (industry) level of capital stocks, and pay above the median (industry) level of wages. Using the value added estimates, exporting is generally positively correlated with productivity. For most industries exporters are found to be more productive than non-exporters, with exporters appearing to be 15% more productive in Chile and 45% more productive in Colombia across all industries. Once we account for intermediate input differences using the gross output specification, these estimates of productivity differences fall, and actually turn negative (although not statistically different from zero) in about half of the cases. Looking at all industries together, in Chile the productivity difference falls from a statistically significant 15% to a statistically insignificant -1%, and in Colombia the difference falls from a statistically significant 45% to a borderline significant 1%. A similar pattern exists when looking at importers of intermediate inputs. In all but one case, importers appear more productive than non-importers under value added. In Chile the average productivity difference is 35% across the five individual industries with a difference of 41% for all industries. For Colombia the corresponding differences are 6% and 14%. However under gross output, these numbers fall to 4% and 9% in Chile and 0% and 4% in Colombia.

A similar pattern emerges for differences in productivity based on advertising expenditures. Moving from value added to gross output, the estimated difference in productivity drops in each case for Chile, and in two cases becomes statistically insignificant. In Colom-

6

bia, a positive estimated correlation between advertising and productivity becomes negative in all but one industry.

The most striking contrast arises when we compare productivity between plants that pay wages above versus below the median. Using the productivity estimates from a value-added specification, firms that pay wages above the median wage are found to be substantially more productive, with the estimated differences ranging from 45%-111% in Chile and from 33%-56% in Colombia. In every case the estimates are statistically significant. Using the gross output specification, these estimates fall to 2%-23% in Chile and 6%-13% in Colombia, representing an average fall of approximately 75% in both countries.

While the evidence we have presented above is illustrative of important economic differences in productivity measurement that arise between value added and gross output, these results are only suggestive since we have not corrected for the standard endogeneity problem arising from the correlation of inputs and productivity. In what follows, we first show how one can identify the gross output production function accounting for the endogeneity of inputs by exploiting the firm's first order condition for the flexible inputs. Our solution offers a new way to identify the production when some inputs are "flexible" and others are inflexible, which is the case if one considers intermediate inputs as being flexibly chosen whereas capital and labor are subject to possible adjustment costs and timing restrictions. If we subtract out the flexible inputs and all remaining inputs in the production function have some amount of inflexibility built into their adjustment, then the well known proxy variable methods of OP/LP/ACF can be used to identify a value added production function. After we present the methodology, we will return to the data and compare the structural estimates of productivity generated by our gross output empirical strategy with the value added estimates generated under OP/LP/ACF. As will be shown, the main qualitative findings of the present section remain robust: the value-added bias is considerable and arguably more economically significant than the "transmission bias" itself.

# 3 Nonparametric Identification of Flexible Input Elasticities

Let $Y$ denote a firm's output and the vector $(L, K, M)$ denote a firm's inputs, with $L$ denoting labor, $K$ denoting capital, and $M$ denoting all intermediate inputs. In addition each firm has a productivity $\omega_{jt} \in \mathbb{R}$. We observe a cross section of firms $j = 1, \ldots J$ over a panel of periods $t = 1, \ldots, T$.[2] In period $t$, we will assume that firm $j$ takes its productivity level $\omega_{jt}$, capital stock $K_{jt}$, and labor $L_{jt}$ as state variables that are fixed for period $t$. That is, capital and labor are taken to be "inflexible" inputs, which is the standard assumption used to identify value added production functions (see ACF and B&S). We will describe the evolution of the state variables in the next section. The focus of this section is the intermediate inputs $M_{jt}$ (raw materials, energy, etc), which is a "flexible" input firm $j$ can control in period $t$.

We assume that productivity differences among firms are driven by heterogeneity of a Hicks-neutral form. The relationship between a firm $j$'s input and output in period $t$ is expressed as

$$
\begin{aligned}
Q_{jt} &= F(L_{jt}, K_{jt}, M_{jt})e^{\omega_{jt}} \\
Y_{jt} &= Q_{jt}e^{\varepsilon_{jt}},
\end{aligned}
\tag{3}
$$

where $F$ is the production function, $Q_{jt}$ is the output anticipated by the firm for a given vector of inputs $(L_{jt}, K_{jt}, M_{jt})$, and $Y_{jt}$ is the measured output that is actually observed by the econometrician. The difference between the firm's anticipated output and the measured output is caused by $\varepsilon_{jt}$, which can be interpreted as either an unanticipated productivity shock (in contrast to the anticipated Hicks neutral shock $\omega_{jt}$) or measurement error. We will refer to $\varepsilon_{jt}$ simply as measurement error.

---

[2]For notational simplicity we assume a balanced panel, but unbalanced panels caused by attrition can be addressed using standard selection corrections.

