# Some Unpleasant Markup Arithmetic: Production Function Elasticities and their Estimation from Production Data$^\star$

Steve Bond[a], Arshia Hashemi[b], Greg Kaplan[b,c,*], Piotr Zoch[d,e]

[a]*University of Oxford, Oxford, UK*
[b]*The University of Chicago, Chicago, USA*
[c]*National Bureau of Economic Research, Cambridge, USA*
[d]*University of Warsaw, Warsaw, Poland*
[e]*FAME/GRAPE, Warsaw, Poland*

---

**Abstract**

The ratio estimator of the markup is the ratio of the output elasticity for a flexible input to that input's cost share in total revenue. We highlight identification and estimation issues pertaining to this ratio estimator, when firm-level output prices are not observed. If the revenue elasticity for a flexible input is used in place of the output elasticity, then profit maximization implies that the ratio estimator is identically equal to one, and thus is uninformative about markups. Concerning estimation of output elasticities: with only revenue data, profit maximization also implies that the output

---

elasticity is not identified non-parametrically from estimation of the revenue production function, if firms have market power. Even with separate output price and quantity data, it is challenging to estimate the output elasticity consistently if there are non-linear productivity dynamics and firms face heterogeneous demand schedules, with unobserved variation in a demand shifter.

## 1. Introduction

The production approach to markup estimation identifies a firm's markup as the ratio of the output elasticity for a flexible input to that input's cost share in total revenue. We refer to this estimator of the markup as the *ratio estimator*.[1] This paper evaluates the usefulness of the ratio estimator of the markup in settings in which the empirical measure of output is revenue, rather than physical quantity, and firms have market power in output markets.

The production approach was pioneered by Hall (1986, 1988), in his estimates of aggregate industry-level markups. The recent literature extends the Hall methodology to estimate microeconomic firm- or establishment-level markups (see De Loecker and Warzynski (2012), De Loecker et al. (2020), and many others). The microeconomic ratio estimator is widely used across the IO, trade, and macro literatures, and also serves as a popular tool for characterizing the distribution of markups in several economic models, including granular business cycles (Burstein et al., 2020), misallocation in production networks (Baqaee and Farhi, 2020), and monopolistic competition with heterogeneous firms (Mrázová et al., 2018). Many important pitfalls of the ratio estimator have already been discussed (see Traina (2018), Basu (2019), and Syverson (2019)).[2] The issues that we raise in this paper should serve as a

---

[1]Strictly speaking, this is the estimand of the ratio estimator, since the output elasticity is not measured directly and typically has to be estimated.

[2]De Loecker and Goldberg (2014) also discuss the implications of unobserved output and input price heterogeneity across firms for production function estimation in settings

caution against drawing inferences from firm-level markup estimates based on the production approach in settings in which firm-level output prices are unobserved.

When physical output quantities are unobserved, as is the case in most of the papers cited above, it is common practice to proxy output with revenues or value added, deflated with common industry-level price deflators. This approach uses the *revenue* elasticity for a flexible input, in place of the *output* elasticity, in the numerator of the ratio estimator. Klette and Griliches (1996) show that, if firm-level output prices are unobserved and correlated with firms' input choices, then estimators of the revenue elasticity are downward-biased estimators of the output elasticity. We show that the implications of this so-called omitted price bias for identifying markups are much more severe than just generating downward bias in the ratio estimator. Under the standard assumption that the flexible input and the output price are determined from a static profit maximization problem, the ratio estimator that uses the revenue elasticity in place of the output elasticity is identically equal to one, and therefore contains no useful information about markups. We pursue the implications of this observation for the identification and estimation of markups using the ratio estimator of the production approach.

---

with limited heterogeneity in markups (e.g. monopolistic competition with constant elasticity of substitution demand) and discuss partial solutions in these cases. We extend this line of research by studying the implications of omitted output prices for (i) identifying markups using the ratio estimator under general demand conditions and market structures and (ii) estimating output elasticities in settings with imperfect competition and unobserved heterogeneity in demand.

The first part of our paper concerns the identification of markups using the ratio estimator. In Section 2.1, we abstract from estimation issues and suppose that the revenue elasticity and output elasticity are known. We then assess the implications for identifying markups of using one elasticity versus the other in the numerator of the ratio estimator. The main takeaway from this section is that it is essential that the output elasticity, rather than the revenue elasticity, is used in the numerator of the ratio estimator. Even in this best case scenario in which population elasticities are known, replacing the output elasticity with the revenue elasticity removes all information about the markup from the ratio estimator. This result follows from imposing firms' static profit maximization conditions in addition to cost minimization.

In Section 2.2, we raise two additional challenges for identifying markups that arise even when the output elasticity is used in the numerator of the ratio estimator. First, we show that if the input that is used to construct the ratio estimator incurs costs of adjustment, then the ratio estimator reflects the shadow cost of adjusting the input as well as the markup. Second, we show that if the input that is used to construct the ratio estimator is used by firms both to produce output and to influence demand, then the ratio estimator generates a downward-biased estimate of the markup. Such inputs include labor and materials used for marketing, product design, or other sales-related purposes (see Syverson (2011) for a related discussion in the context of productivity estimation).

The second part of our paper concerns the estimation of the output elas-

ticity that is required to identify markups using the ratio estimator. In Section 3.1, we show that in the usual setting in which the researcher observes only revenue, and does not have separate information on the price and quantity of output, the output elasticity for a flexible input is not identified non-parametrically from estimation of the revenue production function, if the flexible input and the output price are determined from a static profit maximization problem. There exist parametric restrictions on the forms of the quantity production function and the inverse demand curve under which the output elasticity for a flexible input may be estimated consistently at one point in the parameter space, but these special cases appear to be of limited empirical relevance. The main takeaway from this section is that firm-level data on output prices is needed to obtain credible estimates of the output elasticity for a flexible input from the estimation of a production function when firms have market power.

We also show that even with firm-level data on output prices, it is still challenging to obtain consistent estimates of output elasticities for flexible inputs, particularly if there are non-linear dynamics in total factor productivity. With such non-linearity, the estimators that are widely used in this context do not estimate the output elasticity for a flexible input consistently if firms face heterogeneous demand curves with *unobserved* variation across firms in a demand shifter, or if only a firm-level price *index* is available. In Section 3.2 we briefly discuss the problem of estimating revenue elasticities; standard estimators do not estimate revenue elasticities consistently in panels

with many firms and few time periods, if there is unobserved heterogeneity across firms in markups.

Overall, the identification and estimation issues that we highlight cast serious doubt over whether anything useful can be learned about trends or heterogeneity in markups from applying the ratio estimator in settings in which output prices and quantities are unobserved.

## 2. Difficulties in Identifying Markups from Production Function Elasticities

In this section, we clarify the conditions under which markups can be identified from knowledge of production function elasticities and the cost shares of flexible inputs in total revenue.

In Section 2.1, we emphasize that knowledge of the *output* elasticity for a flexible input, as opposed to knowledge of the *revenue* elasticity for that input, is essential in this regard. In Section 2.2, we highlight two key assumptions that are required to identify markups. Throughout this section, we abstract from firm heterogeneity in productivity, demand, and input prices; we consider these features in Section 3, where we discuss challenges to estimating the population elasticities that are treated as known in this section.

### 2.1. Output elasticities versus revenue elasticities

We begin by describing the cost minimization problem upon which the production approach to markup estimation is based. Each firm $i$ in period $t$

7

produces output $Q_{it}$ using a production technology with $J$ flexible inputs:

$$Q_{it} = \mathcal{F}\left(X_{it}^1, \ldots, X_{it}^J\right).$$

The only restriction that we place on the production function $\mathcal{F} : \mathbb{R}_+^J \to \mathbb{R}_+$ is that it is twice continuously differentiable in each of the inputs $\left(X_{it}^1, \ldots, X_{it}^J\right)$.[3] Denote by $\boldsymbol{W}_t := \left(W_t^1, \ldots, W_t^J\right)$ the corresponding vector of input prices, over which firms have no influence and take as given. The output elasticity of input $X_{it}^j$ is defined as

$$\theta_{it}^{Q,j} \quad := \quad \frac{X_{it}^j}{Q_{it}} \frac{\partial \mathcal{F}(\cdot)}{\partial X_{it}^j}.$$

We define the estimand of the ratio estimator of the markup using the *output elasticity* $\theta_{it}^{Q,j}$ in the numerator as

$$\mu_{it}^{Q,j} \quad := \quad \frac{\theta_{it}^{Q,j}}{\alpha_{it}^j}$$

where $\alpha_{it}^j := \left(W_t^j X_{it}^j\right) / (P_{it} Q_{it})$ denotes the cost share of input $X_{it}^j$ in total revenue $R_{it} := P_{it} Q_{it}$.

The firm's static cost minimization problem involves choosing its inputs

---

[3]For simplicity, we treat all inputs $\left\{X_{it}^j\right\}_{j=1}^J$ as fully flexible, but this is not essential to the points we make in this section. If a subset of the inputs were fully fixed or predetermined, we could work with the conditional cost function. In Appendix A.2, we show that our main results are robust to a subset of the inputs being partially fixed and subject to adjustment costs.

$\left\{X_{it}^j\right\}_{j=1}^J$ to minimize total (variable) cost subject to producing a target level of output $Q_{it}$

$$\mathcal{C}\left(Q_{it};\boldsymbol{W}_t\right) \quad := \quad \min_{\left\{X_{it}^j\right\}_{j=1}^J}\left\{\sum_{j=1}^J W_t^j X_{it}^j\right\}$$
$$\text{s.t.} \quad \mathcal{F}\left(X_{it}^1,\ldots,X_{it}^J\right) \geq Q_{it}\,.$$

The total cost function $\mathcal{C}\left(Q_{it};\boldsymbol{W}_t\right)$ is the solution to this cost minimization problem. Let $\lambda_{it}$ denote the Lagrange multiplier on the output constraint. The first order condition with respect to the input $X_{it}^j$ is

$$W_t^j \quad = \quad \lambda_{it}\frac{\partial\mathcal{F}\left(\cdot\right)}{\partial X_{it}^j}\,. \tag{1}$$

By the envelope theorem, the Lagrange multiplier, which measures the shadow value of relaxing the output constraint, equals marginal cost

$$\lambda_{it} \quad = \quad \frac{\partial\mathcal{C}\left(\cdot\right)}{\partial Q_{it}}\,.$$

The markup $\mu_{it} := P_{it}/\lambda_{it}$ is defined as the ratio of the output price $P_{it}$ to marginal cost $\lambda_{it}$. Invoking the envelope theorem, multiplying both sides of the cost minimization first order condition (1) by $X_{it}^j$, dividing both sides by $P_{it}Q_{it}$, and rearranging yields

$$\mu_{it}^{Q,j} \quad = \quad \mu_{it}$$

9

for each flexible input $X_{it}^j$. That is, the estimand of the ratio estimator using the output elasticity is equal to the markup, as shown by De Loecker and Warzynski (2012).

Suppose now that the researcher does not have knowledge of the output elasticity $\theta_{it}^{Q,j}$, but rather only has knowledge of the revenue elasticity, defined as

$$\theta_{it}^{R,j} \ := \ \frac{X_{it}^j}{R_{it}}\frac{\partial R_{it}}{\partial X_{it}^j}.$$

We define the estimand of the analogous ratio estimator of the markup using the *revenue elasticity* $\theta_{it}^{R,j}$ in the numerator as

$$\mu_{it}^{R,j} \ := \ \frac{\theta_{it}^{R,j}}{\alpha_{it}^j}.$$

We assess whether the ratio estimator using the revenue elasticity $\theta_{it}^{R,j}$ in place of the output elasticity $\theta_{it}^{Q,j}$ is informative about the markup. In this section, we focus on monopolistic competition to illustrate our main theoretical result in the simplest of market structures with imperfect competition. In Appendix A.1, we show that our result is robust to both Bertrand and Cournot forms of oligopolistic competition with strategic interactions between firms.

A monopolistic firm faces an arbitrary inverse demand schedule:

$$P_{it} \ = \ \mathcal{P}\left(Q_{it}\right).$$

10

The only restrictions that we impose on the function $\mathcal{P} : \mathbb{R}_+ \to \mathbb{R}_+$ are that it is twice continuously differentiable and non-increasing in $Q_{it}$. The absolute value of the price elasticity of demand is defined as

$$\eta_{it} \quad := \quad \left| \frac{P_{it}}{Q_{it}} \frac{\partial Q_{it}}{\partial P_{it}} \right| > 1 \, .$$

We can explicitly write the revenue elasticity $\theta_{it}^{R,j}$ in terms of the demand elasticity $\eta_{it}$ and the output elasticity $\theta_{it}^{Q,j}$

$$\theta_{it}^{R,j} \quad = \quad \left( \frac{\eta_{it} - 1}{\eta_{it}} \right) \theta_{it}^{Q,j} \, . \tag{2}$$

Monopolistic firms with market power in output markets face a finite demand elasticity $\eta_{it} < \infty$. It is then apparent from equation (2) that the revenue elasticity $\theta_{it}^{R,j}$ is strictly less than the output elasticity $\theta_{it}^{Q,j}$.[4]

Taking the total cost function $\mathcal{C}(Q_{it}; \boldsymbol{W}_t)$ from cost minimization as given, the static profit maximization problem involves choosing the output quantity $Q_{it}$ to maximize profits subject to the demand schedule

$$\max_{Q_{it}} \{ P_{it} Q_{it} - \mathcal{C}(Q_{it}; \boldsymbol{W}_t) \}$$

$$\text{s.t.} \quad P_{it} = \mathcal{P}(Q_{it}) \, .$$

---

[4]Intuitively, a monopolistic firm facing a strictly downward-sloping demand schedule must reduce its output quantity to increase its price. Hence, the revenue elasticity is strictly less than the output elasticity.

