

WP/14/236

IMF Working Paper

The Size Distribution of Manufacturing Plants and Development

Siddharth Kothari

IMF Working Paper

Research Department

The Size Distribution of Manufacturing Plants and Development¹

Prepared by Siddharth Kothari

Authorized for distribution by Andrew Berg

December 2014

This Working Paper should not be reported as representing the views of the IMF.

The views expressed in this Working Paper are those of the author(s) and do not necessarily represent those of the IMF or IMF policy. Working Papers describe research in progress by the author(s) and are published to elicit comments and to further debate.

Abstract

The typical size distribution of manufacturing plants in developing countries has a thick left tail compared to developed countries. The same holds across Indian states, with richer states having a much smaller share of their manufacturing employment in small plants. In this paper, I explore the hypothesis that this income-size relation arises from the fact that low income countries and states have high demand for low quality products which can be produced efficiently in small plants. I provide evidence which is consistent with this hypothesis from both the consumer and producer side. In particular, I show empirically that richer households buy higher price goods while larger plants produce higher price products (and use higher price inputs). I develop a model which matches these cross-sectional facts. The model features non-homothetic preferences with respect to quality on the consumer side. On the producer side, high quality production has higher marginal costs and requires higher fixed costs. These two features imply that high quality producers are larger on average and charge higher prices. The model can explain about forty percent of the cross-state variation in the left tail of manufacturing plants in India.

JEL Classification Numbers: O11, O17, E26, O53

Keywords: India, size distribution, manufacturing, non-homothetic preferences, quality, informal sector

Author's E-Mail Address: skothari@imf.org

* I am especially grateful to my advisor Pete Klenow for his advice on this project. I would also like to thank Manuel Amador, Nicholas Bloom, Pascaline Dupas, Robert Hall, Chad Jones, Pablo Kurlat, Kalina Manova, Monika Piazzesi, Martin Schneider, and Christopher Tonetti for helpful comments. I gratefully acknowledge support from the Leonard W. Ely and Shirley R. Ely Graduate Student Fund Fellowship (SIEPR), B.F. Haley and E.S. Shaw Fellowship for Economics (SIEPR), and the SEED Fellowship from the Stanford Institute for Innovation in Developing Economies.

	Contents	Page
1.	Introduction	3
2.	Empirical Results	8
	2.1. Richer Households Buy Higher Price Goods.....	9
	2.2. Larger Plants Produce Higher Price Goods	12
	2.3. Larger Plants Use Higher Price Inputs.....	15
3.	Model.....	17
	3.1. Households.....	17
	3.2. Final Goods Producers	21
	3.3. Intermediate Goods Producers	23
	3.4. Equilibrium	24
4.	Calibration	25
	4.1. Production Parameters.....	25
	4.2. Utility Parameters.....	28
5.	Results	30
	5.1. Cross-section of Indian States	30
	5.2. India Over Time	33
	5.3. Parameter Sensitivity: Love of Variety	35
	5.4. Indian vs US	37
6.	Inter-State Trade.....	37
7.	Conclusion.....	41
	References.....	42
A.	Appendix.....	45
	A.1. Annual Survey of Industries	45
	A.2. Survey of Unorganized Manufacturing.....	46
	A.3. Consumer Expenditure Surveys.....	48
	A.4. Employment-Unemployment Survey	49
	A.5. County Business Patterns Database (US)	50
B.	Inter-State Trade: Concordances.....	50
	B.1. NAICS 2002 to NIC 2004 Concordance for Herfindahl Index.....	50
	B.2. HS Product Classification to NIC 2004 Concordance for Export-Import Index	51
	B.3. Concordances Across Different NIC Revisions.....	51
C.	Units Misreporting Problem in the ASI	51
D.	Calibrating Production Parameters - θ_{q_n}	53

1. Introduction

The typical size distribution of manufacturing establishments in developing countries has a thick left tail compared to developed countries. Figure 1 plots the share of total workers in establishments of different size categories for India and the US for 2005-06. While about 60 percent of the workers are employed by establishments of size less than five in India, the corresponding number for the US is less than 2 percent.¹

This size-income relation also holds across Indian states. Figure 2 plots the share of employment in establishments of size five or less in 2005-06 for different Indian states against the per-capita Net Domestic Product (NDP) of the state relative to the poorest state (Bihar).² The richest Indian states have about four times the per-capita NDP of the poorest states. While the poorest states have almost 90 percent of their manufacturing workforce employed in establishments of size five or less, the richer states have only about 40 percent of their workforce working in small establishments.³

What explains this negative correlation between income levels and the share of employment in small establishments? Starting from the work of De Soto (1989), the previous literature has focused on size-dependent policies (regulatory burden faced by large firms, small scale reservation policies, etc) as an explanation for the size-income relation. These policies create distortions which can lead to misallocation of resources, lower income levels, and smaller establishment sizes.

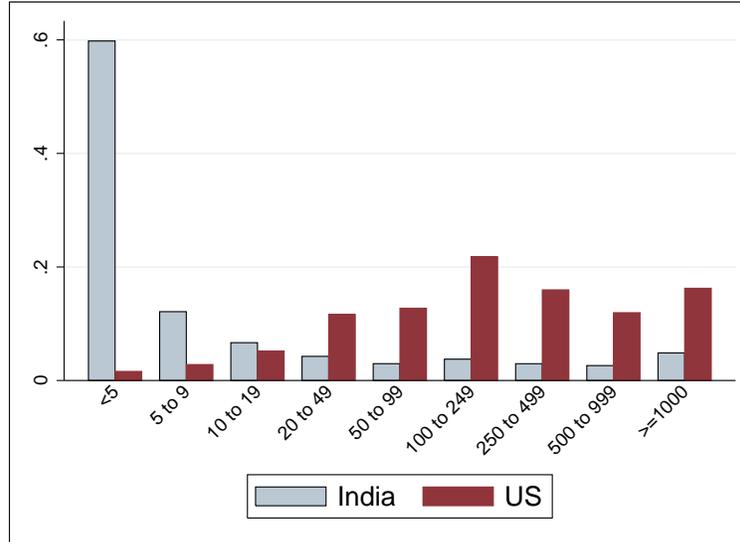
This paper explores an alternative (though potentially complementary) explanation for this size-income relation which is driven by preferences and technology rather than distortions. The hypothesis is that poor households have high demand for low quality products, which can be produced efficiently in small establishments as they require small fixed investments (no research and development expenditure, or no need for large investments in fixed capital). On the other hand, richer households tend to demand higher quality goods, whose production requires a larger scale due to the need for larger fixed investments. This relation between income levels and demand for quality implies that poor countries or states have demand skewed

¹The US data is taken from the US Business County Patterns Database maintained by the US Census Bureau. The Indian data combines two surveys, the Annual Survey of Industries (ASI) and the Survey of Unorganized Manufacturing (SUM). Appendix Sections A.1, A.2, and A.5 give more details regarding these datasets.

²There are a total of 28 states in India. Figure 2 plots only the 15 largest states in order to keep the graph readable. These 15 states cover 96.5 percent of the manufacturing workforce. The negative relation between share of employment in small plants and per-capita state NDP is robust to including all the states. In a regression of share of employment in plants of size five or less on log of per-capita state NDP, the coefficient (standard error) on log state NDP is -0.320 (0.0553) when restricting to 15 states and -0.319 (0.0568) when including all the states. A possible concern with the relation seen in Figure 2 is that it might be driven by differences in industry composition across states. However, a large part of the differences in share of employment in small plants across states is actually driven by within industry differences in size. Controlling for industry composition (weighting the size distribution in every state by the all India industry composition instead of the state specific industry composition) at the 2-digit level causes the slope coefficient on log of state per-capita NDP to fall from -0.320 (0.0553) to -0.274 (0.0417).

³The differences in share of employment in small plants also reflects in differences in average plant size across states. The average plant size in the richest states is about two times the average in the poorest states.

Figure 1: Share of Employment by Size Category: India vs. US



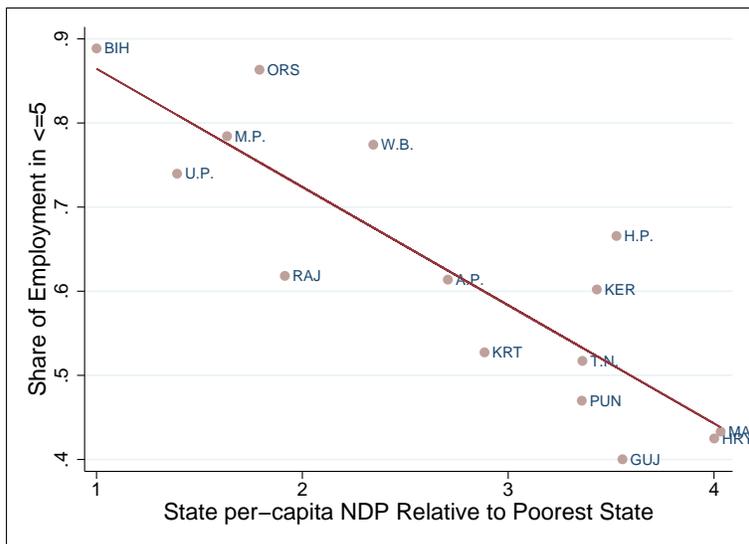
Notes: The graph plots the share of total employment in establishments of different size categories for India and the US. The data for India combines two sources, the Annual Survey of Industries (ASI) and the Survey of Unorganized Manufacturing (SUM) for 2005-06. The data for the US is taken from the County Business Patterns Database for 2006.

towards goods which require a small scale of production, which in turn causes the size distribution to be dominated by small plants. As a region develops and income levels increase, demand shifts towards high quality products, which in turn leads to a shift on the production side towards higher quality goods. This shift in production causes the share of employment in small plants to decrease, and thus can generate the negative relation between the share of employment in small plants and income levels seen in the data.

I provide empirical evidence in support of this hypothesis using Indian data from consumer and producer surveys:

1. Using data from Consumer Expenditure Surveys, I show that in the cross-section richer households tend to pay a higher unit price for the same good, which is consistent with the hypothesis that richer households buy higher quality products.
2. On the producer side, I show that larger plants tend to charge a higher unit price for the same good as compared to smaller plants, which is consistent with larger plants producing higher quality products. To show this, I combine data from the Annual Survey of Industries (ASI), which covers plants employing ten or more workers (twenty or more workers if not using power), and the Survey of Unorganized Manufacturing (SUM), which covers plants employing less than ten workers. The positive relation between prices and plant size holds not just within the formal sector (ASI plants), but also

Figure 2: Size Distribution of Manufacturing Establishments: Across Indian States



Notes: The graph plots the share of employment in plants of size five or less in a state against per-capita NDP of the state relative to the poorest state. The data for the states combines two sources, the Annual Survey of Industries (ASI) and the Survey of Unorganized Manufacturing (SUM). Only the 15 largest states are included to keep the graph readable.

when pooling together the formal and informal plants.

- Using ASI and SUM data, I show that larger plants use higher price material inputs, consistent with them using higher quality inputs. Using data from Household Surveys, I also find that larger plants hire more skilled workers.

I develop a general equilibrium model which matches these cross-sectional facts. Households choose from a finite number of quality levels. The choice over quality levels is modeled as a discrete-choice problem with households choosing to consume one quality level out of those available in the economy. Their preferences exhibit non-homotheticity with respect to quality: richer households are more likely to choose higher quality levels. The non-homotheticity arises because the utility function features complementarity between quality and quantity consumed (the marginal increase in utility from a given increase in quantity consumed is larger for higher quality goods) and richer households can consume more quantity of whichever quality level they choose.

On the producer side, production of high quality goods uses skilled labor more intensively. Also, starting a higher quality plant requires higher fixed costs, which combined with a free entry condition implies that producers of high quality goods will be larger on average (in order to recover their larger fixed costs).

The model parameters are chosen to match the micro-facts documented on the consumer and producer

side. The quality-size relation on the producer side is matched to the relation between prices and plant size from the producer surveys, while the degree of non-homotheticity is chosen to match the price-income relation seen in the consumer surveys.

I then ask the question: How much of the cross-state variation in the size distribution seen in Figure 2 can be explained by the model? In particular I conduct counterfactual exercises in which I simulate changes in per-capita income levels in the model (by varying productivity and the skill level of the population) and see what is the effect on the size distribution. As income levels increase in the model, demand shifts to high quality goods due to the non-homotheticity of preferences. This shift in demand towards higher quality leads to a shift on the production side, with a fall in the number of low quality producers and an increase in the number of high quality producers. As high quality producers are larger on average compared to low quality producers, there is also a shift in the size distribution towards larger plants. I find that the share of employment in plants of size five or less goes down by 19.3 percentage points (which is about 43 percent of the difference seen across Indian states) when income in the model varies by the same extent as it does across Indian states. I also document that the share of employment in plants of size five or less has gone down by about 20 percentage points in India between 1989 and 2009, and show that the model can explain about 65 percent of this change. While most of the results presented in the paper focus on the share of employment in plants which employ five or less people, Section 5.1 also explores the implications of the model on the entire size distribution.

The model and the counterfactual exercises make the implicit assumption that each state can be treated as a closed economy in which local demand is met by local production. How would the possibility of inter-state trade affect the hypothesis presented in the paper? A potential confounding effect of inter-state trade could come through the location choice of large plants. For example, if the richer states are more suited for operating large plants (due to availability of skilled labor, less stringent labor laws etc), then larger plants might choose to locate in these states (and ship their goods to the poor states) and this might be driving the negative relation between income and size that we see in Figure 2. If inter-state trade was an important force, then we would expect the more tradable industries within manufacturing to have a stronger negative relation between size and income levels across states. To test this, I construct two measures of tradability at the 3-digit level of industrial classification. I find that the size-income relation across states is not stronger for tradables as compared to non-tradables (for one of the measures, the non-tradables actually have a stronger negative relation as compared to tradables) indicating that inter-state trade is unlikely to be an important force behind the relation seen in Figure 2. I discuss the issue of inter-state trade in more detail in Section 6.