We focus on the classical case of perfect competition.[3] Let $\rho_t$ equal the intermediate input price and $P_t$ equal the output price, which are competitively set. The firm's first order condition with respect to $M$ yields,

$$P_t F_M(L_{jt}, K_{jt}, M_{jt})e^{\omega_{jt}} = \rho_t. \tag{4}$$

Multiplying the LHS of (4) by $\frac{F(L_{jt}, K_{jt}, M_{jt})}{F(L_{jt}, K_{jt}, M_{jt})}$, using the definition of $Q_{jt}$ in (3), and multiplying both sides of (4) by $\frac{M_{jt}}{P_t Q_{jt}}$ gives

$$\frac{F_M(L_{jt}, K_{jt}, M_{jt})M_{jt}}{F(L_{jt}, K_{jt}, M_{jt})} = \frac{\rho_t M_{jt}}{P_t Q_{jt}}.$$

Observe that the first order condition has been transformed so that $\omega_{jt}$ no longer appears in it. The firm's productivity $\omega_{jt}$ has been subsumed in the profit maximizing (anticipated) output, $Q_{jt}$. Defining the firm's anticipated revenue share of the intermediate input to be $\tilde{S}_{jt} = \frac{\rho_t M_{jt}}{P_t Q_{jt}}$, we have that

$$\tilde{S}_{jt} = \frac{F_M(L_{jt}, K_{jt}, M_{jt})M_{jt}}{F(L_{jt}, K_{jt}, M_{jt})},$$

What arises is the well known fact that the anticipated revenue share, $\tilde{S}_{jt}$, nonparametrically identifies the firm's elasticity of output with respect to the intermediate input. To see this, let $\xi_{jt}$ denote the elasticity and observe that

$$
\begin{aligned}
\xi_{jt} &= \frac{\partial Q_{jt}}{\partial M_{jt}} * \frac{M_{jt}}{Q_{jt}} \\
&= \frac{\partial F(K_{jt}, L_{jt}, M_{jt})e^{\omega_{jt}}}{\partial M_{jt}} * \frac{M_{jt}}{Q_{jt}} \\
&= F_M(K_{jt}, L_{jt}, M_{jt})e^{\omega_{jt}} * \frac{M_{jt}}{F(K_{jt}, L_{jt}, M_{jt})e^{\omega_{jt}}} \\
&= F_M(K_{jt}, L_{jt}, M_{jt}) * \frac{M_{jt}}{F(K_{jt}, L_{jt}, M_{jt})}.
\end{aligned}
$$

---

[3]As we show in Section 5 our framework can be extended to the case of imperfect competition.

Hence we have that

$$\tilde{S}_{jt} = \xi_{jt} = G(K_{jt}, L_{jt}, M_{jt}),$$

where $G$ is a non-parametric function of the inputs, that represents a known transformation of the underlying production function, $F(\bullet)$.

We can now show that, even though the anticipated revenue share is not observed by the econometrician, the data we do observe non-parametrically identifies both the elasticity $\xi_{jt}$ and measurement error $\varepsilon_{jt}$. Whereas the anticipated share $\tilde{S}_{jt}$ is not observed in the data, the realized share $S_{jt} = \tilde{S}_{jt}(e^{\varepsilon_{jt}})^{-1}$ is observed. Hence letting $s_{jt} = \ln S_{jt}$, we have that

$$s_{jt} = \ln G(K_{jt}, L_{jt}, M_{jt}) - \varepsilon_{jt} \tag{5}$$

We refer to equation (5) as the share equation. Since measurement error $\varepsilon_{jt}$ is by construction a stochastic error term that is independent of the inputs $(L_{jt}, K_{jt}, M_{jt})$, the non-parametric regression of $s_{jt}$ on $(L_{jt}, K_{jt}, M_{jt})$ identifies both the log anticipated elasticity $\ln \xi_{jt} = \ln G(L_{jt}, K_{jt}, M_{jt})$ and the measurement error $\varepsilon_{jt} = s_{jt} + \ln G(L_{jt}, K_{jt}, M_{jt})$.

The revenue share of intermediates is equal to the elasticity of output with respect to intermediates times measurement error. The share equation (5) generalizes the index number approach, in which productivity is calculated as the difference between log output and a weighted sum of inputs, with the weights being the cost shares of the inputs.[4] It is more general in several dimensions. First, this approach is robust to measurement error.. Second, as will become apparent later, it does not rely on the assumption of constant returns to scale.[5] Third, it is robust to alternative assumptions regarding the nature of output market competition.[6] Fourth, we do not need to assume that the other inputs are competitively or flexibly chosen.

---

[4]Under the assumption of perfect competition in the output market, revenue shares equal costs shares.

[5]The assumption of constant returns to scale is common in practice with the index number approach. The reason is that measuring the cost share of capital requires a measure of the rental rate of capital, which is often not observed and can be difficult to estimate. Assuming constant returns to scale implies that the cost share is equal to one minus the sum of the other input shares.

[6]See Section 5

# 4    Estimation

To see how the non-parametric identification of the flexible input elasticities enables us to identify a gross output production function, it is useful to recall the source of the identification problem, which is covered in greater depth by Bond and Söderbom (2005) and Ackerberg et al. (2006). Observe that under perfect competition in the output and intermediate input market, we have that intermediate inputs are optimally chosen such that $M_{jt} = \mathbb{M}_t(L_{jt}, K_{jt}, \omega_{jt})$ for some time varying function $\mathbb{M}_t$. Thus $M_{jt}$ is "colinear" with the other inputs (including productivity) that appear in the production function. To see the consequence of this "colinearity" problem, take a a flexible second-order parametric approximation to $F$ (i.e. translog),[7]

$$
\begin{aligned}
y_{j,t} = {} & \alpha_l l_{j,t} + \alpha_k k_{j,t} + \alpha_m m_{j,t} \\
& + \alpha_{ll} l_{j,t}^2 + \alpha_{kk} k_{j,t}^2 + \alpha_{mm} m_{j,t}^2 \\
& + \alpha_{lk} l_{j,t} k_{j,t} + \alpha_{lm} l_{j,t} m_{j,t} \\
& + \alpha_{km} k_{j,t} m_{j,t} + \omega_{jt} + \varepsilon_{jt}.
\end{aligned}
\tag{6}
$$