The first order condition from profit maximization recovers the markup $\mu_{it}$ as a function of the demand elasticity $\eta_{it}$

$$\mu_{it} \;=\; \frac{\eta_{it}}{\eta_{it} - 1} \,.\tag{3}$$

Imposing both cost minimization and profit maximization, we obtain

$$
\begin{aligned}
\mu_{it}^{R,j} \;&=\; \frac{\theta_{it}^{R,j}}{\alpha_{it}^{j}} \\
&=\; \left(\frac{\eta_{it} - 1}{\eta_{it}}\right) \frac{\theta_{it}^{Q,j}}{\alpha_{it}^{j}} \\
&=\; \frac{1}{\mu_{it}} \mu_{it} \\
&=\; 1 \,.
\end{aligned}
$$

The first equality follows from the definition of $\mu_{it}^{R,j}$, the second from equation (2), and the third from the first order conditions from cost minimization (1) and profit maximization (3). It is apparent that the ratio estimator using the revenue elasticity contains no useful information about the markup, except in the very special case of perfect competition under which the markup equals one.[5]

The result $\mu_{it}^{R,j} = 1$ is a consequence of profit maximization and, importantly, does not depend on the particular details of the profit maximization problem, such as the functional form of the demand schedule or the mar-

---

[5]Under perfect competition, the revenue elasticity equals the output elasticity because firms have no influence over output prices.

ket structure. To understand why, consider an industry with $N$ competing firms indexed by $i \in \{1, \ldots, N\}$. Let $\boldsymbol{Q}_{-it} := \{Q_{kt}\}_{k \neq i}$ denote the vector of quantities of all $(N-1)$ competitors of firm $i$, and let $\boldsymbol{Q}_t := \left( Q_{it}, \boldsymbol{Q}_{-it} \right)$. Consider an arbitrary inverse demand schedule $P_{it} = \mathcal{P}_i \left( Q_{it}, \boldsymbol{Q}_{-it} \right)$, for $i \in \{1, \ldots, N\}$, constituting a one-to-one mapping between any vector of quantities $\boldsymbol{Q}_t$ and corresponding vector of prices $\boldsymbol{P}_t := (P_{1t}, \ldots, P_{Nt})$.[6] Let $R_{it} = \mathcal{R}_i \left( Q_{it}, \boldsymbol{Q}_{-it} \right) := \mathcal{P}_i \left( Q_{it}, \boldsymbol{Q}_{-it} \right) Q_{it}$ denote the revenue function for firm $i$ in period $t$. This formulation allows for a range of market structures, including monopolistic competition (as $N \to \infty$) and both Bertrand and Cournot forms of oligopolistic competition (for a finite $N$).

Imposing only cost minimization, we have

$$
\begin{aligned}
\mu_{it}^{R,j} &= \frac{\theta_{it}^{R,j}}{\alpha_{it}^j} \\
&= \frac{X_{it}^j}{R_{it}} \frac{\partial \mathcal{R}_i \left( \cdot \right)}{\partial X_{it}^j} \frac{1}{\alpha_{it}^j} \\
&= \frac{d\mathcal{R}_i \left( \cdot \right)}{dQ_{it}} \frac{1}{P_{it}} \frac{X_{it}^j}{Q_{it}} \frac{\partial \mathcal{F} \left( \cdot \right)}{\partial X_{it}^j} \frac{1}{\alpha_{it}^j} \\
&= \frac{d\mathcal{R}_i \left( \cdot \right)}{dQ_{it}} \frac{1}{P_{it}} \frac{\theta_{it}^{Q,j}}{\alpha_{it}^j} \\
&= \frac{d\mathcal{R}_i \left( \cdot \right)}{dQ_{it}} \frac{\mu_{it}}{P_{it}} \\
&= \frac{d\mathcal{R}_i \left( \cdot \right)}{dQ_{it}} \left[ \frac{\partial \mathcal{C} \left( \cdot \right)}{\partial Q_{it}} \right]^{-1}
\end{aligned}
$$

---

[6]We rule out demand schedules for which a flexible input $X_{it}^j$ enters as an additional argument of the function $\mathcal{P}_i \left( \cdot \right)$. We study the implications of such cases for identifying markups in Section 2.2.

for each flexible input $X_{it}^j$. The first equality follows from the definition of $\mu_{it}^{R,j}$, the second from the definition of $\theta_{it}^{R,j}$, the third from the chain rule and the definition of $R_{it} = P_{it}Q_{it}$, the fourth from the definition of $\theta_{it}^{Q,j}$, the fifth from cost minimization, and the sixth from the definition of the markup $\mu_{it}$. Under cost minimization, but not profit maximization, the estimand $\mu_{it}^{R,j}$ is equal to the ratio of marginal revenue $\frac{d\mathcal{R}_i(\cdot)}{dQ_{it}}$ to marginal cost $\frac{\partial\mathcal{C}(\cdot)}{\partial Q_{it}}$. Additionally imposing profit maximization $\frac{d\mathcal{R}_i(\cdot)}{dQ_{it}} = \frac{\partial\mathcal{C}(\cdot)}{\partial Q_{it}}$ then gives $\mu_{it}^R = 1$, for all values of the true underlying markup $\mu_{it}$.

*Discussion.* For profit maximizing firms, we have established that the estimand of the ratio estimator using the revenue elasticity for a flexible input, in place of the output elasticity, contains no information about the markup. Intuitively, the output elasticity and the revenue elasticity are only equal when a firm is not able to influence its output price by varying its output quantity. But the ability to affect price by changing quantity is the reason why firms with market power charge markups above one.

Our contribution is closely related to Klette and Griliches (1996), who showed that using revenue in place of output to estimate an output elasticity results in a downward bias when firms have market power. In our simple example, this effect is readily seen from equation (2), together with the typical assumption that demand curves slope downward. Since the ratio estimator should use the output elasticity in the numerator, Klette and Griliches (1996) is often cited as a reason why using revenue elasticities instead of output elasticities leads to downward-biased estimates of the markup (see for

example De Loecker and Warzynski (2012), Section VI). While this is true in a technical sense if the true markup is above one, our result shows that the problem is more fundamental. The bias in the ratio estimator from using the revenue elasticity, in place of the output elasticity, removes all the information about the markup, so that the biased estimator is not informative about the markup at all.

Unfortunately, output quantities $Q_{it}$ are rarely observed for individual firms. Instead, researchers typically only have access to measures of revenues $R_{it}$. As we explain in Section 3, when firms have market power, it is not possible to learn about the output elasticity $\theta_{it}^{Q,j}$ by estimating a production function specification that uses revenue as the dependent variable, under any reasonable assumptions (and it is challenging even with data on output quantities). With only data on revenues, it is not clear that we can learn anything about the level of markups using the ratio estimator.

Finally, it is useful to bear in mind that if it were somehow possible to recover the output elasticity from knowledge of the revenue elasticity, then it would not be necessary to use the ratio estimator to learn about markups. Under monopolistic competition, one could simply estimate both the output elasticity and the revenue elasticity, and note from equations (2) and (3) that the ratio of these two elasticities is an estimator of the markup. This observation is a reminder that the problem with revenue elasticities that we are highlighting in this section is not one of estimation, but one of identification: any attempt to learn about the output elasticity from the revenue elastic-

ity must implicitly have assumed knowledge of the markup. The resulting output elasticity therefore cannot contain any additional information that is useful in identifying markups.

Since the estimand underlying the ratio estimator is unity when the revenue elasticity is used in the numerator, it is natural to ask why existing empirical work using this approach does not find estimates that are centered around one. In the following section, we mention two additional sources of bias in the ratio estimator that are likely to be reflected in these estimates. In Section 3.2, we explain why even estimates of the revenue elasticity are likely to be biased. Given these sources of bias, it is not surprising that estimates using the ratio estimator obtained with revenue data are not centered around one.

*2.2. Two additional difficulties in the interpretation of the ratio estimator*

In Section 2.1, we showed that if the revenue elasticity is used in the numerator of the ratio estimator, then the resulting estimand is equal to unity, and contains no information about the markup under the assumption of profit maximization. But when the output elasticity is used in the numerator of the ratio estimator, the resulting estimand correctly recovers the markup under the assumption of cost minimization. In this section, we offer two caveats to this result that apply even in the more favorable case when the output elasticity is known: (i) input adjustment costs, and (ii) inputs that are used not only for production, but also to influence demand.

16

*Input adjustment costs.* For the ratio estimator to recover the markup, it is crucial that the input $X_{it}^j$ whose output elasticity and cost share are combined is perfectly flexible. Alternatively, as explained in Basu (2019), $X_{it}^j$ could be a bundle of inputs, of which at least one component is perfectly flexible, with the other components being fully fixed. However, in reality, inputs rarely fall into one of these two extreme cases. A more realistic intermediate case is to assume that inputs are partially adjustable, in the sense that firms incur costs to adjust their input choices. If the ratio estimator is constructed using an input $X_{it}^j$ that is partially adjustable, or using a bundle that contains partially adjustable inputs, then the estimand of the ratio estimator will reflect both the markup and the shadow cost of adjusting those inputs.

To illustrate this point, assume instead that each input $X_{it}^j$ is associated with a baseline quantity $\overline{X}_{it}^j$ and that the firm incurs adjustment costs when it chooses a quantity of input $X_{it}^j \neq \overline{X}_{it}^j$. The baseline quantity $\overline{X}_{it}^j$ might reflect the input choice from the previous period in a dynamic version of the model. For simplicity, we assume that these costs are given by the smooth convex function $\kappa^j\left(X_{it}^j\right)$, which satisfies $\kappa^j\left(\overline{X}_{it}^j\right) = 0$ and $\frac{d\kappa^j\left(\overline{X}_{it}^j\right)}{dX_{it}^j} = 0$. In Appendix A.2 we show that the estimand using the revenue elasticity is then

$$\mu_{it}^{R,j} = \frac{\theta_{it}^{R,j}}{\alpha_{it}^j} = 1 + \frac{d\kappa^j\left(X_{it}^j\right)}{dX_{it}^j},$$

and the estimand using the output elasticity is

$$\mu_{it}^{Q,j} = \frac{\theta_{it}^{Q,j}}{\alpha_{it}^{j}} = \mu_{it} \left[ 1 + \frac{d\kappa^{j}\left(X_{it}^{j}\right)}{dX_{it}^{j}} \right].$$

Thus, even if the output elasticity for an input were known, it is crucial that none of the inputs in the bundle incur adjustment costs, in order for the ratio estimator to recover the markup.[7]

*Inputs that influence demand.* The framework in Section 2.1 assumed that the inputs $X_{it}^{j}$ are only used to produce output and not also to influence demand. Assume instead that the firm's revenue is given by

$$R_{it} := \mathcal{P}\left(Q_{it}, D_{it}\right) Q_{it}$$

where $D_{it}$ is an endogenous demand shifter that the firm can influence through the use of inputs according to the function

$$D_{it} = \mathcal{D}\left(X_{it}^{D,1}, \ldots, X_{it}^{D,J}\right)$$

---

[7]These results assume that observed input costs are $W_t^{j} X_{it}^{j}$ rather than $W_t^{j} X_{it}^{j} + W_t^{j} \kappa^{j}\left(X_{it}^{j}\right)$. If observed input costs also include the adjustment costs, then we would obtain $\mu_{it}^{Q,j} = \mu_{it} \left( \frac{X_{it}^{j} + \frac{d\kappa^{j}\left(X_{it}^{j}\right)}{dX_{it}^{j}}}{X_{it}^{j} + \kappa^{j}\left(X_{it}^{j}\right)} \right)$, which also does not recover the true markup.