This paper is related to several strands of literature. A large literature has studied the question of why

the size distribution differs markedly across countries. The role of distortionary policies and the regulatory environment in determining the size distribution of plants (and the extent of informality) has been studied in Little, Mazumdar, and Page Jr (1987), De Soto (1989), Loayza (1996), Djankov and others (2002), Loayza, Oviedo, and Serven (2005), Loayza, Serven, and Sugawara (2009), Garicano, LeLarge, and Van Reenen (2013) among others. While size-dependent policies are potentially an important determinant of the size distribution, these policies are unlikely to explain all the differences in size distribution seen between developing and developed countries. Tybout (2000) notes that all developing countries tend to have a large share of their population in small plants, irrespective of whether they have policies which discriminate against large plants or not. This suggests that these policies cannot be the only factor driving plant size. Gollin (1995) and Hsieh and Klenow (2012) conduct quantitative exercises in which they find that size-dependent policies leave a large part of the differences in size across countries unexplained. Hsieh and Olken (2014) document that the “missing middle” in the size distribution in developing countries actually does not exist and that regulatory obstacles which become binding at particular threshold levels do not seem to lead to discontinuities in the size distribution in developing countries.⁴ This paper suggests that a large part of the differences in size distribution that we see across countries and states is a natural consequence of the low levels of income in developing countries and is not necessarily caused by policies which discriminate against large productive plants in favor of small unproductive plants. The hypothesis considered in the paper is closer to the dual-sector view of the informal sector in La Porta and Shleifer (2008) according to which the informal sector does not compete directly with the formal sector. Also related is the idea in Banerjee and Duflo (2011) which considers the informal economy to be employing poor individuals and using a different production technology characterized by small fixed costs. I focus on the heterogeneity of quality levels being produced by plants of different sizes and how the demand for low quality falls with development.⁵

Some of the empirical results documented here have been studied in different contexts (or for different countries) in other papers. Deaton and Dupriez (2011) and Dikhanov (2010) document that richer Indian households buy higher price goods. However, these papers focus on spatial differences in prices within India and not the price income relation itself and its implication for the size distribution. Bils and Klenow (2001) show that richer households in the US also buy higher priced durable products. The fact that larger

⁴There is also a recent quantitative literature which looks at the role of distortionary policies in explaining cross-country differences in Total Factor Productivity. See Guner, Ventura, and Yi (2008), Alfaro, Charlton, and Kanczuk (2009), García-Santana and Pijoan-Mas (2010), DiCecio and Barseghyan (2010), Hsieh and Klenow (2012), and Restuccia and Rogerson (2013). A smaller literature consider the effect of trust and social capital in determining firm size (Bloom, Sadun, and Van Reenen (2012)).

⁵The idea of quality dualism between the formal and the informal sector has been looked at by Banerji and Jain (2007), who develop a partial equilibrium model in which formal sector establishments have a comparative advantage in producing higher quality goods due to differences in factor prices across the two sectors. However, their partial equilibrium model does not have implications for the size distribution of firms and its relation to income levels.

plants produce higher price goods and use higher price inputs is shown using Colombian data by Kugler and Verhoogen (2012). They also interpret these price differences as representing quality differences and develop a model in which more productive firms choose to produce higher quality goods at a higher unit cost. I document similar facts for India. Unlike Kugler and Verhoogen (2012), I combine data from the formal and informal sector to show that the price size relation also holds when we include very small plants in the sample (the Colombian data only has plants of size ten or more).⁶ On the modeling front, I focus on non-homothetic preferences and its effect on the size distribution which is not explored in Kugler and Verhoogen (2012). Faber (2012) documents similar consumer and producer side facts as in this paper using Mexican data, but focuses on the effect of trade liberalization on income inequality.

A number of papers, especially related to international trade, have developed models of non-homothetic preferences with respect to quality. These include Flam and Helpman (1987), Mitra and Trindade (2005), Dalgin, Mitra, and Trindade (2008), and Choi, Hummels, and Xiang (2009). The model I develop is most closely related to the model in Fajgelbaum, Grossman, and Helpman (2011). Their model features non-homothetic preferences with respect to quality where the non-homotheticity arises due to complementarity between the homogenous good and quality. The non-homotheticity with respect to quality in my model arises due to complementarity between the quantity of the good consumed and quality.

The rest of the paper is structured as follows: Section 2 documents that richer households buy higher price goods and that larger plants produce higher price goods and use higher price inputs. Section 3 presents the model and Section 4 discusses the calibration. Section 5 presents the results for the counterfactual exercises and explores the sensitivity of the results to some key parameters. Section 6 considers the role of inter-state trade in explaining the cross-state relation seen in Figure 2 and Section 7 concludes.

2. Empirical Results

In this section, I provide empirical evidence which is consistent with my hypothesis of richer households consuming higher quality products which are produced by larger plants. In particular I show the following facts:

1. Richer households buy higher price goods
2. Larger plants produce higher price goods

⁶There is a large international trade literature which documents heterogeneity in prices either at the product or the firm level for exports and imports and interprets these price differences as quality differences. Some papers in this literature include Schott (2004), Hummels and Klenow (2005), Hallak (2006), Mandel (2010), Hallak and Sivadasan (2011) Manova and Zhang (2012), and Iacovone and Javorcik (2012).

Table 1: Household Regressions: Richer Households Buy Higher Price Goods

Dependent Variable: log(price)			
	(1)	(2)	(3)
log(per-capita expenditure)	0.112*** (0.0006)	0.111*** (0.0006)	0.106*** (0.0006)
Price Ratio (75th to 25th %tile)	1.091	1.090	1.086
Price Ratio (95th to 5th %tile)	1.249	1.246	1.234
Winsorize 1%		Y	Y
Exclude product from RHS			Y
Observations	5,348,463	5,348,463	5,348,463
Number of products	188	188	188
Clusters	124,635	124,635	124,635

Notes: The data is from the Consumer Expenditure Survey of 2004-05. Column 1 reports results for the regression of log of price paid by households for different goods on log of per-capita expenditure of the households. Column 2 winsorizes 1 percent tails of per-capita expenditure and goods prices. Column 3 excludes the expenditure on the good itself from the independent variable. Regressions include fixed effects for the interaction of each good, state, rural-urban cell. The price ratio implied by the coefficient estimates for different percentiles of per-capita expenditure are reported in the rows called "Price Ratio". Standard errors are clustered at the household level. ***p<0.01.

3. Larger plants use higher price material inputs and hire more skilled labor

The facts are documented using four Indian surveys. I give a brief description of each survey along with the main results in the sections that follow.

2.1. Richer Households Buy Higher Price Goods

This sections shows that richer households buy higher price goods, which is consistent with them consuming higher quality products. I use data from the Consumer Expenditure Survey of 2004-05 conducted by the National Sample Survey Office (NSS) of India. About 125,000 households from all Indian states and union-territories were interviewed for the survey. The survey asks households to report the value of consumption for 339 different goods. Households report quantities and rupee values separately for 209 goods, which can be used to compute prices for these goods. More details about the survey can be found in Appendix A.3.

I run regressions of the form

$$\ln(P_{h,g}) = \alpha_{g,state,rural} + \beta \ln(c_h) + \varepsilon_{h,g},$$

where $P_{h,g}$ is the price paid by household h for good g , c_h is per-capita expenditure of the household excluding durables, and $\alpha_{g,state,rural}$ represents fixed effects for each product, state, and urban-rural cell. c_h is a

proxy for the income level of the household, adjusting for household size.⁷ $\alpha_{g,state,rural}$ controls for the fact that different goods have different average price levels and that these price levels can vary across rural and urban areas and across states. For example, real estate prices might differ across rural and urban areas or across states with different levels of per-capita income and this can drive differences in cost of living and all prices. The fixed effects ensure that the price-income relation is not identified out of differences in average price levels across states of different income levels or across rural-urban area. Intuitively, the coefficient β is the elasticity of price with respect to per-capita consumption level and is identified out of variation in prices paid for the same good by households of different income levels within a state and urban-rural sector.

Column 1 of Table 1 reports the estimate of β , the elasticity of price with respect to per-capita consumption, based on 188 goods.⁸ The point estimate for β is 0.112 which implies that the average price paid by the 95th percentile household in terms of per-capita expenditure is 24.9 percent more than the price paid by the 5th percentile household (the 95th percentile household’s per-capita expenditure is about seven times that of the 5th percentile household). Column 2 shows that winsorizing 1 percent tails for per-capita expenditure and prices (for a good within a state and urban-rural cell) doesn’t change the results substantially.

A possible concern with the results in columns 1 and 2 in Table 1 is that the independent variable is itself a function of the dependent variable as per-capita expenditure sums the expenditure of the household across all goods, i.e., $c_h = \frac{\sum_g P_{h,g} Q_{h,g}}{\text{household size}}$ where $Q_{h,g}$ is the quantity consumed by household h of good g . This can give rise to a mechanical correlation and also cause a bias if the variables are measured with error. To account for this, column 3 repeats the regression from column 2 with the independent variable replaced by $\log\left(c_h - \frac{P_{h,g} Q_{h,g}}{\text{household size}}\right)$, i.e., the expenditure on good g is subtracted from per-capita expenditure. The results in column 3 of Table 1 are very similar to columns 1 and 2.

Figure 3 plots the non-parametric equivalent of the the regression in column 3 of Table 1. It estimates a kernel-smoothed local linear regression of residualized log prices (removes good, state, and urban-rural fixed effects) on residualized log of per-capita expenditures.⁹ As seen in the figure, a constant elasticity of price with respect to per-capita expenditure is a very good fit for the data.

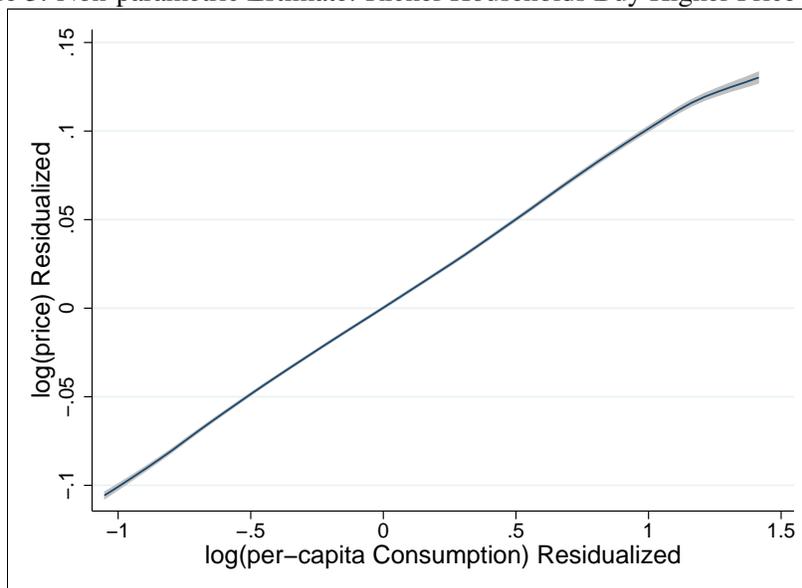
The results in Table 1 show that richer households buy goods at a higher unit price which is consistent with the hypothesis that they buy higher quality goods. However, as documented by Aguiar and Hurst (2007),

⁷Purchase of durables is excluded as these are lumpy, infrequent purchases. Two households with the same level of permanent income might have very different levels of durable expenditure in any particular year simply because of differences in timing of durable purchases.

⁸Although prices can be computed for 209 goods, only 188 were included in the regression. The goods excluded were a) all heavy durables, b) all goods with the word “other” mentioned in the description. The results do not change substantially if these goods are included.

⁹Log price and log of per-capita expenditure are demeaned within each good, state, and urban-rural cell. The residuals from this procedure are used to run a kernel-smoothed local linear regression with an Epanechnikov kernel and a bandwidth of 0.13. The top and bottom 1 percent of residualized log of per-capita expenditure are excluded.

Figure 3: Non-parametric Estimate: Richer Households Buy Higher Price Goods



Notes: The data is from the Consumer Expenditure survey of 2004-05. The graph plots the kernel-smoothed local linear regression of residualized log prices on residualized log per-capita expenditures (removes the interaction of good, state, and urban-rural fixed effects). As in column 3 of Table 1, the goods own value of consumption is subtracted from per-capita expenditure. 1 percent tails of residualized log per-capita expenditure are excluded. An Epanechnikov kernel with a bandwidth of 0.13 is used. The grey regions is the 95 percent confidence interval for the non-parametric estimate.