Since $m_{jt} = m_t(k_{jt}, l_{jt}, \omega_{jt})$, it is clear that $m_{jt}$ is an endogenous regressor since it is determined by the "anticipated" part of the residual $\omega_{jt}$. However, since there is no source of cross sectional variation in $m_{jt}$ other than the firm's remaining productive inputs $(l_{jt}, k_{jt}, \omega_{jt})$, there does not exist any exclusion restriction to vary the intermediate input from outside of the production function. This can be seen as an "impossibility" result concerning the existence of instruments for $m_{jt}$, which poses a fundamental identification problem for the coefficients associated with $m_{jt}$, i.e., $\theta_1 = (\alpha_m, \alpha_{mm}, \alpha_{lm}, \alpha_{km})$.[8]

To see how we solve the identification problem associated with these coefficients, observe that from our result in Section 3 we have a consistent non-parametric estimate of the output

---

[7]There is nothing special about the translog approximation for our purposes. Any other approximation (CES, higher order polynomials, etc.) would work just as well.

[8]See Ackerberg et al. (2006).

elasticity with respect to $m_{jt}$ for each observation in the data, i.e., $\hat{\xi}_{jt}$. For the second order (in logs) approximation to the production function, the implied elasticity with respect to the intermediate input is

$$e_{j,t}(\theta_1) = \alpha_m + 2\alpha_{mm}m_{j,t} + \alpha_{lm}l_{j,t} + \alpha_{km}k_{j,t}.$$

Observe that this implied elasticity only depends on the problematic parameters $\theta_1$. Thus we can consistently estimate $\theta_1$ using a minimum distance estimator that minimizes the distance between the implied elasticities and the non-parametric estimated elasticities,

$$\min_{\alpha_{\cdot m}} \Sigma \left( \hat{\xi}_{j,t} - e_{j,t}(\theta_1) \right)^2. \tag{7}$$

Solving (7) gives us consistent estimates of the parameters $\theta_1$, the parameters that could not be identified with instrumental variables and that are the original source of the identification problem commonly solved with the value-added approximation.

The remaining parameters $\theta_2 = (\alpha_l, \alpha_k, \alpha_{ll}, \alpha_{kk}, \alpha_{lk})$ of the production function are coefficients on terms only involving capital and labor. To see how these parameters are identified, we must now consider how capital and labor are set by the firm. If both capital and labor are "sticky" in the sense of having adjustment costs (which is the source of what makes them inflexible), then $l_{jt} = l(l_{jt-1}, \mathcal{I}_{jt-1})$ and $k_{jt} = k(k_{jt-1}, \mathcal{I}_{jt-1})$ where $\mathcal{I}_{jt-1}$ is the firm $j$'s information set at the start of period $t-1$ (which includes, among other possible things, productivity $\omega_{jt-1}$). This has the implication that the firm's capital and labor in period $t$ are known to the firm in period $t-1$, and are thus an element of firm $j$'s period $t-1$ information set, i.e., $k_{jt}, l_{jt} \in \mathcal{I}_{jt-1}$.[9] Assuming that productivity evolves according to a first-order Markovian process, $\omega_{jt}$ can be expressed as $\omega_{jt} = g(\omega_{jt-1}) + \eta_{jt}$. The term $\eta_{jt}$ represents an innovation to the firm's productivity that, by construction, is orthogonal to the firm's information set at period $t-1$. Thus we have the conditional moment restriction

---

[9]An alternative is to follow ACF and also treat $l_{jt}$ as a static and variable input under an additional timing restriction that it is chosen at a point between $t$ and $t-1$ so that it is not colinear with $k_{jt}$ and $\omega_{jt}$.

$E[\eta_{jt} \mid k_{jt}, l_{jt}] = 0$ that can be used to identify $\theta_2$.

To consistently estimate $\theta_2$, we can follow the following steps that summarize our entire estimation procedure:

1. In the first step, we nonparametrically recover consistent estimates of elasticities $\hat{\xi}_{jt}$ and measurement error $\hat{\varepsilon}_{jt}$ for all $(j, t)$.

2. In the second step, we use our elasticity estimates $\hat{\xi}_{jt}$ to estimate the parameters $\theta_1$ associated with the intermediate input $m_{jt}$ via the minimum distance objective function (7).

3. Finally for any value of the parameter vector $\theta_1$, we can use our estimates $\hat{\theta}_1$ and $\hat{\varepsilon}_{jt}$ from steps (1)-(2) to construct productivity

$$
\begin{aligned}
\omega_{j,t}(\theta_2) &= y_{jt} - \alpha_l l_{j,t} - \alpha_k k_{j,t} - \hat{\alpha}_m m_{j,t} \\
&\quad -\alpha_{ll} l_{j,t}^2 - \alpha_{kk} k_{j,t}^2 - \hat{\alpha}_{mm} m_{j,t}^2 \\
&\quad -\alpha_{lk} l_{j,t} k_{j,t} - \hat{\alpha}_{lm} l_{j,t} m_{j,t} \\
&\quad -\hat{\alpha}_{km} k_{j,t} m_{j,t} - \hat{\varepsilon}_{jt}.
\end{aligned}
$$