18

where $X_{it}^{D,j}$ is the amount of input $j$ used in influencing demand, and $X_{it}^{Q,j}$ is the amount of input $j$ used in production. We assume that we can observe only the total quantity of input $j$ used by the firm $X_{it}^j = X_{it}^{Q,j} + X_{it}^{D,j}$. In Appendix A.3, we show that the estimand underlying the ratio estimator using the output elasticity then becomes

$$\mu_{it}^{Q,j} = \mu_{it} \left[ \frac{\psi_{it}^{Q,j}}{1 + \frac{X_{it}^{D,j}}{X_{it}^{Q,j}}} \right]$$

where $\psi_{it}^{Q,j}$ is the elasticity of $X_{it}^{Q,j}$ with respect to $X_{it}^j$ evaluated at the optimum, which shows how an additional unit of $X_{it}^j$ is allocated between $X_{it}^{Q,j}$ and $X_{it}^{D,j}$. So if the flexible input is only used for production and not to influence demand (i.e. $\psi_{it}^{Q,j} = 1, X_{it}^{D,j} = 0$), then the ratio estimator using the output elasticity recovers the markup. But if some of the input is used to influence demand, and this component cannot be separated out, then the ratio estimator will be biased. If the firm uses a constant fraction of the input $X_{it}^j$ for production, then $\psi_{it}^{Q,j} = 1$ and the ratio estimator is biased downward. For example, if, over time, the input $X_{it}^j$ is increasingly being used to influence demand, then the ratio estimator will fall, without any change in the true markup.

## 3. Difficulties in Estimating Production Function Elasticities when Firms have Market Power

In Section 2, we established that when using the ratio estimator to estimate markups, it is critical to use the *output* elasticity with respect to a flexible input in the numerator, rather than the *revenue* elasticity. In this section, we highlight several difficulties that arise when attempting to estimate the required output elasticity when firms have market power. We also note that it is not straightforward to obtain consistent estimates of the revenue elasticity, particularly if there is unobserved heterogeneity across firms in markups.

### 3.1. Estimation of the Output Elasticity for a Flexible Input

We start in Section 3.1.1 by considering the case in which the researcher observes only revenue, and does not have separate information on the price and quantity of output. We show that in this case the output elasticity for a flexible input is not identified non-parametrically from estimation of the revenue production function. There exist parametric restrictions on the forms of the quantity production function and the inverse demand curve under which the output elasticity for a flexible input may be estimated consistently at one point in the parameter space, but these special cases appear to be of limited empirical relevance for studying heterogeneity in markups.

In Section 3.1.2, we then consider the case in which the researcher observes both revenue and the output price for individual firms, or equivalently

has data on output quantities. In this case the output elasticity for a flexible input is identified under reasonable conditions if there is no measurement error in the data on output, or if total factor productivity follows a linear ARMA process. In these cases, output elasticities can be estimated consistently using moment conditions for the quantity production function of the kind suggested by Blundell and Bond (2000).

Even with output quantity data, consistent estimation of the output elasticity for a flexible input is more challenging if output is measured with error and total factor productivity follows a non-linear process. Two stage estimators, of the type suggested by Ackerberg et al. (2015) for the estimation of value added production functions for price-taking firms, have often been used in this context.[8] The measurement error in observed output is eliminated using a first stage regression, which allows non-linear dynamic processes for unobserved total factor productivity to be considered in the second stage. The first stage specification requires a valid control function for total factor productivity, which is obtained by inverting a demand function for the flexible input in which total factor productivity is the *only* unobserved component. This approach cannot be used if the demand curves are firm-specific and there is some unobserved heterogeneity across firms in a demand shifter, as well as in total factor productivity, unless the researcher can also control for variation across firms in marginal costs.[9]

_____

[8]See, for example, De Loecker and Warzynski (2012) and De Loecker et al. (2020).

[9]This point has also been made in contemporaneous work by Doraszelski and Jauman-

In Section 3.1.3, we consider the case in which the researcher observes both revenue and a firm-specific output price index, but does not have data on output price levels for individual firms. Deflating revenue using the firm-specific output price index results in a measure of output which differs from the true level of output by an unknown multiplicative firm-specific constant, reflecting differences across firms in output prices in the base year. In logarithmic specifications, this measurement error can be accounted for by firm-specific fixed effects, but obtaining consistent estimates of output elasticities then requires these fixed effects to be taken into account. This is also problematic if we need to deal with non-linearity in the dynamic process for total factor productivity. The presence of unobserved firm-specific fixed effects can however be handled if total factor productivity follows a linear ARMA process, using the kind of dynamic panel data estimator for production functions suggested by Blundell and Bond (2000).

### 3.1.1. Data on Revenue

In this section we consider a three factor Hicks-neutral gross output production function for firm $i$ in period $t$ of the form

$$q_{it} = f(k_{it}, l_{it}, m_{it}) + \omega_{it} \qquad (4)$$

---

dreu (2019).

in which $q_{it}$ is the log of gross output, $k_{it}$, $l_{it}$ and $m_{it}$ are the logs of observed capital, labor and intermediate inputs respectively, and $\omega_{it}$ is the log of total factor productivity, which is observed by the firm but not by the researcher. We treat capital and labor as predetermined inputs, for which the input levels are chosen before the firm has observed $\omega_{it}$.[10] We assume that the level of intermediate inputs is chosen after the firm has observed $\omega_{it}$, and that intermediate inputs do not incur adjustment costs of any kind; that is, we consider intermediate inputs as our example of an input which is flexible in the sense required to construct the ratio estimator of the markup. The object of interest is thus the output elasticity $\theta_{it}^{Q,M} := \partial f\left(\cdot\right)/\partial m_{it}$.

The researcher observes neither gross output nor the output price, but only sales revenue or the value of gross output, the log of which is $r_{it} := p_{it} + q_{it}$.[11] To analyze this further, we assume that each firm faces a downward-sloping inverse demand curve of the form

$$p_{it} = p(q_{it}, \xi_{it}) \tag{5}$$

in which $\xi_{it}$ is a demand shifter, which is observed by the firm and may be observed or unobserved by the researcher.

The revenue production function which can be estimated in this setting

---

[10]The predetermined inputs may also be subject to adjustment costs. If so, these adjustment costs do not take the form of foregone production, and do not depend on the level of intermediate inputs in any time period.

[11]We abstract here from any difference between sales revenue and the value of production, due to changes in inventories.

relates the log of observed revenue to the logs of the observed inputs

$$r_{it} = (p_{it} + q_{it}) = f(k_{it}, l_{it}, m_{it}) + (p_{it} + \omega_{it}) . \qquad (6)$$

The dependence of intermediate inputs $(m_{it})$ on unobserved total factor productivity $(\omega_{it})$ raises issues for the consistent estimation of the output elasticity $\theta_{it}^{Q,M}$ from the quantity production function (4) that are well known in the context of price-taking firms; we discuss some additional issues which arise when firms have market power in section 3.1.2 below.

The presence of the output price $(p_{it})$ in the error term of the revenue production function (6) raises more fundamental issues when firms have market power, and their output price depends on $q_{it}$ from (5), and hence on each of the inputs. This additional source of inconsistency has been analyzed by Klette and Griliches (1996) and termed the 'omitted price bias'.[12] Our contribution here is to show that if the output price and the level of the flexible input are chosen at the same time to maximize the same objective, then the output elasticity $\theta_{it}^{Q,M}$ is not identified non-parametrically from estimation of the revenue production function (6).

The intuition for this result is straightforward in the special case in which all firms face the same inverse demand curve, and we have only common shocks $(\xi_{it} = \xi_t$ for all $i)$ in (5). In this case, with observations on $(p_{it}, q_{it})$ constrained to lie along this downward-sloping demand curve, any firm-specific

---

[12]See also De Loecker (2011) and De Loecker and Goldberg (2014).

shock which increases $m_{it}$ and hence $q_{it}$ also reduces $p_{it}$. In other words, any informative instrument for $m_{it}$ is correlated with $p_{it}$ and so not a valid instrument in the revenue production function (6). With heterogeneity across firms in the inverse demand curves, the same still applies, except in special cases in which there is no pass through of demand shocks $(\xi_{it})$ to the output price. In these special cases, if informative proxies for the demand shifter are observed by the researcher, and these are uncorrelated with $\omega_{it}$, then these would provide valid and informative instruments for $m_{it}$ in (6). However, the special cases with zero pass through of demand shocks to the output price require strong parametric restrictions on the form of both the quantity production function (4) and the inverse demand curve (5), such that at best the output elasticity is identified only at one point in the parameter space.

To illustrate this, we assume that the firm chooses its output price $(P_{it})$ and level of intermediate inputs $(M_{it})$ to maximize profits in period $t$, taking the costs of the predetermined inputs as given, or equivalently to maximize revenue net of variable costs

$$\Pi_{it} := P_{it}Q_{it} - P_{it}^M M_{it} \tag{7}$$

subject to the constraints in (4) and (5). Here $P_{it}^M$ is the price of one unit of intermediate inputs for firm $i$ in period $t$, and $p_{it}^M$ is the log of this price; the input price is observed by the firm, and may be observed or unobserved by the researcher. We assume that the firm takes total factor productivity

25

$(\omega_{it})$, the demand shifter $(\xi_{it})$, and the flexible input price $(p_{it}^M)$ as given.

The solution equates marginal revenue and marginal variable cost. We can either find the level of intermediate inputs which maximizes net revenue in period $t$ and infer the output price from the inverse demand curve at the resulting level of output, or we can find the output price and quantity which maximize net revenue in period $t$ and infer the required level of intermediate inputs. In either case, we obtain decision rules or policy functions which express both $m_{it}$ and $p_{it}$ as functions of the *same* state variables $(k_{it}, l_{it})$ and the *same* primitives $(\omega_{it}, \xi_{it}, p_{it}^M)$:

$$m_{it} = m^*(k_{it}, l_{it}, \omega_{it}, \xi_{it}, p_{it}^M) \tag{8}$$
$$p_{it} = p^*(k_{it}, l_{it}, \omega_{it}, \xi_{it}, p_{it}^M).$$

These decision rules then indicate that any informative instrument for $m_{it}$ in (6) will necessarily be correlated with the $p_{it}$ component of the error term, while any instrument that is uncorrelated with $p_{it}$ will not be an informative instrument for $m_{it}$. Equivalently, if we were able to control adequately for the $p_{it}$ component of the error term in (6) we would have exhausted all the sources of variation in the explanatory variable $m_{it}$. The explanatory variable $m_{it}$ and the error component $p_{it}$ are 'functionally dependent' in the sense of Ackerberg et al. (2015). Without parametric restrictions, we cannot separately identify the contributions of $m_{it}$ and $p_{it}$ to the log of observed

revenue $r_{it}$.[13]

In this context, variation in the input price $p_{it}^M$ shifts the marginal variable cost schedule; if the demand and marginal revenue schedules are downward-sloping, this variation necessarily also affects the output price. As a result, there are no parametric restrictions that lead to the exclusion of $p_{it}^M$ from the decision rule for the output price in (8). The demand shocks $\xi_{it}$ shift the marginal revenue schedule, and there are admissible parametric restrictions under which there is zero pass through of the demand shocks to the output price. This would be the case if we have both constant marginal variable cost and the markup does not depend on the level of output.

For example, we may have a Cobb-Douglas gross output production function with increasing returns to scale and a unit output elasticity for the flexible input, and a Constant Elasticity of Substitution (CES) demand curve for each firm.[14] In this case, the demand shocks $\xi_{it}$ affect the level of intermediate inputs but not the output price, and observed proxies for the demand shocks would provide valid and informative instruments for $m_{it}$ in a log-linear version of (6), provided they are also uncorrelated with $\omega_{it}$. This *requires* heterogeneity across firms in the inverse demand curves, and the

---

[13]The dependence of the output price on the predetermined inputs also indicates that when firms have market power, we do not have moment conditions of the form $\mathrm{E}[(p_{it} + \omega_{it})|k_{it}, l_{it}] = 0$, versions of which have typically been used in the estimation of revenue production functions.

[14]That is, we have a gross output production function of the form $q_{it} = \beta_K k_{it} + \beta_L l_{it} + \beta_M m_{it} + \omega_{it}$ with $\beta_M = \theta_{it}^{Q,M} = 1$ for all $i, t$, and returns to scale $\nu = \beta_K + \beta_L + 1 > 1$; and an inverse demand curve of the form $p_{it} = \xi_{it} - \eta^{-1} q_{it}$, where $\eta = \eta_{it} > 1$ is the absolute value of the price elasticity of demand for all $i, t$.

output elasticity parameter for the flexible input ($\beta_M$) is identified *only* at one point ($\beta_M = 1$) in the parameter space. This requirement for the output elasticity to be unity here suggests that these parametric special cases are likely to be of limited empirical relevance. Moreover, since identification here relies on shifts in the demand curve, and shifts in the demand curve would affect the demand for two or more flexible inputs in the same way, the parametric special cases in which this approach could be applied are limited to specifications with a single flexible input, as in the example that we have considered here.[15]

### 3.1.2. Data on Revenue and Output Price Levels

Our result in the previous section indicates that, when firms have market power, data on firm-level output prices is fundamental to obtaining credible estimates of the output elasticity for a flexible input from estimation of a production function. Here we show that even with a quantity measure of output, it is still challenging to estimate this output elasticity consistently, particularly if output is measured with error and total factor productivity follows a non-linear dynamic process.