Table 2: Household Regressions: Controlling for Opportunity Cost of Time

Dependent Variable: log(price)						
	(1)	(2)	(3)	(4)	(5)	(6)
log(per-capita expenditure)	0.102*** (0.0010)	0.102*** (0.0010)	0.094*** (0.0015)	0.105*** (0.0029)	0.104*** (0.0029)	0.099*** (0.0036)
non-worker present		0.020*** (0.0011)	-0.059*** (0.0113)		0.017*** (0.0027)	-0.065** (0.0318)
(non-worker present)*pce			0.012*** (0.0017)			0.011** (0.0045)
Household Size	All	All	All	1 and 2	1 and 2	1 and 2
Observations	1,822,762	1,822,762	1,822,762	219,390	219,390	219,390
Number of products	169	169	169	169	169	169
Clusters	41,013	41,013	41,013	6,161	6,161	6,161

Notes: The data is from the Consumer Expenditure Survey of 2003. Column 1 reports results for the regression of log of price paid by households for different goods on log of per-capita expenditure (replicating Column 1 of Table 1). Column 2 includes a control for opportunity cost of time, namely a variable which takes value 1 if there is at least one non-working adult in the household. Column 3 also includes the interaction of this variable with per-capita expenditure. Columns 4, 5, and 6 repeat the specifications in 1,2, and 3 but restrict the sample to households of size 1 and 2 only. Regressions include fixed effects for each good, state, rural-urban cell. Standard errors are clustered at the household level. ***p<0.01, **p<0.05.

households might be paying different prices for the same good because households with higher opportunity cost of time tend to shop around less for lower prices. If richer households have a higher opportunity cost of time, then the findings in Table 1 might be a result of less time spent shopping by richer households and not because of purchase of higher quality goods.¹⁰

The 2003 Consumer Expenditure Survey asked each individual in the household the main activity they were engaged in (whether they were employed, studying, attending to domestic duties, retired etc).¹¹ I use this to construct a proxy variable which takes value 1 if the household has at least one member between the age of 15 and 70 who is only attending to domestic duties or is retired, and 0 otherwise.¹² I interpret households with a non-worker present as households with low opportunity cost of time and include this variable as a control in the regressions. Column 1 of Table 2 repeats the regression from Column 1 of Table 1, but with the 2003 data instead of the 2004-05 data. Column 2 of Table 2 now adds the measure of “non-worker present” as an additional control. Although the coefficient on the “non-worker present” variable is positive, the key point is that the coefficient of per-capita expenditure does not change substantially. Column 3 also includes the interaction of the “non-worker present” variable with per-capita expenditure and this does not change the results substantially either. Columns 4, 5, and 6 repeat the regressions from columns 1, 2, and 3 respectively, but restrict the sample to include households with one or two members only. This controls for the fact that larger households are more likely to have non-working adults. Again, the coefficient on per-capita expenditure does not change substantially when including the “non-worker present” variable as a control.

The results in this section indicate that richer households tend to buy higher price goods, which is consistent with the hypothesis that they are consuming higher quality products.

2.2. Larger Plants Produce Higher Price Goods

This section shows that larger plants produce higher price goods, which is consistent with the hypothesis that high quality goods are produced in large plants. To show this, I combine data from the Annual Survey of Industries (ASI) of 2005-06 and the Survey of Unorganized Manufacturing (SUM) of 2005-06. The ASI

¹⁰For developing countries, there is evidence that poorer households might in fact be paying more for the same product as opposed to rich households which would imply that the estimates for β are a lower bound for the quality-income relation. For example, Attanasio and Frayne (2006) find that poor people in rural Columbia are less likely to avail of bulk discounts and thus end up paying more for the same product as compared to richer households.

¹¹Unfortunately, the 2004-05 Consumer Expenditure Survey does not ask this question so this exercise cannot be conducted using the same data used in Table 1. The 2003 survey has only one fourth the number of households as the 2004-05 survey. However, the point estimates for the elasticity of price with respect to per-capita expenditure (β) are quite similar across the two surveys.

¹²Table A.3 in the appendix lists the possible responses for the question regarding main activity of the individual. People who reported codes 92, 93, 94, or 97 were classified as non-workers.

Table 3: Plant Regressions: Larger Plants Produce Higher Price Goods

Dependent Variable: log(output price)			
	(1)	(2)	(3)
log(labor)	0.096*** (0.0087)	0.053*** (0.0192)	0.106*** (0.0133)
Price Ratio (Size 50 to 5)	1.247	1.130	1.276
Price Ratio (Size 500 to 5)	1.556	1.276	1.629
Sample	ASI	SUM	BOTH
Winsorize 1%	Y	Y	Y
Observations	46,704	28,457	75,161
Number of products	1,217	2,739	3,181
Number of clusters	1,078	2,731	3,042

Notes: The data is from the ASI and SUM for 2005-06. All columns report results for regressions of log price charged by plants for their products on log of number of employees hired by the plant. Column 1 restricts the sample to the ASI, Column 2 restricts the sample to the SUM, while column 3 combines the two. 1 percent tails of prices (within a product) and plant size are winsorized. Regressions include product fixed effects and state times urban-rural fixed effects. Standard errors are clustered at the product level. The number of product fixed effects exceed the number of clusters because of the units problem discussed in the Appendix as the misreported units are treated as a different product category for fixed effects but not for clustering. The price ratio for different sized plants implied by the coefficient estimates are reported in the rows called "Price Ratio". ***p<0.01.

covers all manufacturing plants registered under the Factories Act, 1948. This includes manufacturing plants employing twenty or more workers and not using electricity or employing ten or more workers and using electricity. The SUM on the other hand covers the smaller manufacturing plants not covered by the ASI. The two surveys together should provide a representative sample of the manufacturing sector as a whole.¹³

Both the surveys ask manufacturing establishments detailed questions about the products they produce and inputs they use. Each establishment reports the quantity of the product it produces (for a 5-digit product classification, which has about 5,500 possible products) and its value (before taxes and distribution expenses) which can be used to compute prices. For the ASI, each products quantity is supposed to be reported for a standardized unit (kilograms, numbers, etc). In the SUM, different plants can report the same products price in different units. I concord units across the two survey so that the price of the same product is not getting compared for different units.¹⁴

¹³A number of recent papers have combined these two surveys to construct a dataset which is representative of the manufacturing sector as a whole. These include Hasan and Jandoc (2010), Nataraj (2011), Hsieh and Klenow (2012), and Ghani, Goswami, and Kerr (2012).

¹⁴In the ASI all plants reporting a certain product are supposed to report quantities in the same units. However, there are clear cases in which plants are misreporting quantity units. For example, all plant which produce milk are supposed to report quantities in terms of kiloliters which means that the price computed by dividing the rupee value by the quantity should yield prices per kiloliter. However, there is a group of plants whose prices are approximately 1000 times lower than others. This is clearly a case of some plants reporting quantities in liters instead of kiloliters. I have manually gone through all product categories and identified products with this problem and split these into two separate categories based on a sensible price cutoff. In addition to this manual

I run regressions of the form

$$\ln(P_{f,g}) = \alpha_g + \alpha_{state,rural} + \gamma \ln(L_f) + \varepsilon_{f,g},$$

where $P_{f,g}$ is the price charged by plant f for product g , L_f is the number of workers employed by plant f , α_g is a product fixed effect, and $\alpha_{state,rural}$ is a state times urban-rural fixed effect. Intuitively, the coefficient γ is the elasticity of the price of output produced with respect to plant size and it is identified out of variation in prices charged by plants of different sizes producing the same product (reported in the same units) and allowing for differences in average price levels across states and urban and rural areas.

Column 1 of Table 3 reports results when the sample is restricted to the ASI only. The estimate for the elasticity of price with respect to size, γ , is 0.096 and is statistically significant at the 1 percent level. The point estimate implies that a plant which employs 500 people on average charges a price which is 55.6 percent more than a plant employing 5 workers.¹⁵

Column 2 report results when the sample is restricted to the SUM only. The point estimate for the coefficient γ (elasticity of price with respect to size) is still positive but smaller. This is not surprising as the variation in employment levels within the SUM is small with 95 percent of the plants employing 16 workers or less.

Column 3 reports results when the two surveys are combined. The estimate for the elasticity of price with respect to size implies that a plant which employs 500 people on average charges a price which is 62.9 percent more than a plant employing 5 workers.

Figure 4 plots the non-parametric equivalent of the the regression in column 3 of Table 3. In particular, it estimates a kernel-smoothed local linear regression of residualized log prices (after removing product fixed effects and state times urban-rural fixed effects) on residualized log of plant size.¹⁶ Again, the non-parametric estimates suggest that the price size relation across plants is close to log-linear.

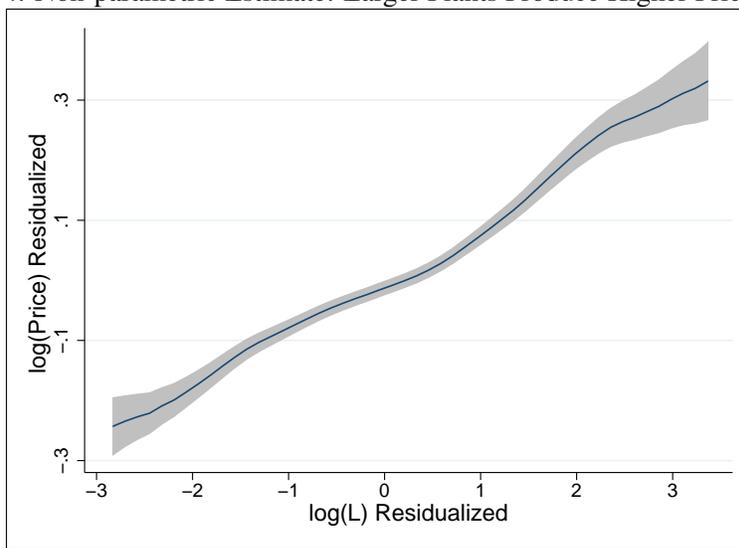
The fact that larger plants produce goods which they sell at a higher price is consistent with the hypothesis that larger plants produce higher quality products.

check, I have also implemented an algorithm to identify these problem products and used the algorithm generated cutoff's to split problematic products. The results are similar to the ones reported here. Appendix C gives more details regarding this problem and how it is being tackled.

¹⁵Note that the formal plants surveyed in the ASI report the value of output before taxes and distribution costs. Therefore, the price-size relation documented here is not driven mechanically by the fact that larger plants might be paying taxes while the smaller plants are not.

¹⁶Log price and log of employment of each plant is regressed on product and state times urban rural fixed effects. The residuals from this procedure are used to run a kernel-smoothed local linear regression with an Epanechnikov kernel and a bandwidth of 0.502. The top and bottom 1 percent of residualized log of employment are excluded.

Figure 4: Non-parametric Estimate: Larger Plants Produce Higher Price Goods



Notes: The data is from the ASI and the SUM of 2005-06. The graph plots the kernel-smoothed local linear regression of residualized log prices charged by a plant for its products on residualized log employment of that plant (removes product fixed effects and the interaction of state and urban-rural fixed effects). Products which have the units problem discussed in footnote 14 and in Appendix C are split into two product categories. 1 percent tails of residualized log employment are excluded. An Epanechnikov kernel with a bandwidth of 0.502 used. The grey regions is the 95 percent confidence interval for the non-parametric estimate.

2.3. Larger Plants Use Higher Price Inputs

This section looks at the relation between the size of a plant and the inputs it uses. First I show that larger plants pay a higher price for the same material input as compared to smaller plants. This is consistent with the idea that larger plants produce higher quality products which require higher quality inputs. I then show that larger plants hire more educated workers as compared to small plants.

As in the last section, the ASI and SUM are used to show that larger plants use higher price material inputs. Each establishment reports the material inputs it uses (for a 5-digit product classification, which has about 5,500 possible products) and the price it pays for the input. The units between the surveys are again concorded.¹⁷

I run a regression of the form

$$\ln(P_{f,i}) = \alpha_i + \alpha_{state,rural} + \gamma \ln(L_f) + \varepsilon_{f,i},$$

where $P_{f,i}$ is the price paid by plant f for input i , L_f is the number of workers employed by plant f , α_i is a

¹⁷The same problem of unit misreporting in the ASI discussed in footnote 14 is also present for inputs. I perform the same correction for this problem as I did in the previous section. The data appendix provides more details.

Table 4: Plant Regressions: Larger Plants Use Higher Price Inputs

Dependent Variable: log(input price)			
	(1)	(2)	(3)
log(labor)	0.077*** (0.0072)	0.033* (0.0193)	0.050*** (0.0104)
Price Ratio (Size 50 to 5)	1.194	1.079	1.122
Price Ratio (Size 500 to 5)	1.426	1.164	1.259
Sample	ASI	SUM	BOTH
Winsorize 1%	Y	Y	Y
Observations	107,325	105,422	212,747
Number of products	2,189	4,316	5,257
Number of clusters	1,502	4,241	4,569

Notes: The data is from the ASI and SUM for 2005-06. All columns report results for regressions of log of price paid by establishments for material inputs used on log of number of employees hired by the establishment. Column 1 restricts the sample to the ASI only. Column 2 restricts the sample to the SUM only while column 3 combines the ASI and the SUM. 1 percent tails of prices (within a product) and plant size are winsorized. All regressions include product fixed effects and state times urban-rural fixed effects. Standard errors are clustered at the product level. The number of product fixed effects exceed the number of clusters because of the units problem discussed in the Appendix as misreported units are treated as a different input category for fixed effects but not for clustering. The price ratio implied by the coefficient estimates for different sized plants are reported in the rows called "Price Ratio ". ***p<0.01, *p<0.1.

product fixed effect, and $\alpha_{state,rural}$ is a state times urban-rural fixed effect. Intuitively, the coefficient γ is the elasticity of the price paid for inputs with respect to plant size and it is identified out of variation in prices paid by plants of different sizes for the same inputs (reported in the same units), controlling for differences in average prices across states and urban-rural sectors.

Column 1 of Table 4 reports results when the sample is restricted to the ASI only. The estimate for the elasticity of input prices with respect to plant size, γ , is 0.077 and is statistically significant at the 1 percent level. The point estimate implies that a plant which employs 500 people on average pays prices for inputs which are 42.6 percent more than a plant employing 5 workers. Column 2 reports results when the sample is restricted to the SUM only. The coefficient γ is positive but smaller.

Column 3 reports results when the two surveys are combined. When combining the two surveys, the estimate for the elasticity of input prices with respect to size implies that a plant which employs 500 people on average pays a price for inputs which is 25.9 percent more than a plant employing 5 workers.