Then nonparametrically regressing $\omega_{jt}(\theta_2)$ on $\omega_{jt-1}(\theta_2)$, we can recover innovation $\eta_{jt}(\theta_2)$ as a function of the parameter to be estimated. Finally we use the orthogonality conditions $\eta_{jt} \perp k_{jt}$, $\eta_{jt} \perp l_{jt}$, $\eta_{jt} \perp k_{jt}^2$, $\eta_{jt} \perp l_{jt}^2$, and $\eta_{jt} \perp k_{jt} l_{jt}$ implied by the conditional moment restriction $E[\eta_{jt} \mid k_{jt}, l_{jt}] = 0$ to estimate $\theta_2$.[10]

---

[10]The fact that we can separate our procedure into 3 steps relies on the property that the implied elasticity depends only on $\theta_1$. Other approximations (like the CES) do not have this property and the elasticity depends not only on parameters related to $m_{jt}$, but parameters related to other inputs as well. In this case, we can either estimate the additional parameters as part of step 2 or we could estimate all of the parameters jointly, i.e., by doing steps 2 and 3 together where we stack the "moment" conditions implied by step 2 and estimate all parameters as a GMM problem.

# 5 Recovering Industry Markups under Imperfect Competition using Revenue Production Functions

In this section we relax the assumption that firms operate in a perfectly competitive environment, and show how our approach can be extended to accommodate imperfect competition. Not only can we still control for the endogeneity of inputs and work with gross output specifications of production in this setting, but we can also use our approach to recover industry-specific time-varying markups. Relaxing the assumption of perfect competition has two important implications. First, deflated revenue is no longer a valid proxy for quantity produced, as under imperfect competition firms will no longer all necessarily charge the same price. As a result, variation in firm-specific prices needs to be accounted for. Second, a firm's first order condition will depend on their ability to markup prices over marginal cost. We will deal with the latter consideration first.

Let $\Lambda_{jt}$ denote a firm's marginal cost. The first-order condition with respect to $M_{jt}$ for a cost minimizing firm will be

$$\Lambda_{jt} F_M \left(L_{jt}, K_{jt}, M_{jt}\right) e^{\omega_{jt}} = \rho_{jt},$$

where recall that $\rho_{jt}$ is the price of $M_{jt}$. Using a similar transformation as in the perfectly competitive case we obtain

$$\frac{\Lambda_{jt}}{P_{jt}} \frac{F_M(L_{jt}, K_{jt}, M_{jt})M_{jt}}{F(L_{jt}, K_{jt}, M_{jt})} = \frac{\rho_t M_{jt}}{P_{jt}Q_{jt}},$$

or

$$\tilde{S}_{jt} = \xi_{jt} \frac{P_{jt}}{\Lambda_{jt}} = \frac{G\left(K_{jt}, L_{jt}, M_{jt}\right)}{\Psi_{jt}},$$

where $\Psi_{jt} = \frac{P_{jt}}{\Lambda_{jt}}$ denotes the markup.

The two key differences between the perfectly competitive case and this case are that

a) we no longer restrict the firm's price to be constant, and b) the firm's revenue share no longer equals the input elasticity directly, but rather it equals the input elasticity divided by the inverse markup charged by the firm. As before, we can rewrite this expression above in terms of the observed share as

$$s_{jt} = -\psi_{jt} + \ln G\left(K_{jt}, L_{jt}, M_{jt}\right) - \varepsilon_{jt}, \tag{8}$$

where $\psi_{jt} = \ln \Psi_{jt}$. Notice that equation (8) nests the one obtained for the perfectly competitive case in (5), the only difference being the addition of the log markup $\psi_{jt}$ which is equal to 0 under perfect competition.

We now develop the use of the share equation (8) to estimate production functions among imperfectly competitive firms under the restriction that all firms have the same markup, i.e. $\psi_{jt} = \psi_t$. We further justify why this case is of interest below when we introduce a demand system. For the moment, we simply note that in this case (8) becomes

$$s_{jt} = -\psi_t + \ln G\left(K_{jt}, L_{jt}, M_{jt}\right) - \varepsilon_{jt}. \tag{9}$$

Notice that, once we control for a time-varying intercept, measurement $\varepsilon_{jt}$ can be recovered as before.

More importantly, an additional implication of equation (9) is that the growth pattern in the markups can be recovered without further assumptions. To see why, rewrite the intermediate input elasticity so that we can break it into two parts: a component that varies with inputs and a constant $\mu$

$$\ln \xi_{jt} = \ln G\left(K_{jt}, L_{jt}, M_{jt}\right) = \Phi\left(K_{jt}, L_{jt}, M_{jt}\right) + \mu.$$

Then, simply rewrite equation (9)

$$
\begin{aligned}
s_{jt} &= (-\psi_t + \mu) + \Phi\left(L_{jt}, K_{jt}, M_{jt}\right) - \varepsilon_{jt} \\
&= -\gamma_t + \Phi\left(L_{jt}, K_{jt}, M_{jt}\right) - \varepsilon_{jt}
\end{aligned}
\tag{10}
$$

Equation (10) immediately shows that, on top of recovering measurement error $\varepsilon_{jt}$, we can recover logmarkups up to a constant, $\gamma_t = \psi_t - \mu$, as well as the input elasticity sans the constant,

$$
\ln \xi_{jt}^{\mu} = \Phi(L_{jt}, K_{jt}, M_{jt}) = \ln G(L_{jt}, K_{jt}, M_{jt}) - \mu.
$$

If the intercept of equation (9) ($\mu$) can be identified (which, as we show below, it can), then the level of the markups can also be recovered.