---

[15]Note that our results in this section apply to a revenue production function which relates revenue to input *quantities*, as in (6). If the specification relates revenue to *expenditures* on (some of) the inputs, so that (some of) the input prices are introduced as additional components of the error term, there may also be parametric special cases in which the output elasticities could be estimated consistently. One example has a constant returns to scale Cobb-Douglas gross output production function, a CES demand curve with the same demand elasticity for all firms, all inputs fully flexible, and all inputs with heterogeneous input prices measured as expenditures.

To simplify the exposition, we now focus on a Cobb-Douglas gross output production function, although the issues we highlight apply for any continuously differentiable gross output production function (see Appendix B.1 for details). We further assume that gross output is measured with a multiplicative error, such that the log of observed output is $y_{it} := q_{it} + \varepsilon_{it}$, where $\varepsilon_{it}$ is a mean zero measurement error. The quantity production function to be estimated then has the form

$$y_{it} = \beta_K k_{it} + \beta_L l_{it} + \beta_M m_{it} + (\omega_{it} + \varepsilon_{it}). \tag{9}$$

For simplicity, we choose units such that the mean of $\omega_{it}$ is also zero. We assume that the measurement error $\varepsilon_{it}$ is uncorrelated with the observed inputs $(k_{is}, l_{is}, m_{is})$ and with the input price $p_{is}^M$ for any $s, t$, and is independent across firms.[16] The slope parameters $(\beta_K, \beta_L, \beta_M)$ are the output elasticities, which are assumed to be constant over time and common to all the firms in the sample. Our parameter of interest here is the output elasticity for the flexible input $\beta_M$.

We again assume that the firm chooses the level of intermediate inputs to maximize net revenue in (7), subject to the constraints in (9) and (5)

---

[16]An alternative interpretation of the two error components in (9) is that $\omega_{it}$ denotes the log of the component of total factor productivity that is known by the firm when making input decisions in period $t$, and $\varepsilon_{it}$ denotes the log of an unforecastable productivity shock that is not known by the firm when making input decisions in period $t$. The presence of the second component ($\varepsilon_{it}$) of the error term here is more important than the particular way we introduce it.

and taking $(\omega_{it}, \xi_{it}, p_{it}^M)$ as given. Without specifying the form of the inverse demand curve (5), we show in Appendix B.1 that the optimal choice of intermediate inputs satisfies the first order condition

$$m_{it} = \frac{\ln \beta_M}{1 - \beta_M} + \left( \frac{\beta_K}{1 - \beta_M} \right) k_{it} + \left( \frac{\beta_L}{1 - \beta_M} \right) l_{it} + \left( \frac{1}{1 - \beta_M} \right) \left( p_{it} - \ln \mu_{it} - p_{it}^M + \omega_{it} \right)$$

(10)

where $\mu_{it}$ is the markup of price over marginal cost as in Section 2, and we can note that $z_{it} := p_{it} - \ln \mu_{it}$ is the log of marginal cost. The only restriction that we place on the demand curve here is that the output price $p_{it}$ is a weakly decreasing function of gross output $q_{it}$.

We assume that total factor productivity $\omega_{it}$ is independent across firms, and start by considering the special case in which $\omega_{it}$ is serially uncorrelated; extensions to more realistic cases in which the unobserved heterogeneity across firms in productivity is persistent over time will be considered below. We consider a setting in which panel data is observed for a large number of firms for a small number of time periods, and asymptotic properties are stated for the case in which the number of firms increases, with the number of time periods treated as fixed.

Under these assumptions, we have the moment conditions $\mathrm{E}[(k_{it}, l_{it})u_{it}] = 0$ where $u_{it} := \omega_{it} + \varepsilon_{it}$ is the error term in (9). If the researcher has data on the input price $p_{it}^M$, and if these input prices vary across firms in a way that is uncorrelated with $\omega_{it}$, then the price of the flexible input provides

30

a valid and informative instrument for the explanatory variable $m_{it}$ in (9).
In that case we have the additional moment condition $\text{E}[p_{it}^M u_{it}] = 0$, and
the parameter vector $(\beta_K, \beta_L, \beta_M)$ is identified from the estimation of the
quantity production function (9).

If the researcher does not have data on the price of the flexible input,
or if the variation across firms in these input prices is correlated with $\omega_{it}$,
the parameter vector $(\beta_K, \beta_L, \beta_M)$ will still be identified here if either: (i)
there is variation across firms in the input price $p_{it}^M$ which is persistent over
time; or (ii) there is variation across firms in the demand shifter $\xi_{it}$ which is
persistent over time and results in persistent variation in the log of marginal
cost $z_{it}$. With persistent variation in either $p_{it}^M$ or $z_{it}$, the first order condition
(10) implies that the lagged input $m_{i,t-1}$ provides a valid and informative
instrument for the explanatory variable $m_{it}$ in (9), and in this case we have
the additional (informative) moment condition $\text{E}[m_{i,t-1} u_{it}] = 0$.[17]

For price-taking firms, it is well known that identification of the output
elasticity for a flexible input from estimation of the quantity production
function requires variation across firms in the price of the flexible input.[18]
For firms with market power and a single flexible input, persistent variation

---

[17]This can also be seen from the decision rule for $m_{it}$ given in (8). We assume here that
the researcher does not observe the demand shifter. If the researcher observes $\xi_{it}$, and $\xi_{it}$
varies across firms in a way which is uncorrelated with $\omega_{it}$, then $\xi_{it}$ could be used as an
instrument for $m_{it}$ in (9), and we would not require the variation across firms in $\xi_{it}$ to be
persistent. The same would apply if the researcher observes an informative proxy for $\xi_{it}$
that is uncorrelated with $\omega_{it}$.

[18]See Bond and Söderbom (2005), Ackerberg et al. (2015) and Gandhi et al. (2020).

across firms in demand provides a second mechanism through which the lagged input may be an informative instrument. This could be useful in applications where the researcher has data on expenditure on the flexible input, but does not have firm-level data on the price of the flexible input. Expenditure on the flexible input, deflated using a common price index, provides a suitable measure of the input quantity only if the input price does not vary across firms. This requirement rules out identification of the output elasticity from estimation of the production function for price-taking firms, but may not do so when firms have market power.

We now extend our discussion to consider more realistic cases in which the variation across firms in unobserved total factor productivity is persistent over time, distinguishing between the cases in which $\omega_{it}$ follows linear and non-linear dynamic processes. In both cases the dynamic process for $\omega_{it}$ has to be correctly specified by the researcher.

*Linear TFP process.* The moment conditions discussed above extend straightforwardly to cases in which $\omega_{it}$ follows a low order ARMA process. Suppose, for example, that $\omega_{it}$ follows an AR(1) process

$$\omega_{it} = \rho\omega_{i,t-1} + \upsilon_{it} \tag{11}$$

with $|\rho| < 1$, in which the productivity innovations $\upsilon_{it}$ are independent across firms and serially uncorrelated. Substituting for $\omega_{it}$ and $\omega_{i,t-1}$ in (11) from (9) results in a quasi-differenced representation of the quantity production

32

function in which the error term is now $u_{it} := v_{it} + \varepsilon_{it} - \rho \varepsilon_{i,t-1}$:

$$(y_{it} - \beta_K k_{it} - \beta_L l_{it} - \beta_M m_{it} - \varepsilon_{it}) = \rho(y_{i,t-1} - \beta_K k_{i,t-1} - \beta_L l_{i,t-1} - \beta_M m_{i,t-1} - \varepsilon_{i,t-1}) + v_{it}$$

$$\Leftrightarrow (y_{it} - \rho y_{i,t-1}) = \beta_K(k_{it} - \rho k_{i,t-1}) + \beta_L(l_{it} - \rho l_{i,t-1}) + \beta_M(m_{it} - \rho m_{i,t-1}) + (v_{it} + \varepsilon_{it} - \rho \varepsilon_{i,t-1}).$$

Here we still have moment conditions of the form $\mathrm{E}[(k_{is}, l_{is})u_{it}] = 0$ for $s \leqslant t$. If the researcher has data on the input price, and the input price is uncorrelated with $\omega_{it}$, we have additional moment conditions $\mathrm{E}[p_{is}^M u_{it}] = 0$ for $s \leqslant t$. If the researcher does not have data on the input price, or if the variation across firms in these input prices is correlated with $\omega_{it}$, but we have persistent variation across firms in either $p_{it}^M$ or $\xi_{it}$, we have additional (informative) moment conditions $\mathrm{E}[m_{is}u_{it}] = 0$ for $s \leqslant t-1$. If the measurement error $\varepsilon_{it}$ is serially uncorrelated, we also have additional moment conditions $\mathrm{E}[y_{is}u_{it}] = 0$ for $s \leqslant t - 2$.[19] These moment conditions can be used to estimate the parameter vector $(\beta_K, \beta_L, \beta_M, \rho)$ consistently in the quasi-differenced quantity production function , following the approach suggested by Blundell and Bond (2000).

*Non-linear TFP process.* Similar moment conditions could be used to estimate the output elasticity parameters consistently with non-linear processes for $\omega_{it}$, if gross output is measured without error and $\omega_{it}$ is the only compo-

_____

[19]This restriction follows naturally if the $\varepsilon_{it}$ component of the error term in (9) is interpreted as a shock to productivity that is not known by the firm when making input decisions in period $t$.

nent of the error term in the quantity production function (9). Otherwise, if we replace the linear AR(1) process (11) by the first-order Markov process

$$\omega_{it} = g(\omega_{i,t-1}) + v_{it} \,, \tag{12}$$

the presence of the unobserved $\varepsilon_{i,t-1}$ inside the non-linear function $g(\omega_{i,t-1})$, when we substitute for $\omega_{i,t-1}$ using (9), will invalidate moment conditions of the kind considered in the previous sub-section.

With a non-linear process for $\omega_{it}$ and measurement error in output, Flynn et al. (2019) have shown that when firms have market power, even with a quantity measure of output and all inputs, gross output production functions with a flexible input are not identified if the decision rule for the flexible input has the form $m_{it} = m_t^*(k_{it}, l_{it}, \omega_{it})$. Comparison to the decision rule for $m_{it}$ in (8) indicates that this assumption rules out variation across firms in both the input price $(p_{it}^M)$ and the demand shifter $(\xi_{it})$.[20] These are the sources of variation that we relied on in the previous sub-section, for identification of the output elasticity for the flexible input $(\beta_M)$ in specifications with linear processes for $\omega_{it}$. Our contribution in this sub-section is to consider whether variation across firms in either $p_{it}^M$ or $\xi_{it}$ would allow this key output elasticity parameter to be estimated consistently, in specifications with a non-linear process for $\omega_{it}$ and measurement error in output.

---

[20]The dependence of $m_t^*(k_{it}, l_{it}, \omega_{it})$ on the time period $t$ allows for common variation over time in both $p_{it}^M$ and $\xi_{it}$.

We still have moment conditions of the form $\mathrm{E}[(k_{is}, l_{is})v_{it}] = 0$ for $s \leqslant t$ and, for example, $\mathrm{E}[m_{is}v_{it}] = 0$ for $s \leqslant t-1$. To exploit these moment conditions, we would first need to eliminate the measurement error component $\varepsilon_{it}$ from the error term of the quantity production function (9), before we substitute for $\omega_{i,t-1}$ in the non-linear function $g(\omega_{i,t-1})$.

A two stage estimation procedure of this kind was proposed by Ackerberg et al. (2015) for the estimation of a value added production function for price-taking firms, and with no flexible inputs. Similar two stage estimators are commonly used in the empirical literature that uses the ratio estimator to study markups.[21] De Loecker and Warzynski (2012) proposed an estimator of this type which can be used when we observe firm-level prices for both output and the flexible input, and we have persistent variation across firms in the price of the flexible input, and no unobserved variation across firms in the demand shifter. However, there are problems in applying this approach to the estimation of a gross output production function when firms have market power and there is unobserved heterogeneity across firms in the demand shifter $\xi_{it}$.

The first stage of these two stage procedures relies on having a valid control function which expresses the unobserved $\omega_{it}$ in (9) as a function of observed variables only. This is obtained by expressing the firm's optimal choice of the flexible input $m_{it}$ as a function of observed variables and the

---

[21]See, for example, De Loecker and Warzynski (2012) and De Loecker et al. (2020).

single unobserved component $\omega_{it}$. We also require that this function is strictly monotonic in $\omega_{it}$, so that it can be inverted to provide the control function. A (possibly non-parametric) regression of $y_{it}$ on the observed inputs and any additional observed variables included in the control function then has the error term $\varepsilon_{it}$. The predicted values of $y_{it}$ from the estimated first stage regression can then be used in place of the actual values of $y_{it}$ when we substitute for $\omega_{it}$ and $\omega_{i,t-1}$ in the specified non-linear dynamic process (12).