Not only do larger plants use higher price inputs, but they also employ more skilled labor. To show this I use the Employment-Unemployment Survey of 2004-05 conducted by the National Sample Survey Office (NSS) of India. Note that plants in the ASI and SUM do not report the education level of their workers, hence they cannot be used to look at the relation between plant size and education levels of workers.

The Employment-Unemployment Survey records demographic information (including education levels)

Table 5: Larger Plants Hire More Educated Workers

	No School	Grade 1 to 9	Grade 10 to 12	> Grade 12
$L \leq 5$	0.43	0.41	0.13	0.03
$5 < L \leq 10$	0.34	0.41	0.17	0.08
$10 < L \leq 20$	0.33	0.41	0.16	0.10
$L > 20$	0.23	0.32	0.22	0.22

Notes: The data is from the Employment-Unemployment Survey of 2004-05. The rows of the table represent the size category of the establishment in which an individual works while the columns represent the education level. Each number represents the share of individuals in the given size category who have attained the level of education given by the column.

for about 600,000 individuals. It also asks individuals to report the size category of establishment in which they work where the size category can take five values - establishment of size less than 6, between 6 and 9, between 10 and 19, 20 or greater, and unknown size. Table 5 reports the skill composition of workers for the different size categories. Out of the workers in establishments of size less than 6, 43 percent have never attended school while only 3 percent have graduated from high school. On the other hand, out of workers in establishments of size more than 20, only 23 percent have never attended school while 22 percent have graduated high school. As can be seen, a larger share of workers in big establishments have high levels of education.

3. Model

This section develops a general equilibrium model which matches the facts described in Section 2. In particular, I model consumers choice between different quality levels with richer households more likely to buy high quality goods. On the production side, I assume that production of better quality requires larger fixed costs which along with free entry implies that high quality producers are larger on average.

3.1. Households

There are a mass L of households in the economy indexed by the subscript j . Share h of the households are skilled and earn wage w_S (which is determined endogenously in equilibrium) while share $1 - h$ are unskilled and earn wage w_U . Unskilled wage w_U is assumed to be the numeraire and is normalized to 1.¹⁸

There are N quality levels. $Q = \{q_1, q_2, \dots, q_N\}$ denotes the the set of qualities available in the economy. The quality indexes q_n are arranged in ascending order of quality with $q_n > q_m \forall n > m$. Therefore q_1 is the

¹⁸Having two skill levels with different wages is crucial for my exercise as it generates cross-sectional differences in income levels in the model. This cross-sectional variation in income levels allows me to calibrate the extent of non-homotheticity in the model to match the price-income slope documented in Section 2.1.

quality index of the lowest quality level and q_N is the quality index of the highest quality level.

The utility derived by household j from consuming quality level q_n is given by

$$u_{j,q_n}(c_{j,q_n}, \varepsilon_{j,q_n}) = a_{q_n} + q_n \log(c_{j,q_n}) + \varepsilon_{j,q_n} \quad \forall q_n \in \mathcal{Q}, \quad (1)$$

where a_{q_n} is a constant in the utility function which can vary by quality level, c_{j,q_n} is the quantity consumed of quality level q_n by household j , and ε_{j,q_n} is a random utility component which represents the idiosyncratic valuation of quality level q_n by household j . The fact that higher quality levels have higher indexes q_n implies that for any given level of quantity consumed, households get more utility from consuming higher quality goods.

The random utility component ε_{j,q_n} is assumed to be independently and identically distributed with a Gumbel Type 1 Extreme Value distribution with density

$$f(\varepsilon_{j,q_n}) = e^{-\varepsilon_{j,q_n}} e^{e^{-\varepsilon_{j,q_n}}}.$$

As shown by McFadden (1974) (see also Chapter 3 of Train (2009)), assuming a Gumbel distribution for the random utility component implies simple closed form expressions for demands.

I assume that a household can choose to consume only one quality level and spends its entire income on the quality level that it chooses. This implies that the indirect utility function of household j if it chooses to consume quality level q_n is given by

$$v_{j,q_n}(w_j, P_{q_n}, \varepsilon_{j,q_n}) = a_{q_n} + q_n \log\left(\frac{w_j}{P_{q_n}}\right) + \varepsilon_{j,q_n} \quad \forall q_n \in \mathcal{Q}, \quad (2)$$

where P_{q_n} is the price of quality level q_n , and w_j represents the wage of household j . Equation (2) is simply equation (1) but with $c_{j,q_n} = \frac{w_j}{P_{q_n}}$ reflecting the assumption that each household can only choose to consume one quality level.

Each household j receives draws of the random utility component ε_{j,q_n} for each quality level q_n and given these draws, chooses to consume the quality level which gives it the highest utility level. Therefore, household j chooses to consume quality level q_n if and only if

$$v_{j,q_n}(w_j, P_{q_n}, \varepsilon_{j,q_n}) > v_{j,q_m}(w_j, P_{q_m}, \varepsilon_{j,q_m}) \quad \forall n \neq m.$$

Let $\rho(q_n|w)$ be the share of households with wage w who choose to consume quality level q_n . Given the assumption that ε_{j,q_n} is independently and identically distributed with a Gumbel distribution, this share

takes the simple logit form

$$\begin{aligned}\rho(q_n|w) &= \frac{e^{a_{q_n} + q_n \log\left(\frac{w}{P_{q_n}}\right)}}{\sum_{i=1}^N e^{a_{q_i} + q_i \log\left(\frac{w}{P_{q_i}}\right)}} \quad \forall q_n \in \mathcal{Q} \\ &= \frac{e^{a_{q_n}} \left(\frac{w}{P_{q_n}}\right)^{q_n}}{\sum_{i=1}^N e^{a_{q_i}} \left(\frac{w}{P_{q_i}}\right)^{q_i}} \quad \forall q_n \in \mathcal{Q}.\end{aligned}\tag{3}$$

Analyzing how $\rho(q_n|w)$ changes as wage changes can help understand how this preference structure leads to non-homotheticity with respect to quality choice. Define $\gamma_{\rho(q_n),w}$ to be the elasticity of $\rho(q_n|w)$ with respect to wages w . Taking logs and differentiating equation (3) with respect to $\log(w)$ yields

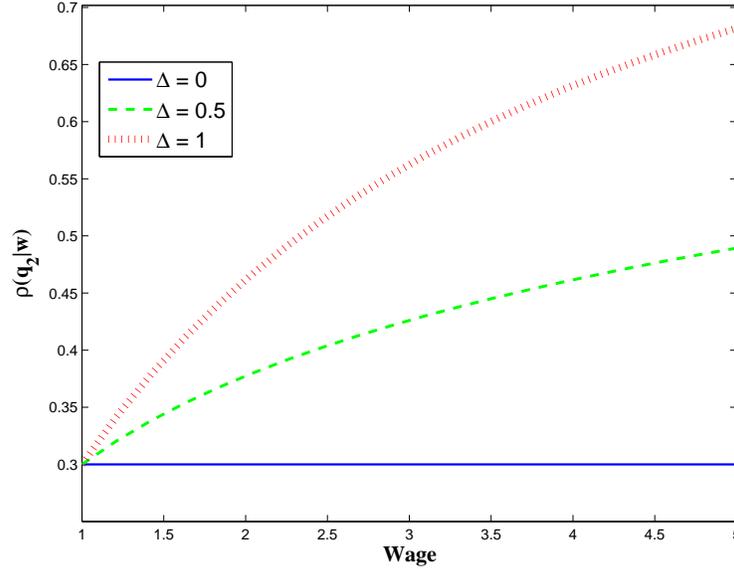
$$\gamma_{\rho(q_n),w} = \frac{\partial \log[\rho(q_n|w)]}{\partial \log(w)} = q_n - \sum_{i=1}^N q_i \rho(q_i|w).$$

The elasticity of $\rho(q_n|w)$ with respect to wages w is simply the quality index q_n minus a weighted average of all the quality indexes where the weights are the share of households with wage w who buy each quality level. A positive elasticity ($q_n > \sum_{i=1}^N q_i \rho(q_i|w)$) implies that as wages increase, a larger share of the households buy the quality q_n . As lower quality goods have a lower quality index ($q_n > q_m \forall n > m$), the lowest quality level will always have a negative elasticity i.e. the share of household who buy the lowest quality level will always go down as wages increase. Furthermore, the highest quality level will always have a positive elasticity implying that the share of households who consume the highest quality always goes up as wage levels increase.

Therefore, the non-homotheticity with respect to quality operates on the extensive margin. As a household becomes richer, it is more likely to choose the higher quality goods. There is a positively sloped “quality Engel curve” where households with higher levels of wages will, on average, spend a larger share of their expenditure on higher quality goods. This arises because the utility function in equation (1) features complementarity between quantity consumed and quality. As wages increase, the household can consume more quantity of whichever quality level that it chooses. Complementarity between quantity and quality implies that the marginal increase in utility from a given increase in wage is larger for higher quality goods which leads to more households choosing higher quality levels as wages increase (given the draw of ε_{j,q_n}).

The steepness of the quality Engel curve is determined by the differences in the quality indexes across quality levels. One way of parameterizing the quality indexes would be to set the index for the lowest quality level to be one and assume that each higher quality level has an index which is a constant Δ larger than the

Figure 5: Quality Engel Curve



Notes: The figure plots the share of households who purchase the high quality product for different wage levels. There are only 2 quality level ($N = 2$) which have prices $P_{q_1} = 1$. Quality index for the low quality is set to one i.e. $q_1 = 1$. The three lines correspond to three different values of Δ where $q_2 = 1 + \Delta$. a_{q_2} , the constant for the high quality is chosen such that 30 percent of households with wage equal to one choose the high quality.

previous quality index i.e. $q_1 = 1$ and $q_n = q_{n-1} + \Delta$. In this case, the size of the constant Δ determines the extent of non-homotheticity with a larger Δ implying that demand shifts to higher quality faster as wages increase.

Consider the following simple example which illustrates this relation between the size of Δ and the extent of the non-homotheticity. Assume that there are only two quality level ($N = 2$) which have prices $P_{q_1} = 1$ and $P_{q_2} = 1.5$ and quality indexes $q_1 = 1$ and $q_2 = 1 + \Delta$.¹⁹ Figure 5 plots the share of households who choose the high quality level q_2 as a function of wages for different value of Δ .²⁰ For each value of Δ , the constant in the utility function a_{q_2} is chosen such that 30 percent of the households with wage equal to one choose the high quality q_2 .²¹

For the case with $\Delta = 0$, there is no change in the share of households who buy the high quality as wage increases. This is expected as $\Delta = 0$ is in effect the case in which there is no quality distinction between the goods. For positive values of Δ , there is an increase in the share of households who buy the high quality

¹⁹In the full calibration done in Section 4, there is a richer quality space with $N = 12$. Here, to illustrate the non-homotheticity, the simplifying assumption of $N = 2$ is made.

²⁰The results in Figure 5 can be viewed as the choice made by an individual *if* they faced a continuous wage profile. However, only two wages will exist in equilibrium (the unskilled and the skilled wages).

²¹Only $N - 1$ constants in the utility function are identified as what matters for consumer choice is the difference in utility across quality levels. Therefore, for the case with $N = 2$, only one constant needs to be calibrated.

good as wages increase, and this increase is larger for higher values of Δ .

Given prices and the wages of skilled and unskilled workers, the total demand for quality level q_n is given by

$$C_{q_n} = \underbrace{Nh\rho(q_n|w_S)\frac{w_S}{P_{q_n}}}_{\text{demand from skilled households}} + \underbrace{N(1-h)\rho(q_n|w_U)\frac{w_U}{P_{q_n}}}_{\text{demand from unskilled households}} \quad \forall q_n \in Q. \quad (4)$$

The first term is the demand for quality q_n from skilled households which is the product of the number of skilled households (Nh), the share of skilled households who choose quality q_n ($\rho(q_n|w_S)$), and the quantity consumed by each skilled household who consumes quality q_n ($\frac{w_S}{P_{q_n}}$). Similarly, the second term is the demand for quality q_n from unskilled households.

In summary, the consumers choose between different quality levels and complementarity between quality and quantity implies that richer households are more likely to consume higher quality. This non-homotheticity with respect to quality will help match the patterns seen in Table 1 (richer households buy higher price goods).

3.2. Final Goods Producers

There are N competitive final goods producers, one for each quality level. In addition to the vertical differentiation across quality levels, there is horizontal differentiation in products within a quality level. The final goods producer of quality q_n combines intermediate varieties (horizontal differentiation) of quality q_n to produce the composite final good of that quality. Each final goods producer has a constant elasticity of substitution (CES) production function given by

$$Y_{q_n}^S = \frac{1}{M_{q_n}^{\frac{\sigma-1}{\sigma}}} \left(\sum_{i=1}^{M_{q_n}} x_{i,q_n}^{\frac{\sigma-1}{\sigma}} \right)^{\frac{\sigma}{\sigma-1}}, \quad \forall q \in Q$$

where i indexes varieties, M_{q_n} is the number of varieties (or plants) of quality q_n present in the economy which will be determined by free entry, x_{i,q_n} is the quantity of variety i of quality q_n used by the final quality producer of quality q_n ,²² and σ is the elasticity of substitution between different varieties of the same quality.