The fact that time-varying markups (up to constant) can be recovered immediately from the share equation is, to the best of our knowledge, a new result. As opposed to the results in Hall (1988) and Basu and Fernald (1995, 1997) we do not need to impose restrictions on the demand for all other inputs (i.e., competitive input markets for all inputs), impose restrictions on the shape of the production function (homogeneity) or compute or estimate the rental rate of capital/profit for the entrepreneur. Even without the ability to recover the constant, the fact that we can recover the pattern of markups is an interesting finding in itself since it allows, for example, to check whether market power has increased over time, or to analyze the behavior of market power with respect to the business cycle.

In order to recover $\mu$, and hence the level of markups and elasticities, an adapted version of the first two steps of our estimation algorithm is employed. As before, we begin by taking a second-order approximation to the production function and notice that the implied intermediate input elasticity is given by:

$$
\begin{aligned}
e_{jt} &= \alpha_m + 2\alpha_{mm}m_{jt} + \alpha_{lm}l_{jt} + \alpha_{km}k_{jt} \\
&= \alpha_m\left(1 + 2\frac{\alpha_{mm}}{\alpha_m}m_{jt} + \frac{\alpha_{lm}}{\alpha_m}l_{jt} + \frac{\alpha_{km}}{\alpha_m}k_{jt}\right),
\end{aligned}
$$

which implies that

$$\ln e_{jt} = \ln\left(1 + 2\frac{\alpha_{mm}}{\alpha_m}m_{jt} + \frac{\alpha_{lm}}{\alpha_m}l_{jt} + \frac{\alpha_{km}}{\alpha_m}k_{jt}\right) + \ln(\alpha_m)$$

and hence $\ln \alpha_m$ is the unidentified constant $\mu$.[11] Hence if we define

$$e_{jt}^\mu(\theta_1^\mu) = \left(1 + 2\frac{\alpha_{mm}}{\alpha_m}m_{jt} + \frac{\alpha_{lm}}{\alpha_m}l_{jt} + \frac{\alpha_{km}}{\alpha_m}k_{jt}\right)$$

and $\theta_1^\mu = \left(\frac{\alpha_{mm}}{\alpha_m}, \frac{\alpha_{lm}}{\alpha_m}, \frac{\alpha_{km}}{\alpha_m}\right)$, the steps of the estimation procedure are very similar to those of Section 4.

1. In the first step, we nonparametrically recover consistent estimates of $\gamma_t$, $\varepsilon_{jt}$ and $\xi_{jt}^\mu$ from 10

2. In the second step, we solve

$$\min_{\frac{\alpha_{\cdot m}}{\alpha_m}} \Sigma \left(\hat{\xi}_{j,t}^\mu - e_{j,t}^\mu(\theta_1^\mu)\right)^2$$

   to recover the parameters related to intermediate inputs, $\theta_1^\mu$, up to the constant $\mu$.

While it may seem that separating the markups and the inherent production function parameters in $\theta_1$ from the constant may require an arbitrary normalization, this is not necesarily the case. To see how one can recover the constant and the remaining parameters of the production function we follow Klette and Griliches (1996) and specify a demand system consistent with our assumption of markups being constant. We let

$$\frac{P_{jt}}{\Pi_t} = \left(\frac{Q_{jt}}{Q_t}\right)^{\tau_t - 1} e^{\Xi_{jt}}, \tag{11}$$

where $\Pi_t$ is the industry price index, $Q_t$ is a quantity index that plays the role of a demand

---

[11]This result, i.e. separating the constant from the production function approximation, is not unique to our second-order approximation. As before, it holds true for other approximations.

shifter as in Klette and Griliches (1996), $\Xi_t$ is a mean zero demand shock observable to the firm at $t$ and $\psi_t = -\ln \tau_t$. [12]

The observed output is now given by the firm's real revenue

$$R_{jt} = \frac{P_{jt}}{\Pi_t} Q_{jr} e^{\varepsilon_{jt}},$$

or in logs

$$r_{jt} = (p_{jt} - \pi_t) + q_{jt} + \varepsilon_{jt}. \tag{12}$$

Replacing 11 into 12 we obtain

$$r_{jt} = \tau_t q_{jt} - (\tau_t - 1) q_t + \Xi_{jt} + \varepsilon_{jt}. \tag{13}$$

>From our definitions of $\gamma_t = -\psi_t + \mu$ and $\psi_t = -\ln \tau_t$ we can write

$$\tau_t = e^{\gamma_t} e^{-\mu},$$

where $\gamma_t$ is known from our analysis above and we seek to recover $\mu$. Replacing back into 13 we get

$$
\begin{aligned}
r_{jt} &= e^{\gamma_t} e^{-\mu} q_{jt} - \left(e^{\gamma_t} e^{-\mu}\right) q_t + \Xi_{jt} + \varepsilon_{jt}. \\
&= e^{\gamma_t} e^{-\mu} \ln F(K_{jt}, L_{jt}, M_{jt}) - \left(e^{\gamma_t} e^{-\mu}\right) q_t + [\tau_t \omega_{jt} + \Xi_{jt}] + \varepsilon_{jt}
\end{aligned}
$$

>From this equation one can already see how the constant will be recovered. It is variation in $q_t$ that will identify it.