The question here is whether we can find a valid control function of this form in settings where we also have informative instruments for $m_{it}$ in the second stage of this procedure. We have the decision rule $m_{it} = m^*(k_{it}, l_{it}, \omega_{it}, \xi_{it}, p_{it}^M)$ obtained in (8). First suppose that the researcher has data on $p_{it}^M$ and all firms face the same demand curve ($\xi_{it} = \xi_t$ for all $i$). Time dummies ($d_t$) can then be used to control for the common demand shocks. The decision rule then depends on the scalar unobservable $\omega_{it}$, and can be inverted to give the valid control function $\omega_{it} = h(k_{it}, l_{it}, m_{it}, p_{it}^M, d_t)$, which can be used in the first stage regression. If the variation in $p_{it}^M$ is un-correlated with $\omega_{it}$, we can also use the observed input prices as instruments for $m_{it}$ in the second stage specification; that is, we have valid and informative moment conditions of the form $\mathrm{E}[p_{is}^M \upsilon_{it}] = 0$ for $s \leqslant t$. If the variation in $p_{it}^M$ is correlated with $\omega_{it}$ but persistent over time, we can instead use lagged intermediate inputs as instruments for $m_{it}$ in the second stage spec-ification; that is, we have valid and informative moment conditions of the form $\mathrm{E}[m_{is} \upsilon_{it}] = 0$ for $s \leqslant t - 1$. Notice that with no heterogeneity across

firms in the demand shifter, we require persistent variation across firms in the input price here; with firm-level data on the input price, this condition can be checked.

Now suppose that the researcher has data on $p_{it}^M$ and there is variation across firms in the demand shifter which is not *perfectly* observed by the researcher (i.e. there is *some* unobserved heterogeneity across firms in $\xi_{it}$). In this case, we can no longer express $m_{it}$ as a function of observed variables and the scalar unobservable $\omega_{it}$. We could still invert the function $m_{it} = m^*(k_{it}, l_{it}, \omega_{it}, \xi_{it}, p_{it}^M)$ to obtain $\omega_{it} = h(k_{it}, l_{it}, m_{it}, \xi_{it}, p_{it}^M)$, but this does not provide a valid control function for $\omega_{it}$ if there is any unobserved variation across firms in the demand shifter $\xi_{it}$.

Similar issues arise if we consider using the first order condition (10) as the basis for obtaining a control function for $\omega_{it}$ in the first stage regression. In this case, we could still invert the function $m_{it} = m(k_{it}, l_{it}, \omega_{it}, z_{it}, p_{it}^M)$ to obtain $\omega_{it} = h(k_{it}, l_{it}, m_{it}, z_{it}, p_{it}^M)$, but with unobserved heterogeneity in $\xi_{it}$, the researcher would need to be able to control for variation in the log of marginal cost $z_{it}$, to obtain a valid control function.[22] Otherwise, with market power and unobserved heterogeneity in demand, we cannot allow for non-linearity in the dynamic process for total factor productivity using a two

---

[22]This has also been noted by Doraszelski and Jaumandreu (2019) in a more general setting than our example here. With no unobserved variation across firms in $\xi_{it}$, we have $z_{it} = z(k_{it}, l_{it}, \omega_{it}, p_{it}^M, d_t)$. Substituting for $z_{it}$ in the first order condition and inverting the resulting function then gives the same control function $\omega_{it} = h(k_{it}, l_{it}, m_{it}, p_{it}^M, d_t)$ that we obtained from the decision rule.

stage procedure of this type, even with firm-level data on the price of the flexible input.[23]

### 3.1.3. Data on Revenue and Output Price Indices

The previous section considered the case in which the researcher has data on both sales revenue and the *level* of the output price for individual firms. An intermediate possibility is that the researcher observes an output price *index* for individual firms, constructed from survey questions about yearly price changes, but does not observe firm-specific price levels in the base year.

If we use these firm-specific output price indices to deflate the value of output in current prices, we obtain

$$P_{i0}Q_{it} := (P_{it}Q_{it}) \times \left( \frac{P_{it}}{P_{i0}} \right)$$

where $(P_{it}/P_{i0})$ is the firm-specific price index, equal to one in the base period $t = 0$, and $P_{i0}$ is the unobserved price of output for firm $i$ in that period.

Deflating revenue in this way measures the true level of output $Q_{it}$ up to the unknown multiplicative firm-specific constant $P_{i0}$, reflecting unobserved differences across firms in the price of output in the base year. In a logarithmic specification, this will introduce firm-specific intercepts. For example,

---

[23]The situation is no better if the researcher does not have data on the price of the flexible input. To obtain a valid control function for $\omega_{it}$ in the first stage regression, we then require no unobserved heterogeneity across firms in $\xi_{it}$ and no variation across firms in $p_{it}^M$. Observed variation in the demand shifter $\xi_{it}$ would then be needed to provide informative instruments for $m_{it}$ in the second stage specification, and this approach could not be used in a specification with two or more flexible inputs.

for the Cobb-Douglas gross output production function considered in the previous section, we obtain from (9)

$$(p_{i0} + y_{it}) = p_{i0} + \beta_K k_{it} + \beta_L l_{it} + \beta_M m_{it} + (\omega_{it} + \varepsilon_{it}) \qquad (13)$$

where again $y_{it} = q_{it} + \varepsilon_{it}$, and $\varepsilon_{it}$ allows for transient measurement error. Persistent differences across firms in the level of the output price will be correlated with input choices, so in the panel data sense these firm-specific intercepts will need to be treated as 'fixed effects' (i.e. correlated with the explanatory variables) rather than 'random effects' (i.e. uncorrelated with the explanatory variables).[24]

In the case where the unobserved total factor productivity component of the error term $\omega_{it}$ follows a low order ARMA process, the 'dynamic panel data' estimator for production functions proposed by Blundell and Bond (2000) can accommodate unobserved firm-specific fixed effects of this form. This allows consistent estimation of the output elasticity parameters $(\beta_K, \beta_L, \beta_M)$ provided that $\omega_{it}$ follows a linear process and either: (i) we have data on $p_{it}^M$, and the input price is uncorrelated with $\omega_{it}$; or (ii) there is persistent variation across firms in either $p_{it}^M$ or $\xi_{it}$, such that lagged inputs provide valid and informative instruments for $m_{it}$. The key point here is that estimation will need to allow for fixed effects if the researcher does not have

---

[24]A similar issue arises if we use an expenditure measure of one or more of the inputs, deflated using a firm-specific input price index, and there is unobserved variation across firms in the level of the input price.

firm-level data on output price levels.

The two stage estimators which have been developed to allow for non-linear dynamics in $\omega_{it}$ cannot allow for unobserved firm-specific fixed effects in $\omega_{it}$, at least in panel data settings with a small number of time periods. It may be possible to extend estimators of this type to allow for unobserved firm-specific fixed effects in the measurement error component of the error term, which is the relevant case here. This could be a useful subject for further research, in settings where we have data on firm-specific price indices but not firm-specific price levels, and are content with the assumption of no unobserved variation across firms in the demand shifter $\xi_{it}$.

*3.2. Estimation of the Revenue Elasticity for a Flexible Input*

In Section 3.1 we showed that the output elasticity for a flexible input is not identified from estimation of the revenue production function without strong parametric restrictions on the forms of both the gross output production function and the inverse demand curve. In this section, we briefly consider conditions under which the revenue elasticity for a flexible input can be estimated consistently.

A useful starting point is the case considered by Klette and Griliches (1996), with a Cobb-Douglas gross output production function (9) and a CES inverse demand curve

$$p_{it} = \delta_t - \eta^{-1} q_{it} + \zeta_{it} \tag{14}$$

in which we have decomposed the demand shifter $\xi_{it}$ into common and idiosyncratic components, such that $\xi_{it} = \delta_t + \zeta_{it}$. Here $\eta > 1$ is the absolute value of the price elasticity of demand. The revenue production function in this case is

$$r_{it}^o = (p_{it} + y_{it}) = \beta_K k_{it} + \beta_L l_{it} + \beta_M m_{it} + (p_{it} + \omega_{it} + \varepsilon_{it}) \qquad (15)$$

where the log of observed revenue $r_{it}^o := r_{it} + \varepsilon_{it}$ differs from the log of true revenue $r_{it}$ by the additive measurement error component $\varepsilon_{it}$.

Substituting for the unobserved output price $p_{it}$ in the error term of (15) from the inverse demand curve (14), we obtain the log-linear equation

$$r_{it}^o = \delta_t + \left(\frac{\beta_K}{\mu}\right) k_{it} + \left(\frac{\beta_L}{\mu}\right) l_{it} + \left(\frac{\beta_M}{\mu}\right) m_{it} + \left[\left(\frac{1}{\mu}\right) \omega_{it} + \zeta_{it} + \varepsilon_{it}\right] \qquad (16)$$

which relates observed revenue to the observed inputs. Here $\mu = \left(1 - \frac{1}{\eta}\right)^{-1} > 1$ is the markup, and the slope parameters are the *revenue* elasticities. The error term contains the idiosyncratic demand shock $\zeta_{it}$, in addition to total factor productivity $\omega_{it}$ and the measurement error $\varepsilon_{it}$.

The revenue elasticity parameters in (16) can then be estimated consistently using the methods discussed in Section 3.1.2, subject to the limitations that we have noted. For example, if both $\omega_{it}$ and $\zeta_{it}$ are assumed to be serially uncorrelated, we have moment conditions $E[(k_{it}, l_{it})u_{it}] = 0$, where now $u_{it} := \left(\frac{1}{\mu}\right) \omega_{it} + \zeta_{it} + \varepsilon_{it}$. With persistent variation across firms in the input

price $p_{it}^M$, the lagged input $m_{i,t-1}$ provides a valid and informative instrument for $m_{it}$, and we have the additional (informative) moment condition $\mathrm{E}[m_{i,t-1}u_{it}] = 0$. This extends straightforwardly to cases in which $\omega_{it}$ follows a low order ARMA process, although not to cases in which $\omega_{it}$ follows a non-linear dynamic process (if we do indeed have both unobserved idiosyncratic demand shocks and measurement error).

In cases where we can estimate these revenue elasticity parameters consistently, we could investigate heterogeneity in the markup parameter $\mu$ across (large) sub-samples of firms by including suitable interaction terms in (16), under the maintained assumption that the output elasticities are common to these sub-samples.[25]

This example also highlights potential problems with estimating the revenue elasticities consistently. Consistent estimation in the example considered above required the researcher to observe a quantity measure of the flexible input.[26] More generally, consistent estimation may be difficult if the sum $\left(\frac{1}{\mu}\right)\omega_{it} + \zeta_{it}$ does not follow a low order ARMA process. Consistent estimation may also be difficult if the markup parameter $\mu$ is not common within (large) sub-samples of firms. The moment conditions that are typically used

---

[25]For example, we could investigate if the revenue elasticity parameters take different values for exporting and non-exporting firms, as in De Loecker and Warzynski (2012).

[26]If the researcher only has data on expenditure on the flexible input, the assumption that the price of the flexible input does not vary across firms then implies that the lagged input is not an informative instrument for the current input, given the levels of the predetermined inputs $k_{it}$ and $l_{it}$, under the maintained assumptions that $\omega_{it}$ and $\zeta_{it}$ are both serially uncorrelated; see (10).

to estimate production functions will not be valid if there is unmodeled heterogeneity in the slope parameters in (16).[27] Finally, consistent estimation of the revenue elasticities is likely to be more difficult if the gross output production function and inverse demand curve do not take the convenient log-linear forms implied by a Cobb-Douglas production technology and a CES demand schedule.

## 4. Conclusion

Our primary objective in this paper is to caution against drawing inferences from firm-level markup estimates based on the production approach, when firm-level output prices are not observed. Static profit maximization conditions imply that when a revenue elasticity is used in place of an output elasticity, the commonly-used ratio estimator contains no useful information about markups. Static profit maximization also implies that the required output elasticity for a flexible input is not identified from estimation of a revenue production function, without placing strong parametric restrictions on the functional forms of both the production function and the demand schedule. We discuss additional problems with the ratio estimator of markups when

---

[27]In the model $y_{it} = \beta x_{it} + u_{it}$ with $\mathrm{E}(u_{it}) = 0$ and $\mathrm{E}(x_{it}u_{it}) \neq 0$, we can obtain consistent estimators of $\beta$ if $\mathrm{E}(x_{i,t-1}u_{it}) = 0$ and $x_{i,t-1}$ is also an informative instrument for $x_{it}$. With heterogeneity across firms in the slope parameter, we have $y_{it} = \beta_i x_{it} + u_{it} = \beta x_{it} + u_{it} + (\beta_i - \beta)x_{it} = \beta x_{it} + e_{it}$, with $e_{it} \coloneqq u_{it} + (\beta_i - \beta)x_{it}$. If the explanatory variable is serially correlated, we then have $\mathrm{E}(x_{i,t-1}e_{it}) \neq 0$, and standard estimators do not estimate $\beta$ consistently. With time-invariant heterogeneity of this form, the $\beta_i$ coefficients (and hence $\beta$) could be estimated consistently if panel data is available for a large number of time periods. See Pesaran and Smith (1995) for further discussion.

the flexible input is used by firms not only to produce output, but also to influence demand; and we show that even with separate data on output prices and quantities, it is still challenging to estimate the output elasticity for a flexible input consistently, if there are non-linear productivity dynamics and firms face heterogeneous demand schedules, with unobserved heterogeneity across firms in a demand shifter.