The multiplicative factor $\frac{1}{M_{q_n}^{\frac{\sigma-1}{\sigma}}}$ in the production function scales out the love of variety from the CES production function. This ensures that the price difference between different quality levels does not reflect differences in number of varieties available. I maintain this assumption of no love of variety in the baseline

²²Note that the pair (i, q_n) together identifies a variety uniquely in the economy. i represents the horizontal differentiation dimension while q_n represents the vertical differentiation dimension. For example, $(i = 1, q_1)$ represents the first variety of lowest quality q_1 while $(i = 1, q_N)$ represents the first variety of the highest quality.

specification for two reasons. Firstly, assuming no love of variety is the conservative choice as changes in the size distribution in the counterfactual exercises are smaller in this case as opposed to the case with love of variety. Secondly, allowing for love of variety makes the changes in size distribution in the counterfactual sensitive to the average level of the quality indexes q_n which is a difficult parameter to calibrate as it represents the own price elasticity of each quality level with respect to the unobserved CES price index of that quality.²³ Therefore, while the baseline results presented in Section 5.1 maintains the assumption of no love of variety, Section 5.3 provides results when allowing for love of variety and further discusses the sensitivity of the results to the average level of the quality indexes q_n .

The final quality producers take the prices of intermediate varieties, p_{i,q_n} , as given and solve their cost minimization problem

$$\begin{aligned} & \min_{x_{i,q_n}} \sum p_{i,q_n} x_{i,q_n} \\ \text{s.t. } & Y_{q_n}^s = \frac{1}{M_{q_n}^{\frac{1}{\sigma-1}}} \left(\sum_{i=1}^{M_{q_n}} x_{i,q_n}^{\frac{\sigma-1}{\sigma}} \right)^{\frac{\sigma}{\sigma-1}}, \quad \forall q_n \in \mathcal{Q}. \end{aligned}$$

This yields their demand curves

$$x_{i,q_n} = p_{i,q_n}^{-\sigma} M_{q_n}^{\frac{1}{\sigma-1}} Y_{q_n}^s \left(\sum_{i=1}^{M_{q_n}} p_{i,q_n}^{1-\sigma} \right)^{\frac{\sigma}{1-\sigma}} \quad \forall q_n \in \mathcal{Q}, \quad (5)$$

which are taken as given by downstream intermediate producers. The final quality producers make zero profits. The price that they charge consumers is given by

$$P_{q_n} = \frac{\sum_{i=1}^{M_{q_n}} p_{i,q_n} x_{i,q_n}}{Y_{q_n}^s}, \quad \forall q_n \in \mathcal{Q}.$$

Given the assumption of no love of variety, P_{q_n} will be independent of the number of varieties M_{q_n} available in the economy.

²³As mentioned in Section 3.1, I parametrize the quality indexes using the recursion $q_n = q_{n-1} + \Delta$, where the size of Δ determines the steepness of the quality Engel curves. With no love of variety, the choice of the level of q_1 (which given a Δ determines the average level of the quality indexes) does not impact the changes in size distribution in the counterfactual. However, when allowing for love of variety, the results become sensitive to the choice of q_1 .

3.3. Intermediate Goods Producers

Each variety of each quality is produced by a monopolistically competitive intermediate producer. The intermediate producers combine skilled and unskilled labor and their production function is given by

$$x(A_i, q_n) = A_{i,q_n} \left(\theta_{q_n} (l_{i,q_n}^U)^{\frac{\sigma_{su}-1}{\sigma_{su}}} + (1 - \theta_{q_n}) (l_{i,q_n}^S)^{\frac{\sigma_{su}-1}{\sigma_{su}}} \right)^{\frac{\sigma_{su}}{\sigma_{su}-1}}, \quad (6)$$

where l_{i,q_n}^U is the quantity of unskilled labor hired by variety i producer of quality q_n , l_{i,q_n}^S is the quantity of skilled labor hired by variety i producer of quality q_n , σ_{su} is the elasticity of substitution between the two types of labor, A_{i,q_n} is the idiosyncratic productivity level of variety i producer of quality q_n , and θ_{q_n} is the share parameter of unskilled labor for quality q_n producers.

Solving the cost minimization problem of the intermediate goods producer subject to the production function given in equation (6) yields the marginal cost of production for variety i of quality q_n which is given by

$$\kappa(A_i, q_n) = \frac{1}{A_{i,q_n} \left(\theta_{q_n}^{\sigma_{su}} \left(\frac{1}{w_U} \right)^{\sigma_{su}-1} + (1 - \theta_{q_n})^{\sigma_{su}} \left(\frac{1}{w_S} \right)^{\sigma_{su}-1} \right)^{\frac{1}{\sigma_{su}-1}}}.$$

The marginal costs is a function of skilled and unskilled wage, and is inversely proportional to the productivity level A_{i,q_n} .

Intermediate quality producers will take the demand curve of final quality producers (equation 5) as given and will maximize profits. As the demand curve of final quality producers is of the constant elasticity form, the optimal price charged by intermediate producers will be a constant markup over marginal cost and is given by

$$p(A_i, q_n) = \frac{\sigma}{\sigma - 1} \kappa(A_i, q_n). \quad (7)$$

To start an intermediate goods plant of quality q_n requires f_{q_n} units of labor. Share α_{q_n} of the entry labor needs to be skilled and this share is different for different quality levels. On paying the fixed cost f_{q_n} , entrant receive a productivity draw from a log normal distribution given by

$$\log(A_{i,q_n}) \sim g_{q_n} \sim N(\mu_{q_n}, \nu^2).$$

Note that the mean of the log of the productivity draw can differ across quality levels but the variance is the same.

Free entry requires that the fixed cost paid must equal the ex-ante expected profit i.e.

$$\alpha_{q_n} f_{q_n} w_S + (1 - \alpha_{q_n}) f_{q_n} w_U = \int \pi(A_i, q_n) g_{q_n}(A_i) dA_i \quad \forall q_n \in Q \quad (8)$$

where $\pi(A_i, q_n)$ is the flow profit earned by an intermediate quality producer of quality q_n with productivity draw A_i and is given by

$$\pi(A_i, q_n) = [p(A_i, q_n) - \kappa(A_i, q_n)] x(A_i, q_n).$$

The number of varieties M_{q_n} will adjust to ensure that the free entry condition holds for all quality levels.

If fixed costs for higher quality levels is larger than for lower quality levels, then for the free entry condition to hold, the scale of production $x(A_i, q_n)$ will have to be larger for higher quality producers. Furthermore, if $\theta_{q_n} > \theta_{q_m} \quad \forall n > m$ then higher quality producers will use skilled labor more intensively and will have a higher cost of production. Finally, differences in μ_{q_n} will also translate into differences in prices between different quality levels as marginal costs and prices are proportional to productivity.

3.4. Equilibrium

The equilibrium in this economy is a set of prices $\left(w_S, \left\{ \{p_{i,q_n}\}_{i \in M_{q_n}}, P_{q_n} \right\}_{q_n \in Q} \right)$, allocations $\left\{ \{c_{j,q_n}\}_{j \in L}, C_{q_n}, \{x_{i,q_n}\}_{i \in M_{q_n}}, Y_{q_n} \right\}_{q_n \in Q}$, and mass of entrants M_{q_n} such that

- Given prices P_{q_n} , wages, and draws of the random utility component (ε_{j,q_n}) , consumers choose their optimal quality level (equations 3 and 4 hold)
- Given prices, final quality producers demand optimal amounts of intermediate goods (demand follows equation 5)
- Intermediate good producers maximize profits (charge the constant markup price given by equation 7)
- Free entry conditions hold for all quality levels (equation 8)
- Markets clear

$$Y_{q_n} = C_{q_n} \quad \forall q_n \in Q$$

$$L(1-h) = \sum_{q_n} M_{q_n} \int l^U(A_i, q_n) g_{q_n}(A_i) dA_i + \sum_{q_n} M_{q_n} (1 - \alpha_{q_n}) f_{q_n}$$

$$Lh = \sum_{q_n} M_{q_n} \int l^S(A_i, q_n) g_{q_n}(A_i) dA_i + \sum_{q_n} M_{q_n} \alpha_{q_n} f_{q_n}$$

The last two equations are the labor market clearing conditions. The second last equation says that the demand for unskilled labor for production by the intermediate producers (summing over all quality levels) and entry requirements must equal the supply of unskilled labor. Similarly, the last equations says that the demand for skilled labor from intermediate producers and entry requirements must equal the supply of skilled workers.

4. Calibration

I now calibrate the model to match the cross-sectional facts documented in Section 2 and some additional moments taken from the Indian data. I then conduct counterfactual exercises in which I simulate differences in per-capita income levels in the model and see how this effects the size distribution. The key parameters which determine the *change* in size distribution in the counterfactual exercises are the degree of non-homotheticity (Δ) on the consumer side and the price-size relation on the producer side. These parameters are calibrated independently of the aggregate relation between the share of employment in small plants and income levels seen across Indian states (which is what I want to explain in the counterfactual). In particular, I use the micro-facts documented in Section 2 (richer households buy higher priced goods and larger plants produce higher priced goods) to discipline these parameters of the model.

4.1. Production Parameters

For the calibration, I define an individual with less than ten years of education as unskilled. h , the share of the labor force which is skilled, is set to 0.24, which is the share of manufacturing workers with at least ten years of education in India in 2004-05. σ_{su} , the elasticity of substitution between skilled and unskilled workers in the intermediate goods production function (equation 6), is assumed to be 1.75 which is in the range of estimates for developing countries in Behar (2009).

The elasticity of substitution between varieties for the final goods producer, σ , is set to 5, which implies a markup over cost of 25 percent for the intermediate producers and is in the range of estimates in Broda and Weinstein (2006).

This leaves five sets of parameters to be calibrated on the production side: (1) f_{q_n} , the fixed cost for each quality level; (2) θ_{q_n} , the share of unskilled workers in the production function for each quality level; (3) μ_{q_n} , the mean of the log of the productivity draw for each quality level; (4) α_q , the share of skilled labor needed for entry for each quality level; and (5) v^2 , the variance of the productivity draw which is common across all quality levels. These parameters (along with the utility parameters) are jointly calibrated as there is no

Table 6: Unskilled to Skilled Ratio for Different Size Categories

	U/S Ratio	Ratio Relative to Smallest
$L \leq 5$	5.05	1.00
$5 < L \leq 20$	2.92	0.58
$L > 20$	1.25	0.25

Notes: The data is from the Employment-Unemployment Survey of 2004-05. The rows of the table represent the size category of the establishment in which an individual works. The first column gives the ratio of skilled to unskilled workers in each size category where the definition of skilled is assumed to be an individual with at least ten years of education. The second column gives the ratio of skilled to unskilled relative to the smallest size category.

one-to-one mapping between the parameters and the target moments. However, for expositional purposes, I explain the calibration of each parameter in terms of the moments which are most informative about the parameter.

The number of quality levels N is set to 12.²⁴

The fixed costs, f_{q_n} , determines the average scale of operation of the intermediate producers of each quality level. A larger fixed cost will mean that the average size (in terms of output and employment) of intermediate producers will need to be larger in order for the free entry condition to hold. As shown in Section 2.2, larger plants tend to produce higher price products, which is indicative of higher quality goods being produced in larger plants. Therefore, the fixed costs are chosen such that the average employment (skilled plus unskilled workers) in intermediate producers of the lowest quality levels is 1.25 workers and each higher quality level has double the average size of the previous quality level i.e. the average employment of the intermediate producers of the different quality levels are $size_{q_n} = \{1.25, 2.5, 5, \dots, 2560\}$.²⁵

The level of $\theta'_{q_n} s$ determine the demand for unskilled labor relative to skilled labor and are informative about the wage premium, w_S , in the economy. The ratio of skilled to unskilled workers in any quality level relative to the lowest quality is also a function of the $\theta'_{q_n} s$ and is given by

$$ratio_{q_n}^{U,S} = \left(\frac{L_{q_n}^U}{L_{q_n}^S} \right) / \left(\frac{L_{q_1}^U}{L_{q_1}^S} \right) = \left(\frac{\theta_{q_n}}{1-\theta_{q_n}} \right)^{\sigma_{us}} / \left(\frac{\theta_{q_1}}{1-\theta_{q_1}} \right)^{\sigma_{us}} \quad \forall q_n \in Q. \quad (9)$$

Therefore, the twelve $\theta'_{q_n} s$ are chosen to match a target for the wage premium and eleven targets for unskilled to skilled ratio in different quality levels relative to the lowest quality level.

²⁴The results discussed in Section 5 are not very sensitive to the choice of N . For example, if I instead choose N to be 6, and choose all the other parameters in the same way as described below, then the model explains 45 percent instead of 43 percent of the differences in share of employment in small plants in rich versus poor states (the baseline results discussed in Section 5.1).

²⁵Different intermediate producers of the same quality will have different levels of employment due to heterogeneity in the productivity draw. Within the same quality level, intermediate producers with higher productivity draws will be larger compared to those with lower productivity draws. The fixed costs are chosen such that the average employment level of the producers within a quality level matches the target $size_{q_n} = \{1.25, 2.5, 5, \dots, 2560\}$.

The targets for these moments are obtained from the Employment-Unemployment Survey conducted by the NSS in 2004-05 (see Section 2.3 and Appendix A.4 for details about the dataset).

The target for the wage premium is set at 1.6, and is obtained from running Mincerian regressions on data from the Employment-Unemployment Survey.²⁶ Table 6 gives the ratio of unskilled to skilled workers for three different size categories, along with the ratio relative to the smallest size category, as computed from the Employment-Unemployment Survey. Smaller plants have a much higher ratio of unskilled to skilled workers indicating that low quality producers have higher θ'_{q_n} s. Unfortunately, the size categories reported in the Employment-Unemployment survey are very coarse, and therefore cannot be used to compute eleven ratios for equation (9) for eleven different quality (size) levels. I use the first two data points reported in Table 6 for the unskilled to skilled ratio (column 1) and extrapolated the relation to larger sizes (with a minimum of 0.5) to compute eleven ratios, one for each quality (size) level.