The rest of the estimation procedure is almost the same as before. The key difference is

---

[12]We can allow for time varying firm specific markups. If we let $\Upsilon_{jt} > 0$ be an independent demand shock that is realized after inputs are chosen, then *expected markups* will be equalized across firms, i.e., $E(\Psi_{jt}) = \Psi_t$ and $\Xi_{jt}$ will enter into the firm's period $t$ input decisions. That is, while actual markups $\Psi_{jt} = \frac{P_{jt}}{\Lambda_{jt}}$ will be firm specific due to the $\Upsilon_{jt}$ demand shocks, firms will still have ex-ante symmetric markups.

that now we cannot recover $\omega_{jt}$ but rather only[13]

$$\tau_t \omega_{jt} + \Xi_{jt},$$

i.e., a linear combination of productivity and the demand shock. The reason is obvious: since we do not observe prices, we have no way of disentangling whether, after controlling for inputs, a firm has higher revenues because it is purely more productive ($\omega_{jt}$) or because it can sell at a higher price ($\Xi_{jt}$). Notice, however that this is in fact the only thing that matters for the firm's decision, i.e., in the maximization problem the firm only cares about $\tau_t \omega_{jt} + \Xi_{jt}$ and not about the separate components. Then, since we can write $\tau_t \omega_{jt} + \Xi_{jt}$ as a function of the parameters that remain to be estimated, $\theta_2^\mu$, by imposing the Markovian assumption on the sum we can use a similar moment restriction (paired with variation in $q_t$) to identify the remaining parameters.

# 6    Application

Given that most datasets do not contain good instruments that can be used to correct for the transmission bias, recent work on production function estimation has focused on the structural techniques developed by Olley and Pakes (1996) and Levinsohn and Petrin (2003).[14] However, recently Bond and Söderbom (2005) and Ackerberg et al. (2006), henceforth B&S and ACF, have uncovered a problem underlying these methods. In particular what they show is that there is a fundamental problem with the identification of the coefficients on flexible inputs using these techniques. For those inputs that are chosen perfectly flexibly there is no independent variation in these inputs that identifies their parameters in the production function. B&S solve this by adding adjustment costs to all inputs. With adjustment costs, previous input levels affect current input decisions and then are an independent source of

---

[13]Or, of course, $\omega_{jt} + \frac{\Xi_{jt}}{\tau_t}$.

[14]See Griliches and Mairesse (1998) for a summary of the various methods that have been developed for dealing with the transmission bias.

variation that leads to identification. ACF solve this problem by imposing restrictions on the timing in which inputs are chosen. If inputs are chosen before productivity for the period is fully realized, then the innovation in productivity between when those inputs are chosen and when production occurs allows for parameters on these inputs to be identified. A key implication of these two papers is that these structural methods cannot allow for inputs that are perfectly flexible.[15] If a researcher believes that some inputs are chosen flexibly, such as intermediate inputs, then these inputs need to be netted out from gross output, and a value-added specification must be used instead. Note that the use of value added in these cases is not due to a belief that the assumptions justifying value added actually hold, but rather results from an inability to estimate a gross output model.

A key contribution of our approach, which is described in Sections 3 and 4, is that it solves this problem identified by B&S and ACF. It allows the researcher to estimate a production function and recover estimates of productivity, while being able to simultaneously control for potential bias introduced by the correlation between unobserved productivity and input decisions (transmission bias) and the bias introduced by netting out flexibly chosen inputs and working instead with value-added specifications of production.

We estimate a gross output production function using our new approach for each of five of the largest industries in both Chile and Colombia: Food Products (311), Textiles (321), Apparel (322), Wood Products (331), and Fabricated Metal Products (381), as well as one for all manufacturing industries grouped together. In order to isolate the effect of controlling for the transmission bias, we first compare the estimate from our model to those from the reduced-form gross output model. A well-known result is that failing to control for the transmission bias leads to overestimates of the coefficients on more flexible inputs. The intuition behind this is that the more flexible the input is, the more it responds to productivity shocks and the higher the degree of correlation between that input and unobserved produc-

---

[15]In Bond and Söderbom (2005) input decisions are constrained by previous input choices. In Ackerberg et al. (2006), inputs cannot adjust to the entire productivity shock, since they are chosen before productivity is fully realized. In both cases there are no perfectly flexible inputs.

tivity. In Table 3 we report the average input elasticity estimates for the two models. The estimates show that in the reduced-form model substantially overestimates the elasticity of intermediate inputs in every case. Often the difference is as much as 40%, which illustrates the importance of controlling for the endogeneity generated by the correlation between input decisions and productivity.

In order to isolate the effect that subtracting out intermediate inputs and using value added specifications has on the results, we also compare our results using a gross output specification to results obtained by using the method developed by ACF for a value-added specification. Both methods control for the transmission bias, and our method controls for any potential bias generated from value added. In Table 4 we report estimates of the average input elasticities for both models. Although the differences between value added and gross output are smaller than with the reduced-form estimation, value-added still overestimates the sum of elasticities in all but one case, by an average of 6% in Chile and 3% in Colombia. In Tables 5A and 5B we summarize our estimation results related to productivity. As can be seen from the tables, after correcting for the transmission bias, the results described in Section 2 persist. Value added continues to substantially overstate the level of productivity dispersion within and across industries. It also generates misleading estimates of the relationship between productivity and other dimensions of plant-level heterogeneity. One other important conclusion that results from these estimates is that the bias induced from value-added has a much larger effect on the productivity estimates than the transmission bias. This suggests that being able to avoid the value-added bias by working with gross output specifications is more important from a policy perspective than controlling for the transmission bias.