These difficulties notwithstanding, the clear implication of our main results is that firm-level data on output prices are required to obtain credible estimates of markups using the ratio estimator. With revenue data alone, we are not aware of any procedures that would allow the level of markups to be recovered, without imposing additional structure on the demand side of the market. If a researcher is reluctant to place structure on demand, an alternative is to focus instead on the difference in mean markups between groups of firms for which one is comfortable assuming that the production function parameters are the same across the groups. In Appendix B.2 we show that this difference can be estimated consistently without knowledge of the output elasticity, using a regression specification for the cost share in revenue for a flexible input. A leading example would be the comparison of mean markups across exporters and non-exporters, considered in De Loecker and Warzynski (2012), provided one is willing to assume that production function elasticities do not vary systematically with export status. However, this approach is not well suited to studying trends in markups, since the maintained assumption that the output elasticity is stable over time cannot

be verified without a way of estimating the output elasticity consistently for different sub-periods of the sample.

# References

**Ackerberg, Daniel A., Kevin Caves, and Garth Frazer**, "Identification Properties of Recent Production Function Estimators," *Econometrica*, 2015, *83* (6), 2411–2451.

**Amiti, Mary, Oleg Itskhoki, and Jozef Konings**, "International Shocks, Variable Markups, and Domestic Prices," *Review of Economic Studies*, 2019, *86* (6), 2356–2402.

**Atkeson, Andrew and Ariel Burnstein**, "Pricing-to-Market, Trade Costs, and International Relative Prices," *American Economic Review*, 2008, *98* (5), 1998–2031.

**Baqaee, David R. and Emmanuel Farhi**, "Productivity and Misallocation in General Equilibrium," *Quarterly Journal of Economics*, 2020, *135* (1), 105–163.

**Basu, Susanto**, "Are Price-Cost Markups Rising in the United States? A Discussion of the Evidence," *Journal of Economic Perspectives*, 2019, *33* (3), 3–22.

**Blundell, Richard and Steve Bond**, "GMM Estimation with Persistent Panel Data: An Application to Production Functions," *Econometric Reviews*, 2000, *19* (3), 321–340.

**Bond, Steve and Måns Söderbom**, "Adjustment Costs and the Identification of Cobb-Douglas Production Functions," Working Paper No. 05/04, Institute for Fiscal Studies, 2005.

**Burstein, Ariel, Vasco M. Carvalho, and Basile Grassi**, "Bottom-Up Markup Fluctuations," Working Paper No. 27958, National Bureau of Economic Research, 2020.

**De Loecker, Jan**, "Product Differentiation, Multiproduct Firms, and Estimating the Impact of Trade Liberalization on Productivity," *Econometrica*, 2011, *79* (5), 1407–1451.

_ **and Frederic Warzynski**, "Markups and Firm-Level Export Status," *American Economic Review*, 2012, *102* (6), 2437–71.

_ **and Pinelopi K. Goldberg**, "Firm Performance in a Global Market," *Annual Review of Economics*, 2014, *6* (1), 201–227.

_ **, Jan Eeckhout, and Gabriel Unger**, "The Rise of Market Power and the Macroeconomic Implications," *Quarterly Journal of Economics*, 2020, *135* (2), 561–644.

**Doraszelski, Ulrich and Jordi Jaumandreu**, "Using Cost Minimization to Estimate Markups," Discussion Paper No. DP14114, CEPR, 2019.

**Flynn, Zach, James Traina, and Amit Gandhi**, "Measuring Markups with Production Data," Working Paper, https://ssrn.com/abstract=3358472, 2019.

**Gandhi, Amit, Salvador Navarro, and David A. Rivers**, "On the Identification of Gross Output Production Functions," *Journal of Political Economy*, 2020, *128* (8), 2973–3016.

**Hall, Robert E.**, "Market Structure and Macroeconomic Fluctuations," *Brookings Papers on Economic Activity*, 1986, *17* (2), 285–338.

_ , "The Relation between Price and Marginal Cost in US Industry," *Journal of Political Economy*, 1988, *96* (5), 921–947.

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

**Mrázová, Monika, J. Peter Neary, and Mathieu Parenti**, "Sales and Markup Dispersion: Theory and Empirics," Working Paper No. 7433, CESifo, 2018.

**Pesaran, M. Hashem and Ron Smith**, "Estimating Long-Run Relationships from Dynamic Heterogeneous Panels," *Journal of Econometrics*, 1995, *68* (1), 79–113.

**Robinson, Peter M.**, "Root-N-Consistent Semiparametric Regression," *Econometrica*, 1988, *56* (4), 931–954.

**Syverson, Chad**, "What Determines Productivity?," *Journal of Economic Literature*, 2011, *49* (2), 326–65.

_ , "Macroeconomics and Market Power: Context, Implications, and Open Questions," *Journal of Economic Perspectives*, 2019, *33* (3), 23–43.

**Traina, James**, "Is Aggregate Market Power Increasing? Production Trends using Financial Statements," Working Paper No. 17, Stigler Center, 2018.

## Online Appendices

## Appendix A. Appendix for Section 2

*Appendix A.1. Oligopolistic competition*

The demand system, expressed in logarithms, facing an oligopolist $i$ in period $t$ is of the form

$$q_{it} = \mathcal{Q}_i(p_{it}, z_t) \tag{A.1}$$

where $q_{it}$ and $p_{it}$ denote the log output quantity and the log output price of firm $i$ in period $t$, respectively, and $z_t$ denotes the log industry expenditure function in period $t$ that depends on the prices $(p_{1t}, \ldots, p_{Jt})$ of all $N$ active firms operating in the industry in period $t$. Oligopolistic firms internalize the fact that they are nonatomistic in their industry and can influence industry expenditure $z_t$ through their decisions. The key assumption underlying the demand system (A.1) is that the industry expenditure function $z_t$ serves as a sufficient statistic for the prices of all competitors $\boldsymbol{p}_{-it} := \{p_{kt}\}_{k \neq i}$ of firm $i$, given firm $i$'s own price $p_{it}$.

The demand system (A.1) includes the popularized log-linear nested CES demand system of Atkeson and Burnstein (2008)

$$\mathcal{Q}_i(p_{it}, z_t) = (\rho - \eta) z_t - \rho p_{it}$$

where $\eta \geq 1$ denotes the constant elasticity of substitution of goods across industries and $\rho > \eta$ denotes the constant elasticity of substitution of goods within an industry. We do not restrict our subsequent analysis to this particular parametric functional form and instead work with the general demand system in equation (A.1).

Shephard's lemma (an envelope condition) identifies the partial change in $z_t$ with respect to $p_{it}$ as

$$\frac{\partial z_t}{\partial p_{it}} = S_{it}$$

where firm $i$'s market share of the total industry revenues is

$$S_{it} := \frac{P_{it} Q_{it}}{\displaystyle\sum_{k=1}^{J} P_{kt} Q_{kt}}$$

The own-price $\eta_{it}$ and cross-price $\delta_{ik,t}$ elasticities of demand are defined as

$$\eta_{it} := -\frac{dq_{it}}{dp_{it}} = -\left[ \frac{\partial q_{it}}{\partial p_{it}} + \frac{\partial q_{it}}{\partial z_t} S_{it} \right],$$

$$\delta_{ik,t} := \frac{dq_{it}}{dp_{kt}} = \frac{\partial q_{it}}{\partial z_t} S_{kt}, \forall k \neq i$$

1

Totally differentiating the demand system (A.1) and the accounting identity for revenue $r_{it} = p_{it} + q_{it}$, and invoking Shephard's lemma, yields an expression for the elasticity of revenue $r_{it}$ with respect to output $q_{it}$

$$\frac{dr_{it}}{dq_{it}} = \left( 1 + \left[ \frac{\partial q_{it}}{\partial p_{it}} + \frac{\partial q_{it}}{\partial z_t} \left( S_{it} + \sum_{k \neq i} S_{kt} \frac{dp_{kt}}{dp_{it}} \right) \right]^{-1} \right)$$

We emphasize that the elasticity $dr_{it}/dq_{it}$ is different under the Bertrand and Cournot models of oligopolistic competition. The reason is that the firm's conjectural elasticities, $dp_{kt}/dp_{it}, \forall k \neq i$, differ between these two models of oligopoly. A useful benchmark is monopolistic competition, under which there are no strategic considerations, and therefore we obtain the simplification $dp_{kt}/dp_{it} = 0, \forall k \neq i$. This yields the familiar result $dr_{it}/dq_{it} = \frac{\eta_{it}-1}{\eta_{it}}$.

Recall that the revenue elasticity of the flexible input $X_{it}^j$ is

$$\theta_{it}^{R,j} = \frac{dr_{it}}{dq_{it}} \theta_{it}^{Q,j}$$

Then, the estimand of the ratio estimator using the revenue elasticity in the numerator is

$$\begin{aligned}
\mu_{it}^R &= \frac{\theta_{it}^{R,j}}{\alpha_{it}^j} \\
&= \frac{dr_{it}}{dq_{it}} \frac{\theta_{it}^{Q,j}}{\alpha_{it}^j} \\
&= \frac{dr_{it}}{dq_{it}} \mu_{it}
\end{aligned}$$

Amiti et al. (2019) show that the firm's first order condition in the static profit maximization problem uniquely characterizes the firm's markup $\mu_{it}$ as a function of the firm's perceived price elasticity of demand $\sigma_{it}$. That is,

$$\mu_{it} = \frac{\sigma_{it}}{\sigma_{it} - 1}$$

The perceived demand elasticity $\sigma_{it}$ differs under Bertrand and Cournot competition. We now consider each in turn.

The Bertrand-Nash equilibrium condition is that all competitors of firm $i$ hold their prices fixed, i.e. $dp_k = 0, \forall k \neq i$. Then, the elasticity of revenue with respect to output simplifies to

$$\frac{dr_{it}}{dq_{it}} = \frac{\eta_{it} - 1}{\eta_{it}}$$

The perceived demand elasticity under Bertrand competition is equal to the own-price

2

demand elasticity.

$$\sigma_{it} \ = \ \eta_{it}$$

Then, the Bertrand markup $\mu_{it}^{Bertrand}$ is

$$
\begin{aligned}
\mu_{it}^{Bertrand} \ &= \ \frac{\eta_{it}}{\eta_{it} - 1} \\
&= \ \left[ \frac{dr_{it}}{dq_{it}} \right]^{-1}
\end{aligned}
$$

It follows that the estimand of the ratio estimator using the revenue elasticity does not identify the markup:

$$
\begin{aligned}
\mu_{it}^{R,Bertrand} \ &= \ \frac{dr_{it}}{dq_{it}} \mu_{it}^{Bertrand} \\
&= \ \frac{dr_{it}}{dq_{it}} \left[ \frac{dr_{it}}{dq_{it}} \right]^{-1} \\
&= \ 1
\end{aligned}
$$

The Cournot-Nash equilibrium condition is that competitors hold their quantities fixed, i.e. $dq_k = 0, \forall k \neq i$. Then, the elasticity of revenue with respect to output simplifies to

$$
\frac{dr_{it}}{dq_{it}} \ = \ \left( 1 + \left[ \frac{\partial q_{it}}{\partial p_{it}} + \frac{\partial q_{it}}{\partial z_t} \left( \frac{S_{it}}{1 - \tilde{S}_{-i}} \right) \right]^{-1} \right)
$$

where the response of competitiors is summarized in the statistic

$$
\tilde{S}_{-i} \ := \ - \sum_{k \neq i} \left( \frac{\partial q_{kt}}{\partial p_{kt}} \right)^{-1} \left( \frac{\partial q_{kt}}{\partial z_t} \right) S_{kt}
$$

The perceived demand elasticity under Cournot competition is equal to

$$
\sigma_{it} \ = \ - \left[ \frac{\partial q_{it}}{\partial p_{it}} + \frac{\partial q_{it}}{\partial z_t} \left( \frac{S_{it}}{1 - \tilde{S}_{-i}} \right) \right]
$$

Then, the Cournot markup $\mu_{it}^{Cournot}$ is

$$
\begin{aligned}
\mu_{it}^{Cournot} \ &= \ \frac{\sigma_{it}}{\sigma_{it} - 1} \\
&= \ \left( 1 + \left[ \frac{\partial q_{it}}{\partial p_{it}} + \frac{\partial q_{it}}{\partial z_t} \left( \frac{S_{it}}{1 - \tilde{S}_{-i}} \right) \right]^{-1} \right)^{-1} \\
&= \ \left( \frac{dr_{it}}{dq_{it}} \right)^{-1}
\end{aligned}
$$

3

Combining everything together, we again establish that the ratio estimator using the revenue elasticity does not identify the markup:

$$
\begin{aligned}
\mu_{it}^{R,Cournot} &= \frac{dr_{it}}{dq_{it}}\mu_{it}^{Cournot} \\
&= \frac{dr_{it}}{dq_{it}}\left(\frac{dr_{it}}{dq_{it}}\right)^{-1} \\
&= 1
\end{aligned}
$$

*Appendix A.2. Input adjustment costs*

We consider the same firm problem from Section 2, but we now assume that each input $j$ is associated with a baseline quantity $\overline{X}_{it}^{j}$ and that the firm incurs adjustment costs when it chooses an input quantity $X_{it}^{j} \neq \overline{X}_{it}^{j}$. The baseline quantity $\overline{X}_{it}^{j}$ might reflect the input choice from the previous period in a dynamic version of the model. For simplicity, we assume that these costs are given by the smooth convex function $\kappa^{j}\left(X_{it}^{j}\right)$, which satisfies $\kappa^{j}\left(\overline{X}_{it}^{j}\right) = \frac{d\kappa^{j}\left(\overline{X}_{it}^{j}\right)}{dX_{it}^{j}} = 0$.