μ_{q_n} , the mean of the log of the productivity draw for each quality level, is informative about the average price of each quality level as $p(A_i, q_n) \propto \frac{1}{A_i}$. If the mean of the productivity draw for a particular quality is high, then the average price of that quality level will be lower. Therefore, the μ_{q_n} for each quality level is chosen to match the price-size relation seen in Table 3.²⁷

The share of skilled labor needed for entry for each quality level, α_{q_n} , is chosen to match the share of skilled labor used in the production of that quality. Therefore, high quality producers use a more skill intensive production process (lower θ_{q_n}) and also have more skill intensive entry requirement.²⁸

Finally, v^2 , the variance of the log of the productivity draw (common across qualities), is chosen to match the standard deviation of the log of employment in the combined ASI and SUM dataset which was 0.64.

²⁶I run a regression of log wages on a dummy of whether the individual is skilled (at least ten years of education) for all manufacturing workers, controlling for potential experience, sex, state, industry, occupation, and whether the individual is residing in a rural or urban area. Individuals with ten or more years of education on average make 56.8 percent more than workers with less than ten years of education which is rounded up to a wage premium of 1.6. Appendix D reports more details and the regression results.

²⁷In particular, the μ'_{q_n} s are chosen to match a price-size slope of 0.1. Note that plants of each higher quality level are calibrated to be two times the size of the previous quality level. Therefore, the μ'_{q_n} s are chosen such that each higher quality level charges a log price which is $0.1 * \log(2)$ higher than the previous quality levels log price.

²⁸The ratio of skilled to unskilled labor used by plants of quality q_n is given by $\frac{l^u_{i,q_n}}{l^s_{i,q_n}} = \left(\frac{w_S}{w_U} \frac{\theta_q}{1-\theta_q} \right)^{\sigma_{us}}$ and is independent of the productivity draw of the plant. Therefore, the share of skilled workers used in production of quality q_n is simply $\frac{1}{1 + \left(\frac{w_S}{w_U} \frac{\theta_q}{1-\theta_q} \right)^{\sigma_{us}}}$.

Table 7: Calibration

Param.	Description	Targets
f_{q_n}	Fixed costs	$size_{q_n} = \{1.25, 2.5, 5, \dots, 2560\}$
μ_{q_n}	Mean of productivity draws	Price-size slope of 0.1
θ_{q_n}	Sh of U in production	$w_S = 1.6$; and $\frac{L_{q_n}^U}{L_{q_n}^S}$ across qualities
α_{q_n}	Sh of skilled in entry	$\frac{L_{q_n}^S}{L_{q_n}^S + L_{q_n}^U}$ in production
v^2	Variance of productivity draw	Std dev of employment = 0.64
$q_n (\Delta)$	Utility from quality	Price-income slope of 0.1
a_{q_n}	Constant in utility function	Size distribution

4.2. Utility Parameters

The utility function in the model takes the form

$$u_{j,q_n}(c_{j,q_n}, \varepsilon_{j,q_n}) = a_{q_n} + q_n \log(c_{j,q_n}) + \varepsilon_{j,q_n} \quad \forall q_n \in \mathcal{Q}. \quad (10)$$

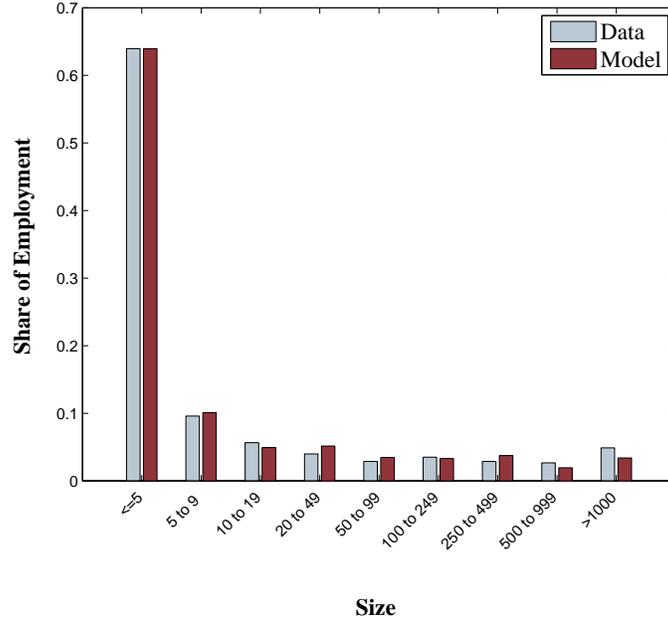
Two sets of parameters need to be calibrated: (1) q_n , the quality indexes; and (2) a_{q_n} , the quality specific constant in the utility function.

As mentioned in Section 3.1, the quality indexes are parametrized as follows: $q_1 = 1$ and $q_n = q_{n-1} + \Delta$.²⁹ The value of Δ determines the steepness of the quality Engel curve i.e. how quickly does demand move to higher quality as income levels increase. In the model, skilled workers earn wage w_S (which is calibrated to be 1.6) and unskilled workers earn wage w_U (which is normalized to one as the numeraire). Δ is chosen to match the price-income relation documented in Table 1 of Section 2.1. In particular, Δ is chosen such that the price-income elasticity in the model is 0.1 i.e. the average log price paid by skilled households is $0.1 * \log\left(\frac{w_S}{w_U}\right)$ more than for unskilled households. As higher quality producers in the model have higher prices, this in effect determines the extent to which demand shifts towards high quality as we move from unskilled wages to skilled wages.

The quality specific constant in the utility function, a_{q_n} , determines the absolute levels of demand for different quality levels i.e. it determines $\rho(q_n|w)$ given in equation (3). A higher a_{q_n} for a specific quality means that a larger share of households are likely to buy that quality (irrespective of income level). There-

²⁹Setting q_1 to be one is not a normalization in the model. However, for the baseline specification with no love of variety, the results are not sensitive to the choice of q_1 . This issue is discussed further in Section 5.3.

Figure 6: Size Distribution - Data vs Model



Notes: The figure plots the share of employment in different size categories in the data and in the calibrated baseline of the model. The data is for the manufacturing sector in India for 2005-06. It combines the ASI and the SUM (same as Figure 1).

fore, I choose a_{q_n} such that the size distribution in the model matches the size distribution for India as a whole in 2005-06.

In summary, a_{q_n} pins down the absolute level of demand for the different qualities and are calibrated to match the size distribution in the model to the Indian data. Δ determines the differences in demand for high versus low quality levels between skilled and unskilled workers and is calibrated to match the price-income elasticity seen in the data.

Table 7 summarizes the calibration. Figure 6 plots the share of workers in plants of different size categories for the calibrated model and the data (combining the ASI and the SUM for 2005-06). As the model parameters were chosen to match the size distribution, it is not surprising to see that the size distribution in the model matches the data very closely. However, the model was not calibrated to match the change in size distribution as income levels change. The extent to which the size distribution changes in the model as income levels change depends crucially on the degree of non-homotheticity (Δ) on the consumer side and the price-size relation on the producer side and these parameters were calibrated using micro-data from consumer and producer surveys.

5. Results

Having calibrated the model, I now conduct counterfactual exercises in which I simulate differences in per-capita income levels in the model and see how this effects the size distribution. In addition to the counterfactual exercises, the sensitivity of the results to some important parameters is also explored.

5.1. Cross-section of Indian States

I now ask the question: How much of the cross-state differences in the size distribution seen in the data can be explained by the model if per-capita income in the model varies by the same amount as it varies across Indian states? To do this I conduct counterfactual exercises in which I vary three sets of parameters in the model while keeping all the other parameters unchanged:

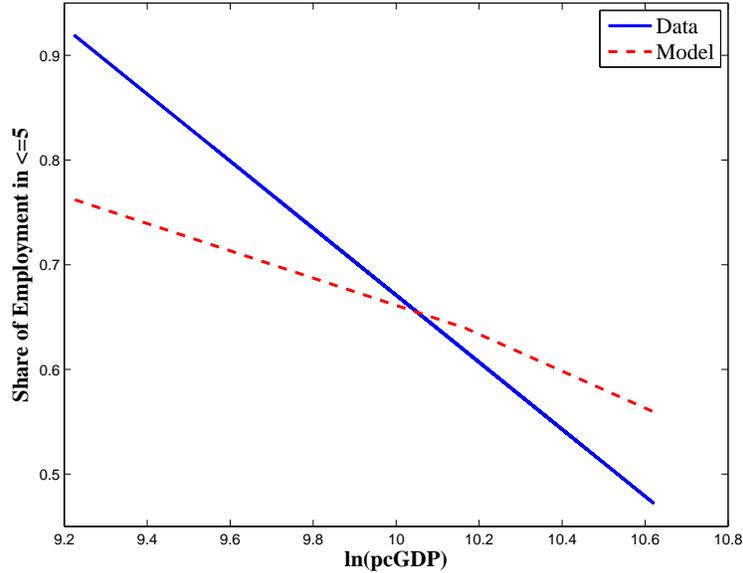
1. The share of the households in the model who are skilled, h , is varied in the counterfactual exercises to match the share of workers with ten or more years of education across rich and poor states. About 13 percent of the manufacturing workers in the poorest states are skilled as compared to 43 percent in the richest states.
2. The share parameter of unskilled labor for intermediate producers, θ_{q_n} , is changed across the counterfactuals to keep the wage premia unchanged.³⁰ This can be viewed as skill biased technical change with richer states having a higher supply of skilled labor and also using skilled labor more intensively in the production of all quality levels.³¹
3. The mean of the productivity draw of intermediate producers, μ_{q_n} , is changed to match the differences in per-capita income across states and to maintain the price-size slope of 0.1 across the counterfactuals.³² Per-capita income of the poorest Indian state (Bihar) is 0.39 times India's per-capita income

³⁰If I do not change θ_{q_n} , then the wage premia falls in the counterfactual for the richer states due to the higher supply of skilled workers. However, in the data, wage premia does not vary systematically across states. In particular, if I run a Mincerian regression of log of wages on a dummy which takes value 1 if the person is skilled and also include the interaction of the dummy with per-capita state NDP (controlling for industry, occupation, sex, experience etc), then the coefficient on the interaction is not significantly different from zero.

³¹In effect, θ_{q_1} for each counterfactual is chosen to maintain the wage premia ($w_S = 1.6$). All the other θ'_{q_n} s are picked as described in equation (9) to match the ratio of skilled to unskilled in different quality levels relative to the worst quality level. Furthermore, in the counterfactual, the share of entry labor which needs to be skilled workers (α_{q_n}) is also changed to match the share of skill in production for each quality level i.e. the richer states do not just use more skill intensive production techniques but also use more skill in the entry process.

³²As mentioned in Section 4.1, μ_{q_n} for each quality level was chosen to match the price-size elasticity of 0.1. In the counterfactual exercises, as the θ'_{q_n} s are changed, this can lead to changes in prices of the high quality relative to low quality even though there is no change in wage premia. These changes in relative prices can cause a shift in demand and thus changes in the size distribution for reasons other than changes in real income which is what I want to focus on. Therefore, in the counterfactual, in addition to scaling

Figure 7: Counterfactual Across Indian States - Data vs Model



Notes: The figure plots the share of employment in plants of size five or less across Indian states in the data and for the counterfactual exercise in the model. The blue line is the linear regression line of share of employment in plants of size five or less in different Indian states on log of per-capita GDP of the state. The red line is the model predicted share of employment in plants of size five or less when conducting the counterfactual exercise.

while that of the richest state (Maharashtra) is 1.57 times India's per-capita income. To generate similar differences in per-capita income in the model, the poorer states in the counterfactual exercise have lower average productivity levels compared to the richer states.³³

To summarize, three sets of parameters are changed in the counterfactual exercises: the share of skilled in the population, the skill intensity of the production process, and the means of the productivity draws of intermediates. These parameters are changed to match the differences in skill composition and per-capita income levels across Indian states while keeping the wage premia and the relative prices of different quality levels unchanged.³⁴

An increase in the productivity of intermediate producers and in the supply of skill translates into an increase in real income levels in the model. The increase in real income level leads to demand shifting

all the $\mu'_{q,s}$ by a constant (to match the differences in per-capita income seen across Indian states), I also change the relative $\mu'_{q,s}$ of different qualities to maintain the same relative prices of different quality levels. This eliminates any substitution effects due to relative price changes and only focuses on changes in demand (caused by the non-homotheticity in the preferences) due to changes in per-capita income levels.

³³In order to define per-capita GDP in the model, I need to define a set of base prices. I use the prices of intermediates in the calibrated baseline as the base prices and value output in the counterfactuals using these prices.

³⁴Much of the variation in income levels across states in the model is captured by differences in productivity. This is consistent with development accounting exercises which also find that residual TFP explains the majority of the differences in per-capita income across Indian states (see Chanda (2011)).

towards higher quality goods due to the non-homotheticity in the preferences. This change in demand leads to a shift in the production side. The number of plants producing low quality goods declines while those producing high quality increases. This in turn implies that there is a shift in the size distribution with the share of employment in small plants falling.

The red dashed line in Figure 7 plots the share of employment in plants of size five or less that is predicted by the model when conducting the counterfactual exercises. In the calibrated baseline, the share of employment in small plants in the model is 63.9 percent. When productivity and supply of skill is lowered such that per-capita income levels decrease by a factor of 0.39 (0.94 log points lower), the share of employment in plants of size five or less increases to 75.6 percent. On the other hand, when productivity and supply of skill is increased such that per-capita income levels increase by a factor of 1.57 (0.43 log points higher) compared to the calibrated baseline, the share of employment in small plants falls to 56.3 percent.