# 7 Conclusion

# References

**Ackerberg, Daniel A., Kevin Caves, and Garth Frazer**, "Structural Identification of Production Functions," 2006. Unpublished Manuscript, UCLA Economics Department.

**Alvarez, Roberto and Ricardo A. López**, "Exporting and Performance: Evidence from Chilean Plants," *Canadian Journal of Economics*, 2005, *38* (4), 1384–1400.

**Baily, Martin Neil, Charles Hulten, and David Campbell**, "Productivity Dynamics in Manufacturing Plants," *Brookings Papers on Economic Activity. Microeconomics*, 1992, pp. 187–267.

**Bartelsman, Eric J. and Mark Doms**, "Understanding productivity: lessons from longitudinal microdata," Finance and Economics Discussion Series 2000-19, Board of Governors of the Federal Reserve System (U.S.) 2000.

**Basu, Susanto and John G Fernald**, "Are Apparent Productive Spillovers a Figment of Specification Error?," *Journal of Monetary Economics*, August 1995, *36*, 165 – 188.

_ **and** _ , "Returns to Scale in U.S. Production: Estimates and Implications," *Journal of Political Economy*, April 1997, *105* (2), 249 – 283.

**Bernard, Andrew B. and J. Bradford Jensen**, "Exporters, Jobs, and Wages in U.S. Manufacturing: 1976-1987," *Brookings Papers on Economic Activity. Microeconomics*, 1995, pp. 67–119.

_ **and** _ , "Exceptional Exporter Performance: Cause, Effect or Both?," *Journal of International Economics*, 1999, *47* (1), 1–25.

_ , **Jonathan Eaton, J. Bradford Jensen, and Samuel Kortum**, "Plants and Productivity in International Trade," *American Economic Review*, 2003, *93* (4), 1268–1290.

**Bond, Stephen and Måns Söderbom**, "Adjustment Costs and the Identification of Cobb Douglas Production Functions," 2005. Unpublished Manuscript, The Institute for Fiscal Studies, Working Paper Series No. 05/4.

**Bruno, M.**, "Duality, Intermediate Inputs, and Value-Added," in M. Fuss and McFadden D., eds., *Production Economics: A Dual Approach to Theory and Practice*, Vol. 2, Amsterdam: North-Holland, 1978.

**Clerides, Sofronis K., Saul Lach, and James R. Tybout**, "Is Learning by Exporting Important? Micro-Dynamic Evidence from Colombia, Mexico, and Morocco," *The Quarterly Journal of Economics*, 1998, *113* (3), 903–947.

**Collard-Wexler, Allan**, "Productivity Dispersion and Plant Selection in the Ready-Mix Concrete Industry," 2010. NYU Stern working paper.

**Dhrymes, Phoebus**, "The Structure of Production Technology: Productivity and Aggregation Effects," 1991. Columbia University.

**Fox, Jeremy and ValÃ©rie Smeets**, "Does Input Quality Drive Measured Differences in Firm Productivity?," *International Economic Review*, 2011. forthcoming.

**Griliches, Zvi and Jacques Mairesse**, "Production Functions: The Search for Identification," in "Econometrics and Economic Theory in the Twentieth Century: The Ragnar Frisch Centennial Symposium," New York: Cambridge University Press, 1998, pp. 169–203.

**Hall, R.E.**, "The relation between price and marginal cost in US industry," *The Journal of Political Economy*, 1988, *96* (5), 921–947.

**Kasahara, Hiroyuki and Joel Rodrigue**, "Does the Use of Imported Intermediates Increase Productivity? Plant-level Evidence," *The Journal of Development Economics*, 2008, *87* (1), 106–118.

**Klette, Tor Jacob and Zvi Griliches**, "The Inconsistency of Common Scale Estimators When Output Prices are Unobserved and Endogenous," *Journal of Applied Econometrics*, July 1996, *11* (4), 343–361.

**Levinsohn, James and Amil Petrin**, "Estimating Production Functions Using Inputs to Control for Unobservables," *Review of Economic Studies*, April 2003, *70* (243), 317–342.

**Marschak, Jacob and W.H. Andrews**, "Random Simultaneous Equations and the Theory of Production," *Econometrica*, 1944, *12*, 143–205.

**Olley, G. Steven and Ariel Pakes**, "The Dynamics of Productivity in the Telecommunications Equipment Industry," *Econometrica*, November 1996, *64* (6), 1263–1297.

**Pavcnik, Nina**, "Trade Liberalization Exit and Productivity Improvements: Evidence from Chilean Plants," *Review of Economic Studies*, 2002, *69*, 245–276.

**Roberts, M.J. and J.R. Tybout**, "The Decision to Export in Colombia: An Empirical Model of Entry with Sunk Costs," *American Economic Review*, 1997, *87* (4), 545–564.

**Syverson, Chad**, "Market Structure and Productivity: A Concrete Example," *The Journal of Political Economy*, 2004, *112* (6), 1181–1222.

‒ , "Product Substitutability and Productivity Dispersion," *The Review of Economics and Statistics*, 2004, *86* (2), 534–550.