The firm's cost function is now given by

$$
\mathcal{C}\left(Q_{it}; \boldsymbol{W}_{t}\right) := \min_{\left\{X_{it}^{j}\right\}_{j=1}^{J}} \left\{\sum_{j=1}^{J} W_{t}^{j} X_{it}^{j} + \sum_{j=1}^{J} \kappa\left(X_{it}^{j}\right) W_{t}^{j}\right\}
$$
$$
\text{s.t.} \quad \mathcal{F}\left(X_{it}^{1}, \ldots, X_{it}^{J}\right) \geq Q_{it},
$$

where we have normalized the adjustment cost functions by the input price $W_{t}^{j}$. Following the same steps as in the previous section, we obtain the FOC

$$
\frac{W_{t}^{j} X_{it}^{j}}{P_{it} Q_{it}}\left[1 + \frac{d\kappa^{j}\left(X_{it}^{j}\right)}{dX_{it}^{j}}\right] = \frac{\lambda_{it}}{P_{it}}\theta_{it}^{Q,j}.
$$

Using $\alpha_{it}^{j}$ to denote the share of input $j$'s cost in revenue and using the envelope condition, this implies

$$
\alpha_{it}^{j}\left[1 + \frac{d\kappa^{j}\left(X_{it}^{j}\right)}{dX_{it}^{j}}\right] = \frac{\partial \mathcal{C}\left(\cdot\right)}{\partial Q_{it}}\frac{\theta_{it}^{Q,j}}{P_{it}}. \tag{A.2}
$$

Hence, the ratio estimator using the revenue elasticity recovers

$$
\mu_{it}^{R,j} = \frac{\theta_{it}^{R,j}}{\alpha_{it}^{j}} = 1 + \frac{d\kappa^{j}\left(X_{it}^{j}\right)}{dX_{it}^{j}},
$$

4

and the ratio estimator using the output elasticity recovers

$$\mu_{it}^{Q,j} = \frac{\theta_{it}^{Q,j}}{\alpha_{it}^j} = \mu_{it} \left[ 1 + \frac{d\kappa^j \left( X_{it}^j \right)}{dX_{it}^j} \right].$$

Why might it be more common to estimate $\mu_{it}^{R,j} > 1$ than $\mu_{it}^{R,j} < 1$ when using firm-level data? One hypothesis is that adjustment costs are asymmetrical. It is less costly to use less of an input than previously planned than to use more of an input. If this is the case then on average we would recover $\mu_{it}^{R,j} > 1$. Similarly if firms are growing on average we would recover $\mu_{it}^{R,j} > 1$ on average.

The argument above effectively assumes that observed input costs are $W_t^j X_{it}^j$ rather than $W_t^j X_{it}^j + W_t^j \kappa^j \left( X_{it}^j \right)$. If this is the measure of observed input costs then

$$\alpha_{it}^j = \frac{W_i^j X_{it}^j + W_t^j \kappa^j \left( X_{it}^j \right)}{P_{it} Q_{it}}$$

and we obtain

$$\frac{W_t^j X_{it}^j + W_t^j \frac{d\kappa^j \left( X_{it}^j \right)}{dX_{it}^j}}{P_{it} Q_{it}} = \frac{\lambda_{it}}{P_{it}} \theta_{it}^{Q,j}$$

$$\mu_{it}^{Q,j} = \frac{\theta_{it}^{Q,j}}{\alpha_{it}^j} = \mu_{it} \left( \frac{X_{it}^j + \frac{d\kappa^j \left( X_{it}^j \right)}{dX_{it}^j}}{X_{it}^j + \kappa^j \left( X_{it}^j \right)} \right)$$

so wedge $> 1$ whenever $\kappa' > \kappa$.

Neither of the two cases that are typically considered in the literature lead to a bias. The flexible input case is $\kappa^j = 0$, in which case the bias disappears. The fixed input case is one in which $X_{it}^j \to \overline{X}_{it}^j$ in which case the bias also disappears. (Note, however that the fixed input case is not the limit as $\kappa^j \to \infty$, and so is not a special case of the model with adjustment cost model. When $\kappa^j \to \infty$ in the adjustment cost model, the bias remains even in the limit, even though $X_{it}^j \to \overline{X}_{it}^j$).

*Appendix A.3. Inputs that influence demand*

In this section we show that even if output elasticities are available, markup estimates are biased whenever the variable factor of production is used partly to affect demand in addition to producing output.

We assume that the firm's production function is as in Section 2, but that its revenue is now given by

$$R_{it} := \mathcal{P} \left( Q_{it}, D_{it} \right) Q_{it}$$

where $D_{it}$ is an endogenous demand shifter that the firm can influence through the use of

inputs according to the function

$$D_{it} = \mathcal{D}\left(X_{it}^{D,1},,\ldots,X_{it}^{D,J}\right).$$

We denote the amount of input $j$ used in production as $X_{it}^{Q,j}$ and the amount used in influencing demand as $X_{it}^{D,j}$. The total quantity of input $j$ used by the firm is $X_{it}^{j} = X_{it}^{Q,j} + X_{it}^{D,j}$.

The profit maximization problem of the firm is now

$$\Pi_{it} \quad := \quad \max_{Q_{it},D_{it}} \left\{ \mathcal{P}\left(Q_{it}, D_{it}\right) Q_{it} - \mathcal{C}_Q\left(Q_{it}; \boldsymbol{W}_t\right) - \mathcal{C}_D\left(D_{it}; \boldsymbol{W}_t\right) \right\} \tag{A.3}$$

where $\mathcal{C}_Q\left(Q_{it}; \boldsymbol{W}_t\right)$ is the firm's cost function for producing output, defined by

$$\mathcal{C}_Q\left(Q_{it}; \boldsymbol{W}_t\right) \quad := \quad \min_{\left\{X_{it}^{Q,j}\right\}_{j=1}^{J}} \left\{ \sum_{j=1}^{J} W_t^j X_{it}^{Q,j} \right\} \tag{A.4}$$

$$\text{s.t.} \quad Q_{it} \leq \mathcal{F}\left(X_{it}^{Q,1},\ldots,X_{it}^{Q,J}\right)$$

and $\mathcal{C}_D\left(D_{it}; \boldsymbol{W}_t\right)$ is the firm's cost function for influencing demand, defined by

$$\mathcal{C}_D\left(D_{it}; \boldsymbol{W}_t\right) \quad := \quad \min_{\left\{X_{it}^{D,j}\right\}_{j=1}^{J}} \left\{ \sum_{j=1}^{J} W_t^j X_{it}^{D,j} \right\} \tag{A.5}$$

$$\text{s.t.} \quad \mathcal{D}\left(X_{it}^{D,1},,\ldots,X_{it}^{D,J}\right) \geq D_{it}$$

The optimality conditions from the profit maximization problem (A.3) are

$$1 - \frac{1}{\eta_{it}} = \frac{\partial \mathcal{C}_Q\left(\cdot\right)}{\partial Q_{it}} \frac{1}{P_{it}} \tag{A.6}$$

$$\varsigma_{it} = \frac{\partial \mathcal{C}_D\left(\cdot\right)}{\partial D_{it}} \frac{D_{it}}{P_{it} Q_{it}} \tag{A.7}$$

where $\varsigma_{it}$ describes the effect of the demand shifter on the price that a firm can charge for a given quantity of output. As in the previous section, the optimal markup of price over marginal production cost is

$$\mu_{it} := \left[\frac{\partial \mathcal{C}_Q\left(\cdot\right)}{\partial Q_{it}} \frac{1}{P_{it}}\right]^{-1} = \left(1 - \frac{1}{\eta_{it}}\right)^{-1}.$$

The FOC for the production cost minimization problem (A.4) yields the relationship

$$\alpha_{it}^{Q,j} = \frac{\partial \mathcal{C}_Q\left(\cdot\right)}{\partial Q_{it}} \frac{1}{P_{it}} \theta_{it}^{Q,j} \tag{A.8}$$

where $\alpha_{it}^{Q,j}$ is the share of revenue paid to input m for use in producing output, and $\theta_{it}^{Q,j}$ is the elasticity of output to the use of input $j$ for production. It follows from equation (A.8) that if one could observe $X_{it}^{Q,j}$ separately from $X_{it}^j$ then the ratio estimator would correctly recover the markup.

However, in practice we observe only the total usage of an input $X_{it}^j = X_{it}^{Q,j} + X_{it}^{D,j}$, rather then the usage in different activities separately. Using the FOC for the cost minimization problem for influencing demand (A.5) yields the relationship

$$\alpha_{it}^{D,j} = \frac{\partial \mathcal{C}_D\left(\cdot\right)}{\partial D_{it}} \frac{D_{it}}{P_{it}Q_{it}} \theta_{it}^{D,j} \tag{A.9}$$

where $\alpha_{it}^{D,j}$ is the share of revenue paid to input $j$ for shifting demand and $\theta_{it}^{D,j}$ is the elasticity of $D_{it}$ with respect to $X_{it}^{D,j}$. Combining (A.6),(A.7), (A.8) and (A.9) yields an expression for the total revenue share of input $X_{it}^j$

$$\alpha_{it}^j = \left(1 - \frac{1}{\eta_{it}}\right)\theta_{it}^{Q,j} + \varsigma_{it}\theta_{it}^{D,j} \tag{A.10}$$

To see what the ratio estimator recovers, note that the optimality condition for allocating an input $j$ between producing goods $X_{it}^{Q,j}$ and influencing demand $X_{it}^{D,j}$ implies that the output elasticity of an input $X_{it}^j$ is

$$\theta_{it}^{Q,j}\rho_{it}^{Q,j} + \frac{\partial \mathcal{F}}{\partial X_{it}^{D,j}} \frac{X_{it}^{D,j}}{Q_{it}} \psi_{it}^{D,j} = \theta_{it}^{Q,j}\rho_{it}^{Q,j} \tag{A.11}$$

where $\psi_{it}^{Q,j}$ is the elasticity of $X_{it}^{Q,j}$ with respect to $X_{it}^j$ evaluated at the optimum. $\psi_{it}^{D,j}$ denotes the elasticity of $X_{it}^{D,j}$ with respect to $X_{it}^j$ evaluated at the optimum. This means that in order to correctly recover the output elasticity of an input $X_{it}^j$, it is necessary to separately observe the part of that input that is actually used in producing goods as long as $\psi_{it}^{Q,j} \neq 1$. The fact that a firm uses inputs partly to influence demand introduces a bias into the estimate of the output elasticity. It also introduces a bias into the estimate of the markup. Combining (A.10) and (A.11) reveals that the estimand is given by

$$\mu_{it}^{Q,j} = \mu_{it} \left[\frac{\psi_{it}^{Q,j}}{1 + \frac{X_{it}^{D,j}}{X_{it}^{Q,j}}}\right].$$

There are however special cases in which $\psi_{it}^{Q,j} = 1$, i.e. the share of $X_{it}^j$ in production and in influencing demand does not depend on the level of $X_{it}^j$. For example it is sufficient that the firm faces an isoelastic demand curve and $\mathcal{F}$ and $\mathcal{D}$ are Cobb-Douglas. If this is the case, there is no bias the estimate of the output elasticity, but the ratio estimator is

still biased. [1]

$$\mu_{it}^{Q,j} = \mu_{it} \left[ \frac{1}{1 + \frac{X_{it}^{D,j}}{X_{it}^{Q,j}}} \right].$$

So if the flexible input is only used for production and not to influence demand ($X_{it}^{D,j} = 0$) then the ratio estimator recovers the markup. But if some of the input is used to influence demand, and this component is not separated out, then the ratio estimand will be biased. If, over time, the input $X_{it}^{j}$ is increasingly being used to influence demand, then the ratio estimand will fall over time, without any change in the true markup.

Casual observation suggests that at least some part of the workforce currently employed in the corporate sector devotes its energy to influencing demand rather than to producing goods. This suggests that using labor as an input for estimating markups will yield estimates that are hard to interpret. When using the ratio estimator, heterogeneity across firms and industries in the extent to which they use labor for production versus marketing and sales-related expenses will thus manifest as heterogeneity in measured markups.