The solid blue line in Figure 7 plots the projection from a linear regression of the share of employment in plants of size five or less on log of per-capita State NDP across Indian states. The share of employment in small plants is computed by combining the ASI and the SUM (the same data as in Figure 2). In the data, the poorest Indian states have about 91.9 percent of employment in small plants while the richest have 47.2 percent employment in small plants.

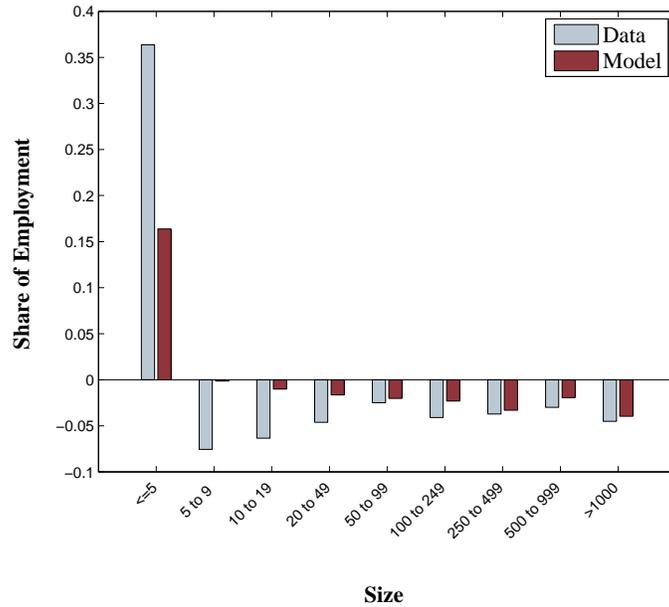
While the share of employment in small plants varies by 44.7 percentage points across Indian states in the data, the model predicts an 19.3 percentage points difference. Therefore, the model explains about 43 percent of the difference in share of employment in small plants seen across Indian states.

Figure 8 compares how the entire size distribution (as opposed to just the share of employment in plants of size five or less) changes in the model as compared to the data as we change income levels. In the data, I pool together the three poorest states and the three richest states and compute the share of employment in different size categories for these groups of states.³⁵

The light blue bars in Figure 8 show the difference (in percentage points) in the share of employment in the three poorest states compared to the three richest states for each size category. The poorest states have about 36 percentage points more employment in plants of size five or less as compared to the richest states. The richer states have a larger share of their employment in all the larger size categories as compared to the poor states, which is why the blue bars lie below zero for all these size categories. The red bars represent the same difference in share of employment for different size categories that the model predicts when productivity and skill levels in the model are varied to match the incomes differences across these groups of states. The model predicts that the share of employment in plants of size five or less is about 15

³⁵I pool the three richest and poorest states in order to avoid having the results being driven by an outlier state. The results are similar if I just compare the richest state to the poorest state.

Figure 8: Counterfactual: Changes in Distribution for 3 Richest vs 3 Poorest States



Notes: The figure plots the share of employment in the three poorest states minus the share in the three richest states for different size categories in the data and in the model (when productivity and skill levels are varied to match the differences in per-capita income across these groups of states). The data is from the ASI and SUM for 2005-06.

percentage points higher in the poorer states as compared to richer states, which again accounts for about 42 percent of the difference seen in the data. Again, like the data, the red bars lie below zero for all the other size categories, indicating that the model predicts a larger share of employment in richer states for these size categories.

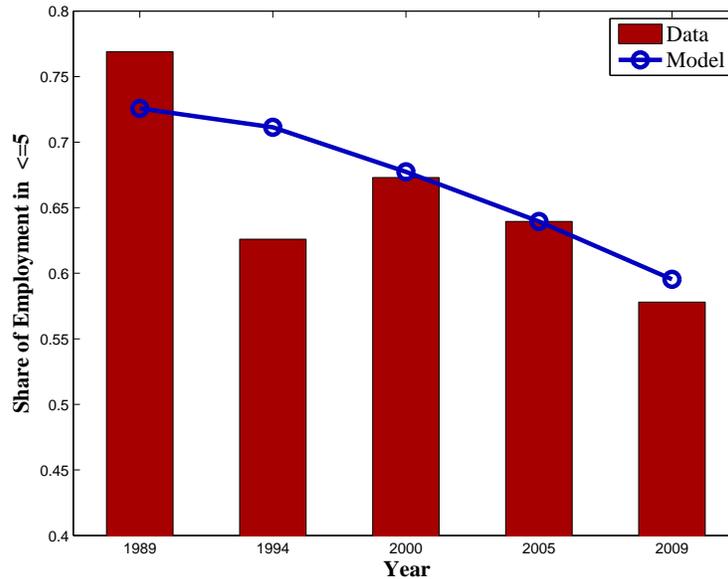
5.2. India Over Time

I now look at how well the model does in explaining the evolution of the size distribution of manufacturing plants in India over time. Five waves of the Survey of Unorganized Manufacturing (SUM) have been conducted in Indian between 1989-90 and 2010-11. These can be combined with the corresponding years of the Annual Survey of Industries (ASI) to get five data points for how the size distribution has evolved over time in India.

The bars in Figure 9 show the share of employment in plants of size five or smaller for 1989, 1994, 2000, 2005, and 2009.³⁶ As can be seen, the share of employment in small plants has decreased from 77 percent

³⁶More details of the surveys are given in Appendix A.1 and A.2. One point to note is that the last wave of the SUM was done in 2010-11. However, the last wave of the ASI that was available at the time of writing was 2009-10. Hence the 2009-10 ASI is pooled with the 2010-11 SUM.

Figure 9: Counterfactual India Over Time - Data vs Model



Notes: The red bars in the figure plot the share of employment in plants of size five or less for five years for India. The data for each year pools the SUM and the the ASI for that year. The blue line plots the model predicted share of employment for each year when productivity and skill levels are varied to match the differences in per-capita income in India over time.

of total employment in 1989 to 58 percent in 2009.

Per-capita income in 1989 was 0.54 times the 2005 level of per-capita income while the share of manufacturing workers with ten or more years of schooling was just 14 percent. In 2009 per-capita income levels were 1.30 times the 2005 level while the share of manufacturing workers with ten or more years of schooling had increased to 31 percent. The blue line in Figure 9 plots the share of employment in plants of size five or less as predicted by the model when productivity and skill supply in the model is varied to the extent required to match the differences in per-capita income levels and share of skilled in the data. The model was calibrated to match the share of employment in small plants in 2005, therefore, the fit in 2005 is very good by construction. The model predicts that 72 percent of employment would be in plants of size five or less in 1989, which is a little less than the 77 percent seen in the data. Similarly, the model under-predicts the change in the size distribution going from 2005 to 2009 by a small amount. Overall, the model predicts 65 percent of the change in share of employment in small plants seen in the data between 1989 to 2009.³⁷

³⁷A related question can also be asked: Did states which grow more over time see a larger drop in share of employment in small plants? If we look at the change in the size distribution between 1989 and 2009, then this does seem to be the case. However, this relation is not robust to looking at the 2000 to 2009 period only. See Figure A.1 in the appendix.

Table 8: Love of Variety: Percent of Cross-State Difference Explained

	$q_1 = 1$	$q_1 = 0.1$
$\eta = \frac{1}{\sigma-1}$	43.1%	43.1%
$\eta = 0$	71.2%	53.1%

Notes: The table shows the percent of cross-state variation in share of employment in plants of size five or less that is explained by the model counterfactual for different parameter values of η and q_1 . $\eta = \frac{1}{\sigma-1}$ is the baseline specification of no love of variety while $\eta = 0$ is the case of full love of variety.

5.3. Parameter Sensitivity: Love of Variety

As mentioned in Section 3.2, the baseline specification of the model assumed that the final goods producers production function had no love of variety. A generalization of the the production function of the final goods producer of quality q_n is given by

$$Y_{q_n}^s = \frac{1}{M_{q_n}^\eta} \left(\sum_{i=1}^{M_{q_n}} x_{i,q_n}^{\frac{\sigma-1}{\sigma}} \right)^{\frac{\sigma}{\sigma-1}} \quad \forall q \in Q.$$

In the baseline specification, η was set equal to $\frac{1}{\sigma-1}$, which corresponded to the case of no love of variety. In this section, I provide results for the case when $\eta = 0$ (the case with full love of variety) and compare this to the baseline. As mentioned in Section 3.2, the no love of variety assumption is the conservative case, with changes in the size distribution in the counterfactual being larger when we allow for love of variety. Furthermore, when allowing for love of variety, the results become more sensitive to the choice of q_1 , the quality index of the lowest quality level (note that given q_1 , all subsequent quality indexes are given by the recursion $q_n = q_{n-1} + \Delta$).

Table 8 shows how much of the cross-state differences in share of employment in small plants is explained by the model for different values of η and q_1 . The first row and first column corresponds to the baseline specification, with $\eta = \frac{1}{\sigma-1}$ (no love of variety) and $q_1 = 1$. As mentioned in Section 5.1, when varying productivity and supply of skill to match the differences in per-capita incomes across states, the model explains 43.1 percent of the difference in the share of employment in small plants as compared to the data.

Now consider the model with love of variety ($\eta = 0$). When allowing for love of variety, all other parameters are recalibrate to match the same moments as in the baseline. I then run the same counterfactual exercises as in Section 5.1. As reported in Table 8, in the case with love of variety, the model can explain 71.2 percent of the differences in size distribution between the rich and poor states.

Why is it that in the case with love of variety, the model generates bigger changes in the size distribution

in the counterfactual? The reason is that in the case with love of variety, relative prices of different quality levels change in the counterfactual, due to changes in the relative varieties of the different qualities. In particular, the CES price index (the price charged by the final producer to the consumer) for quality q_n is given by

$$P_{q_n} = M_{q_n}^{\eta - \frac{1}{\sigma-1}} \left(\int (p(A_i, q_n))^{1-\sigma} g_{q_n}(A_i) dA_{q_i} \right)^{\frac{1}{1-\sigma}} \quad \forall q \in \mathcal{Q}.$$

In the baseline specification, because $\eta = \frac{1}{\sigma-1}$, the price index for q_n was independent of the number of varieties M_{q_n} . However, when $\eta = 0$, the CES price index of a quality level, P_{q_n} , is inversely related to the number of varieties of that quality (M_{q_n}) available in the economy. In the counterfactual, as income levels increase, demand shifts towards higher quality, and this induces more entrants of the higher quality levels. The increase in number of varieties of high quality intermediate producers causes the relative price of high quality goods to fall in the counterfactual when $\eta = 0$. This causes a further shift in demand towards high quality which in turn causes more entry into higher quality goods. The additional increase in demand for high quality which acts through relative price changes due to change in number of varieties does not occur in the baseline specification when $\eta = \frac{1}{\sigma-1}$. Hence, the change in size distribution in the counterfactual in the baseline specification is less than in the case with love of variety. In effect, the baseline specification focuses attention on the changes in demand caused by changes in income levels alone. It abstracts away from any changes in relative prices caused by changes in number of varieties in the counterfactual.

Furthermore, when allowing for love of variety, the change in the size distribution in the counterfactual, becomes more sensitive to the choice of q_1 , the quality index for the lowest quality level. When q_1 is set to 0.1 and $\eta = 0$, the model counterfactual explains only 53.1 percent of the difference in size distribution as opposed to 71.2 percent when $q_1 = 1$. As shown in Section 3.1, the share of households with wage w who choose quality level q_n is given by

$$\rho(q_n|w) = \frac{e^{a_{q_n}} \left(\frac{w}{P_{q_n}} \right)^{q_n}}{\sum_{i=1}^N e^{a_{q_i}} \left(\frac{w}{P_{q_i}} \right)^{q_i}} \quad \forall q_n \in \mathcal{Q}.$$

As P_{q_n} is raised to the power q_n in the numerator, the absolute levels of q_n approximately determine the own price elasticity of demand for a quality level. Lower absolute levels of the quality indexes imply that demand is less sensitive to changes in relative prices (of the CES price indexes). Therefore, a lower value for q_1 (which translates into lower values for all the quality indexes) makes the model less sensitive to the changes in relative prices induced by changes in varieties.

5.4. India vs US

I conduct a counterfactual in which I simulate an economy with per-capita income level equivalent to that of the US in 2005 (seventeen times that of India) and see how the size distribution in the counterfactual compares to that of the calibrated baseline. I vary the levels of productivity, supply of skill, and θ'_s in the model to match per-capita GDP and supply of skill in the US while keeping the wage premium and relative prices unchanged. The share of employment in plants which employ 5 or less people falls from 64 percent in the calibrated baseline to 13 percent in the counterfactual.

It is important to note an important caveat while interpreting this cross-country result. The calibration of the model was local to India's level of development and therefore a simple extrapolation to the US might be potentially biased. The price-size relation on the producer side and the price-income relation on the consumer side are similar across states of different income levels within India but might be very different for the US. For example, in the US it is possible that the producer side relation between price and size is flatter (or even negative) as the lower quality goods might be produced in large factories while the higher quality goods might be produced in small boutique organizations. Similarly, the elasticity of substitution between skilled and unskilled labor might be very different between the two countries. Therefore, although it is an interesting exercise to do the counterfactual for the US, the results should be interpreted with a lot more caution than the cross-state counterfactual.

6. Inter-State Trade

The model presented above implicitly assumed that each state in India can be treated as a closed economy and that differences in income levels across states translate into differences in demand and in the size distribution at the state level. How would the possibility of inter-state trade affect the hypothesis presented in the paper?

A potential confounding effect of inter-state trade could come through the location choice of large plants. For example, if the richer states are more suited for operating large plants (due to availability of skilled labor, better labor laws etc), then all the larger plants might choose to locate in these states and ship their goods to the poor states. In this case, the fact that richer states have a smaller share of employment in small plants would not reflect differences in demand across states but rather just the spatial location choice of large plants.