# Appendix A: Value-Added Bias

Value added is defined as the difference between gross output and expenditures on intermediates:

$$VA_{jt} \;=\; Q_{jt} - \rho M_{jt},$$

where the price of output has been normalized to 1, and where $\rho$ is the price of intermediate inputs and $s_{jt}^m$ is the share of intermediate expenditures for plant $j$ in period $t$. Using the definition of $s_{jt}^m$, we have that

$$VA_{jt} = Q_{jt} \left( 1 - s_{jt}^m \right).$$

Consider a generic production function $F(K, L, M)$. Value added can then be expressed as follows:

$$VA_{jt} = Q_{jt} \left( 1 - s_{jt}^m \right) = F(K_{jt}, L_{jt}, M_{jt}) e^{\omega_{jt}} \left( 1 - s_{jt}^m \right).$$

Taking logs of both sides yields the following relationship:

$$va_{jt} = \ln\left(F(K_{jt}, L_{jt}, M_{jt})\right) + \omega_{jt} + \ln\left(1 - s_{jt}^m\right). \tag{14}$$

If attempt to estimate a value added model, by regressing log value added on labor and capital only, it can be seen that there are two omitted variables. The first is intermediate inputs, and the second is one minus the share of intermediates in total output. Let us first focus on intermediate inputs. The only way in which estimates of $F()$ and productivity $\omega_{jt}$ do not suffer from bias due to the omitted variable $M$, is if one of two conditions are satisfied. The first is that $M$ is not correlated with any of other other inputs. In this case, only estimates of the function $F()$ will be unbiased. Productivity estimates remain biased, as they will be a function of omitted intermediate inputs. This assumption is likely rejected in most datasets, and is certainly rejected in both the Chilean and Colombian datasets, in all industries. The second condition is that $M$ has no explanatory power for output, conditional on $K$ and $L$. One known example of this is when the production function, $F$, is Leontief in intermediates, *i.e.*, $M$ is used in fixed proportions relative to $K$ and $L$.

In order to illustrate this, consider a first-order approximation to $F$. Value added in logs is then given by,

$$va_{jt} = \alpha l_{jt} + \beta k_{jt} + \gamma m_{jt} + \omega_{jt}. \tag{15}$$

We can see from equation (15) that a regression of log value added on log labor and log capital leads to a classic omitted variable bias in the estimates of the labor and capital elasticities ($\alpha$ and $\beta$), as well as productivity, $\omega$. The estimates of the elasticities will suffer from a bias that is proportional to the correlation between log labor and capital and log intermediates:

$$\begin{pmatrix} \hat{\alpha} \\ \hat{\beta} \end{pmatrix} = \begin{pmatrix} \alpha \\ \beta \end{pmatrix} + \begin{pmatrix} \delta_\alpha \\ \delta_\beta \end{pmatrix},$$

where the bias terms, $\delta_\alpha$ and $\delta_\beta$, are equal to $\gamma$ (the output elasticity of intermediate inputs) multiplied by the coefficients of a regression of log intermediates on log labor and log capital. This implies that the resulting estimate of productivity is also biased. The estimated value of productivity is:

$$\widehat{\omega}_{jt} = \omega_{jt} + \gamma \left( m_{jt} - \widehat{m_{jt}} \right),$$

where $\widehat{m_{jt}}$ denotes the predicted value from a regression of log intermediates on log labor and log capital. The difference between $m_{jt}$ and $\widehat{m_{jt}}$ is a measure of relative intermediate input intensity. Positive values indicate that a firm uses more intermediate inputs compared to other firms with the same level of capital and labor.

It is now clear where the bias originates. Value added does not take into account variation in intermediate input intensity. As a result, differences in this intensity will be measured as differences in productivity instead. This comes from the restrictive assumption that value added places on intermediate inputs; in particular, that differences in intermediate input use does not explain any variation in output, conditional on capital and labor use. It is

important to note that in cases in which this assumption is satisfied, productivity estimates based on value added and gross output will be equivalent. If capital and labor perfectly predict intermediate inputs, $m_{jt} - \widehat{m_{jt}} = 0$, and there is no bias.

A second source of bias comes from the second unobserved term in equation (14), $\ln\left(1 - s_{jt}^m\right)$. When all inputs are perfectly flexible in the short run, the share of intermediate inputs in total output will be constant across firms. As we showed earlier in the share equation, the share of intermediate inputs, $s_{jt}^m$, is equal to the elasticity of output with respect to intermediate inputs. This elasticity is a function of the level of other inputs, $K$ and $L$ and the production function, $F\left(\right)$. When all inputs are perfectly flexible, elasticities are just a function of $F\left(\right)$, which is common across firms. Therefore share of intermediate inputs will be the same across firms. As a result, $\ln\left(1 - s_{jt}^m\right)$ will be a constant in the error term of equation (14). As a result, estimates of $F\left(\right)$ remain unbiased and all productivity estimates will be scaled by a constant:$\ln\left(1 - s^m\right)$.[16] This is the well-known result that productivity estimates based on value-added need to be scaled up by the share of value-added in gross output.[17]

This is no longer the case when some inputs are not perfectly flexible. Consider the case of hiring and firing costs for labor. If a firm has more labor than it would otherwise demand due to the presence of firing costs, then it will compensate by demanding less that the otherwise optimal level of intermediate inputs. This will drive down the share, $s_{jt}^m$. Similarly firms with less than the optimal amount of labor with have higher shares, $s_{jt}^m$. Without further assumptions it is not possible to sign the direction of this bias.

---

[16]If input and output prices move over time, then the constant scaling productivity will be time varying. Including time dummies in the regression will control for this.

[17]Note that one minus the share of intermediate inputs is the share of value-added in gross output.