These observations also help shed light on the difference in the trend in markups that one obtains from Compustat data on US firms when one uses only COGS versus when one includes SGA as the flexible input (De Loecker et al., 2020; Traina, 2018). It seems reasonable to assume that in the COGS bundle, a larger fraction of the inputs is used to produce output and a smaller fraction is used to influence demand, than in the SGA bundle. Thus the downward bias in the ratio estimand is likely to be larger when including SGA in the bundle of flexible inputs, versus when using only COGS. Since the cost share of SGA in total revenue has been increasing relative to the cost share of COGS in total revenue, this will manifest as a widening gap between the ratio estimator that uses only COGS and the ratio estimator that also includes SGA. This is precisely what the literature has found.

So far in this section we have proceeded as if output were observed. If only revenue were observed, as in Section 2.1, then the ratio estimator again recovers $\mu_{it}^{R,j} = 1$, regardless of whether the input is being used for production or to influence demand. Given that Compustat data contains only revenue, not output, the aforementioned discussion is relevant only if one believes that the procedures in those papers do successfully recover output elasticities, which we believe they do not.

---

[1]This result does not require that $X_{it}^{D,j}$ and $X_{it}^{Q,j}$ are perfect substitutes, but it does require that they satisfy $X_{it} = h\left(X_{it}^{Q,j}, X_{it}^{D,j}\right)$ where $h$ is a constant-returns-to-scale function. Thanks to Agustin Gutierrez for pointing this out.

## Appendix B. Appendix for Section 3

### Appendix B.1. Optimal input demand functions

This appendix supplies the derivations of the optimal input demand equation for intermediate inputs under two production technology specifications. Section Appendix B.1.1 provides the derivation for a Cobb-Douglas technology, while Section Appendix B.1.2 provides that for a non-parametric technology.

### Appendix B.1.1. Cobb-Douglas technology

The three-factor Cobb-Douglas production function for gross output $Q_{it}$ with Hicks-neutral productivity $\omega_{it}$ is

$$Q_{it} = \exp\left(\omega_{it}\right) K_{it}^{\beta_K} L_{it}^{\beta_L} M_{it}^{\beta_M}$$

Since $M_{it}$ is the single flexible input, the cost minimizing input demand for $M_{it}$ can be obtained by rearranging the Cobb-Douglas production function conditional on a given output quantity $Q_{it}$

$$M_{it} = M^*\left(Q_{it}; K_{it}, L_{it}, \omega_{it}\right) := \exp\left(-\frac{1}{\beta_M}\omega_{it}\right) Q_{it}^{\frac{1}{\beta_M}} K_{it}^{-\frac{\beta_K}{\beta_M}} L_{it}^{-\frac{\beta_L}{\beta_M}} \tag{B.1}$$

Then, the minimized total variable cost function is

$$\mathcal{C}\left(Q_{it}; K_{it}, L_{it}, P_{it}^M, \omega_{it}\right) := P_{it}^M M^*\left(Q_{it}; K_{it}, L_{it}, \omega_{it}\right) \tag{B.2}$$

where $P_{it}^M$ is the unit input price of $M_{it}$ that firm $i$ takes as given. Taking the demand system $P_{it} = \mathcal{P}\left(Q_{it}\right)$ and the total cost function $\mathcal{C}\left(Q_{it}; K_{it}, L_{it}, P_{it}^M, \omega_{it}\right)$ as given, firm $i$ chooses $Q_{it}$ to solve a static profit maximization problem

$$\max_{Q_{it}} \left\{\mathcal{P}\left(Q_{it}\right) Q_{it} - \mathcal{C}\left(Q_{it}; K_{it}, L_{it}, P_{it}^M, \omega_{it}\right)\right\}$$

The first order condition in profit maximization equates marginal revenue to marginal cost

$$\mathcal{P}\left(Q_{it}\right)\left(\frac{\eta_{it} - 1}{\eta_{it}}\right) = \frac{\partial\mathcal{C}\left(Q_{it}; K_{it}, L_{it}, P_{it}^M, \omega_{it}\right)}{\partial Q_{it}} \tag{B.3}$$

where $\eta_{it}$ is the absolute value of the price elasticity of demand. Equation (B.3) identifies the markup $\mu_{it}$ under monopolistic competition as a function of the demand elasticity.

$$\mu_{it} = \frac{\eta_{it}}{\eta_{it} - 1}$$

Applying the functional form in equation (B.1) to the FOC in equation (B.3) and solving for $q_{it} := \ln Q_{it}$ gives

$$q_{it} = \frac{\beta_M}{1 - \beta_M}\ln\beta_M + \frac{\beta_k}{1 - \beta_M}k_{it} + \frac{\beta_l}{1 - \beta_M}l_{it} + \frac{\beta_M}{1 - \beta_M}\left(p_{it} - \ln\mu_{it} - p_{it}^M\right) + \frac{1}{1 - \beta_M}\omega_{it} \tag{B.4}$$

where $p_{it}^M := \ln P_{it}^M$ and $p_{it} := \ln P_{it}$. Using equation (B.4) to substitute for $q_{it}$ in equation (B.1) produces the optimal input demand equation for $m_{it}$ in terms of the state variables $(k_{it}, l_{it}, \omega_{it})$, the exogenous input price $p_{it}^M$, and the endogenous optimal output price $p_{it}$ and markup $\mu_{it}$.

$$m_{it} = \frac{\ln\beta_M}{1 - \beta_M} + \frac{\beta_K}{1 - \beta_M} k_{it} + \frac{\beta_L}{1 - \beta_M} l_{it} + \frac{1}{1 - \beta_M} \left( p_{it} - \ln\mu_{it} - p_{it}^M + \omega_{it} \right)$$

*Appendix B.1.2. Non-parametric technology*

The non-parametric three-factor production function for gross output with productivity $\omega_{it}$ is

$$Q_{it} = \mathcal{F}\left(K_{it}, L_{it}, M_{it}, \omega_{it}\right) \tag{B.5}$$

The only restriction we impose on the function $\mathcal{F}\left(\cdot\right)$ is that it is twice continuously differentiable in each of its arguments. As in Section Appendix B.1.1, $M_{it}$ is the single flexible input. Inverting equation (B.5) produces the cost-minimizing input demand for $M_{it}$.

$$M_{it}^* = \mathcal{F}^{-1}\left(Q_{it}; K_{it}, L_{it}, \omega_{it}\right) \tag{B.6}$$

The minimized total variable cost function is thus

$$\mathcal{C}\left(Q_{it}; K_{it}, L_{it}, P_{it}^M, \omega_{it}\right) := P_{it}^M \mathcal{F}^{-1}\left(Q_{it}; K_{it}, L_{it}, \omega_{it}\right)$$

The first order condition in profit maximization is then

$$\mathcal{P}\left(Q_{it}\right)\left(\frac{\eta_{it} - 1}{\eta_{it}}\right) = P_{it}^M \frac{\partial \mathcal{F}^{-1}\left(Q_{it}; K_{it}, L_{it}, \omega_{it}\right)}{\partial Q_{it}} \tag{B.7}$$

Given a functional form for $\mathcal{F}\left(\cdot\right)$, one can solve equation (B.7) for the optimal output level $Q_{it}$.

$$Q_{it} = Q^*\left(K_{it}, L_{it}, P_{it}^M, \omega_{it}, P_{it}, \mu_{it}\right) \tag{B.8}$$

Using equation (B.8) to substitute for $Q_{it}^*$ in equation (B.6) yields the optimal input demand function for intermediate inputs.

$$\begin{aligned} M_{it} &= \mathcal{F}^{-1}\left(Q^*\left(K_{it}, L_{it}, P_{it}^M, \omega_{it}, P_{it}, \mu_{it}\right); K_{it}, L_{it}, \omega_{it}\right) \\ &:= M^*\left(K_{it}, L_{it}, P_{it}^M, \omega_{it}, P_{it}, \mu_{it}\right) \end{aligned}$$

In the absence of price data on inputs and outputs, the scalar unobservables in the input demand function $M^*\left(\cdot\right)$ are $\left(P_{it}^M, \omega_{it}, P_{it}, \mu_{it}\right)$.

*Appendix B.2. Learning about variation in markups from variation in the cost share only*

Without a way to estimate the output elasticity for a flexible input consistently from typical production data, we cannot use the ratio estimator to learn about the level of markups. We can however still use insights from the production approach to learn about variation in markups across firms. This variation can be studied using a regression model

for the log of the cost share in total revenue for a perfectly flexible input. We sketch this 'cost share approach' to studying markups in this appendix.

As discussed in Section 2, the ratio estimator relies on the relationship $\mu_{it} = \frac{\theta_{it}^{Q,j}}{\alpha_{it}^{j}}$ for a flexible input $X_{it}^{j}$. Taking logs and rearranging, we obviously have $-\ln \alpha_{it}^{j} = -\ln \theta_{it}^{Q,j} + \ln \mu_{it}$. First consider the three factor, Cobb-Douglas case in which intermediate inputs $(M_{it})$ is the perfectly flexible input, as discussed in Section 3. Here $\ln \alpha_{it}^{M} = (p_{it}^{M} + m_{it}) - (p_{it} + q_{it})$ is the log of the true cost share in revenue for intermediate inputs, and $\ln \theta_{it}^{M} = \ln \beta_{M}$ is a constant term. Letting $\ln s_{it}^{M} = (p_{it}^{M} + m_{it}) - (p_{it} + y_{it})$ denote the log of the observed cost share in revenue for firm $i$ in period $t$, we then have

$$-\ln s_{it}^{M} = -\ln \beta_{M} + \ln \mu_{it} + \varepsilon_{it} \tag{B.9}$$

where $y_{it} = q_{it} + \varepsilon_{it}$ as before.[2]

Without a consistent estimate of the output elasticity $(\beta_{M})$, it is clear that the mean of the log of the observed cost shares conflates the log of the output elasticity and the mean of the log of the markups, and does not separately identify the latter. Nevertheless, under the maintained assumption that the output elasticity is common to all the firm-year observations, we can use this relation to study variation in markups. For example, if the binary dummy $D_{it}^{X}$ indicates whether or not firm $i$ in period $t$ is an exporter, we can specify a linear relationship between log markups and export status

$$\ln \mu_{it} = \delta_{0} + \delta_{1} D_{it}^{X} + \nu_{it} \tag{B.10}$$

as in De Loecker and Warzynski (2012). Substituting (B.10) into (B.9), we have the linear specification

$$-\ln s_{it}^{M} = (\delta_{0} - \ln \beta_{M}) + \delta_{1} D_{it}^{X} + (\varepsilon_{it} + \nu_{it}) \tag{B.11}$$

In the Cobb-Douglas case, we can thus learn about the *association* between log markups and export status from a simple regression of the log of the observed cost share in revenue for a flexible input on a constant and the export status dummy.[3]

For more general Hicks-neutral gross output production functions, we can write the

---

[2]For simplicity, we assume here that this is the only source of measurement error in the log of the observed cost share in revenue. In the Cobb-Douglas case, we can easily allow for (multiplicative) measurement error in both the numerator and the denominator of the cost share for intermediate inputs.

[3]As in De Loecker and Warzynski (2012), additional controls can be included in this regression specification, but OLS is still unlikely to estimate the causal effect of exporting on markups consistently. If the sample used to estimate (B.11) pools data for firms in several sectors, sector dummies can be used to allow for heterogeneity in the output elasticity $\beta_{M}$ between sectors.

log of the output elasticity $\ln \theta_{it}^M = f(k_{it}, l_{it}, m_{it})$,[4] in which case (B.11) becomes

$$- \ln s_{it}^M = g(k_{it}, l_{it}, m_{it}) + \delta_1 D_{it}^X + (\varepsilon_{it} + \nu_{it}) \tag{B.12}$$

where $g(k_{it}, l_{it}, m_{it}) = \delta_0 - f(k_{it}, l_{it}, m_{it})$. We can then learn about the association between log markups and export status either by approximating $g(k_{it}, l_{it}, m_{it})$ using a flexible functional form, or by estimating (B.12) using semi-parametric methods for partially linear models (Robinson, 1988).

This cost share approach allows us to learn about some forms of variation across firms in markups under essentially the same assumptions needed for the production approach, but without requiring a consistent estimate of the output elasticity. Except in the Cobb-Douglas case, we could not use this approach to study the association between markups and measures of firm size (e.g. the log of employment, $l_{it}$) or measures of factor intensity (e.g. the log of the capital-labor ratio, $k_{it} - l_{it}$); we may also have low power to detect significant association between markups and observed firm characteristics that are strongly correlated with functions of the production inputs. In principle, this approach could also be used to study trends in markups over time, as in De Loecker et al. (2020). However, it should be emphasized that the trend in the log of the cost share in revenue for a flexible input identifies the trend in the log of the markup only under the maintained assumption that the output elasticity is stable over time, which cannot be verified without a way of estimating the output elasticity consistently for different sub-periods.

---

[4]For example, in the translog case, we have $f(k_{it}, l_{it}, m_{it}) = \ln(\beta_M + \beta_{KM} k_{it} + \beta_{LM} l_{it} + \beta_{MM} m_{it})$.