To address this concern, it would be ideal to have a measure of inter-state trade flows (similar to the Commodity Flow Survey in the US) to see how important this channel could be. Unfortunately, data on

extent of inter-state trade is not collected in India. Here I provide indirect evidence to suggest that inter-state trade is not completely driving the cross-state relation seen in Figure 2.

Firstly, transportation costs in developing countries are often very high which makes it harder for plants to transport goods over large distances to poorer states. Atkin and Donaldson (2012) show that intranational transportation costs in two African countries are seven to fifteen times larger than similar estimates for the US. Furthermore, Hillberry and Hummels (2008) show that even in the US, manufacturing production is extremely localized with local shipments volumes being three times larger than shipments to more distant locations. This suggests that local demand is likely to be an important determinant of the the size distribution in any region, especially in developing countries.

Furthermore, if inter-state trade is driving the cross-state relation seen in Figure 2, then we would expect more tradable industries to exhibit larger differences in share of employment in small plants across states as compared to less tradable industries. On the other hand, if the states are in fact approximated well as closed economies then we would expect the relation between share of employment in small plants and per-capita NDP to be stronger for non-tradables. To test this fact, I construct two measures of tradability (within manufacturing) at the 3-digit level of the National Industrial Classification (NIC) of 2004.³⁸ These are:

1. Herfindahl index of geographical concentration in the US: The County Business Patterns Database of 2005 released by the United States Census Bureau provides information regarding the number of people working in each 6-digit industry of the North American Industry Classification System (NAICS) for each county in the US.³⁹ As the tradability index is to be applied to the Indian industry classification, I first create a concordance from 6-digit NAICS to 3-digit NIC and then construct a Herfindahl Index (H-index) of geographical concentration of each 3-digit NIC across US counties.⁴⁰ The H-index is defined as

$$H_i = \sum_{c=1}^C (sh_{i,c}^L)^2,$$

where ‘*i*’ indexes industry (according to NIC), ‘*c*’ indexed counties, and $sh_{i,c}^L$ represents the share of industry ‘*i*’ employment which is in county ‘*c*’. The H-index for industry ‘*i*’ is simply the sum across counties of the square of the share of the industries employment which is present in county ‘*c*’. The

³⁸Economic activity in India is classified according to the National Industrial Classification (NIC) which closely follows the United Nation’s International Standard Industrial Classification (ISIC). Details regarding different NIC revisions and the concordance used between them are given in Appendix B.3.

³⁹The data can be found at <http://www.census.gov/econ/cbp/>. The exact number of people in many industry-county cells is masked. Instead, the dataset reports an employment size class for that cell. In these cases, the employment in the cell is assigned the midpoint of the size class reported.

⁴⁰The concordance from 6-digit NAICS and 3-digit ISIC Rev 3.1 was based on the Census Bureau’s concordance file available at <http://www.census.gov/eos/www/naics/concordances/concordances.html>. ISIC Rev 3.1 to NIC 2004 is a one to one correspondence at the 3-digit level. Appendix B.1 has more details regarding the concordance.

industries which are highly concentrated in a few counties in the US (have a high value for Herfindahl index) are considered to be tradable industries while industries which have employment spread over lots of counties (have a low value for the Herfindahl index) are considered non-tradable industries. This measure for tradability of an industry based on US levels of concentration is applied to India.

2. Degree of international trade in India: For each 3-digit NIC in the manufacturing sector, I construct a measure of the degree of international trade carried out in the industry as a share of domestic production. In particular, I define this measure of international trade as the exports plus imports in that industry as a share of gross production of that industry carried out by domestic plants in 2005-06. The data for exports and imports for India is taken from the website of the Department of Commerce, Government of India.⁴¹ The imports and exports data is not at the industry level but rather classified according to the Harmonized Commodity Description and Coding System (HS) product classification. This is converted to 3-digit NIC using the products to industry concordance developed by World Integrated Trade Solutions (WITS).⁴² The data on gross domestic production for each industry is computed by combining the ASI and the SUM. Industries in which international trade is a large percent of domestic production are considered to be more tradable.

Table A.5 in the appendix lists the 3-digit industries which lie above and below the median value of the two indexes of tradability. The two measures of tradability are weakly positively correlated with the rank correlation coefficient between them being 0.25.

I run regressions of the form

$$sd_{i,s,t} = \alpha_{i,t} + \alpha_{s,t} + \gamma \ln(SNDP_{s,t}) * tradability_i + \varepsilon_{i,s,t} \quad (11)$$

where $sd_{i,s,t}$ is the share of employment in plants of size five or less in industry 'i' in state 's' at time 't', $SNDP_{s,t}$ is the per-capita NDP of state 's' at time 't', and $tradability_i$ is a dummy variable which takes value 1 if an industry is classified as tradable. $\alpha_{i,t}$ represents fixed effects for industry interacted with time and it controls for the fact that different industries might have different average levels for the share of employment in small plants. $\alpha_{s,t}$ represents fixed effects for state interacted with time and controls for the fact that rich states on average have a lower share of employment in small plants.

The coefficient of interest is γ , the coefficient on the interaction of state per-capita income and the tradability dummy. A positive γ implies that the relation between the share of employment in small plants and

⁴¹The data is available from <http://commerce.nic.in/eidb/default.asp>.

⁴²WITS is based on a collaboration of the World Bank with UNCTAD, WTO and other international organization associated with international trade data. More details regarding the concordance can be found in Appendix B.2.

Table 9: Size Income Relation Across States for Tradables vs. Non-tradables

Dependent Variable: share of employment in <=5 in industry 'i', state 's', time 't'				
	(1)	(2)	(3)	(4)
log(per-capita SNDP)*tradability	0.068* (0.0351)	0.052 (0.0394)	-0.010 (0.0469)	0.000 (0.0498)
Index	H-index	H-index	Exp-Imp	Exp-Imp
Cutoff	Median	Quartile	Median	Quartile
Observations	3,885	1,826	3,899	1,959

Notes: The data is from five rounds of the ASI and SUM. The table reports regression results for the share of employment in plants of size 5 or less in industry 'i' in state 's' at time 't' on log per-capita state NDP interacted with a dummy which takes value 1 if industry 'i' is classified as a tradable industry. Column 1 classifies an industry as tradable if the Herfindahl Index across US counties for the industry was above the median of Herfindahl Indexes, and non-tradable if it was below the median. Column 2 uses top and bottom quartiles of the Herfindahl Index as cutoffs. Column 3 and 4 use the tradability index based on Indian exports and imports and uses the median and the top and bottom quartiles as cutoffs respectively. All regressions include fixed effects for industry interacted with time and state interacted with time. Each observation is weighted by the share of observations in the state-industry cell out of the total observations in the ASI and SUM combined for the given year. Standard errors are clustered at the state level. *p<0.1.

log of per-capita income across states is stronger for non-tradables. This is because the share of employment in small plants and per-capita NDP are negatively related and therefore a positive interaction term implies that the slope for tradable industries is less negative compared to non-tradables. Therefore, a positive value of γ is supportive of the view that inter-state trade is not a major driving force behind the size distribution of plants across states.

An industry is classified as tradable if the tradability index for the industry lies above the median (or in the top quartile) of the index across industries. Data for five waves of the SUM is combined with the corresponding year of the ASI (1989, 1994, 2000, 2005, and 2010). Only the fifteen large Indian states mentioned in footnote 44 are included as the smaller states often have no observations for many industries as the 3-digit level.

Table 9 reports results for equation (11) for both the measures of tradability. Each observation is weighted by the share of observations in the state-industry cell out of the total observations in the ASI and SUM combined for the given year.⁴³ Column 1 uses the Herfindahl index and classifies an industry as tradable if its Herfindahl Index is above the median value of the Herfindahl Index across industries. The coefficient on the interaction of per-capita NDP and the tradability index is positive and marginally significant at the 10

⁴³This weighting scheme accounts for the fact that the size distribution variable (dependent variable) for some industry-state pairs is based on a lot fewer observations than other cells, and are therefore likely to be measured with less precision. Table A.6 in the appendix reports results when all observations are weighted equally. Table A.7 reports results when industrial categories which are residual categories (industry categories with descriptions which include words like "not elsewhere covered" or "others") are excluded.

percent level. Column 2 classifies an industry as tradable if it is in the top quartile in terms of the Herfindahl Index and non-tradable if it is in the bottom quartile. The results are very similar to the first column. Columns 3 and 4 use the median and quartile of the tradability measure based on exports and imports in India. The point estimates of the coefficient on the interaction of per-capita NDP and the tradability index is much smaller in absolute value and statistically insignificant.

The results in Table 9 suggest that the size-income relation across states is not stronger for tradable industries as compared to non-tradable industries.

7. Conclusion

The size distribution in developing countries usually has a thick left tail compared to developed countries. The same holds across Indian states, with richer states usually having a much smaller share of their manufacturing employment in small plants. In this paper, I explore the hypothesis that this income-size relation arises from the fact that low income countries and states have high demand for low quality products which can be produced efficiently in small plants. I provide evidence which is consistent with this hypothesis from both the consumer and producer side. In particular I show that richer households buy higher price goods while larger plants produce higher price products (and use higher price inputs). Finally, a model is developed which features non-homothetic preferences with respect to quality and is calibrated to match the cross-sectional facts from the consumer and producer sides. A calibrated version of the model indicates that up to 41 percent of the cross-state variation seen in the left tail of manufacturing plants in India can be explained by the model.

Therefore, this paper suggests that a large part of the differences in size distribution that we see across countries and states is a natural consequence of the low levels of income in developing countries and is not caused by policies which discriminate against large productive plants in favor of small unproductive plants. The presence of small plants in developing countries should not be viewed as originating necessarily from policy failures.

References

- Aguiar, Mark, and Erik Hurst, 2007, "Life-Cycle Prices and Production," *American Economic Review*, Vol. 97, No. 5, pp. 1533–1559.
- Alfaro, Laura, Andrew Charlton, and Fabio Kanczuk, 2009, "Plant-Size Distribution and Cross-Country Income Differences," in *NBER International Seminar on Macroeconomics 2008*, NBER Chapters, pp. 243–272 (National Bureau of Economic Research, Inc).
- Atkin, David, and Dave Donaldson, 2012, "Who's Getting Globalized? The Size and Nature of Intranational Trade Costs," Techn. rep., Yale University.
- Attanasio, Orazio P., and Christine Frayne, 2006, "Do the Poor Pay More?" Techn. rep.
- Banerjee, A.V., and E. Duflo, 2011, *Poor Economics: A Radical Rethinking of the Way to Fight Global Poverty* (PublicAffairs).
- Banerji, Arup, and Sanjay Jain, 2007, "Quality dualism," *Journal of Development Economics*, Vol. 84, No. 1, pp. 234–250.
- Behar, Alberto, 2009, "Directed technical change, the elasticity of substitution and wage inequality in developing countries," Economics Series Working Papers 467, University of Oxford, Department of Economics.
- Bils, Mark, and Peter J. Klenow, 2001, "Quantifying Quality Growth," *American Economic Review*, Vol. 91, No. 4, pp. 1006–1030.
- Bloom, Nicholas, Raffaella Sadun, and John Van Reenen, 2012, "The Organization of Firms Across Countries," *The Quarterly Journal of Economics*, Vol. 127, No. 4, pp. 1663–1705.
- Broda, Christian, and David E. Weinstein, 2006, "Globalization and the Gains from Variety," *The Quarterly Journal of Economics*, Vol. 121, No. 2, pp. 541–585.
- Chanda, Areendam, 2011, "Accounting for Bihar's Productivity Relative to India's: What can we learn from recent developments in Growth Theory," Techn. Rep. 11/0759, International Growth Centre.
- Choi, Yo Chul, David Hummels, and Chong Xiang, 2009, "Explaining import quality: The role of the income distribution," *Journal of International Economics*, Vol. 77, No. 2, pp. 265–275.
- Dalgin, Muhammed, Devashish Mitra, and Vitor Trindade, 2008, "Inequality, Nonhomothetic Preferences, and Trade: A Gravity Approach," *Southern Economic Journal*, Vol. 74, No. 3, pp. 747–774.
- De Soto, Hernando, 1989, *The other path* (Harper & Row New York).
- Deaton, Angus, and Olivier Dupriez, 2011, "Spatial price differences within large countries," Working Papers 1321, Princeton University, Woodrow Wilson School of Public and International Affairs, Research Program in Development Studies.
- DiCecio, Riccardo, and Levon Barseghyan, 2010, "Entry Costs, Industry Structure, and Cross-Country Income and TFP Differences," 2010 Meeting Papers 964, Society for Economic Dynamics.
- Dikhanov, Yuri, 2010, "Income Effect and Urban-Rural Price Differentials from the Household Survey Perspective," Techn. rep., ICP Global Office.
- Djankov, Simeon, Rafael La Porta, Florencio Lopez-De-Silanes, and Andrei Shleifer, 2002, "The Regulation Of Entry," *The Quarterly Journal of Economics*, Vol. 117, No. 1, pp. 1–37.

- Faber, Benjamin, 2012, "Trade Liberalization, the Price of Quality, and Inequality: Evidence from Mexican Store Prices," Techn. rep., London School of Economics.
- Fajgelbaum, Pablo, Gene M. Grossman, and Elhanan Helpman, 2011, "Income Distribution, Product Quality, and International Trade," *Journal of Political Economy*, Vol. 119, No. 4, pp. 721 – 765.
- Flam, Harry, and Elhanan Helpman, 1987, "Vertical Product Differentiation and North-South Trade," *American Economic Review*, Vol. 77, No. 5, pp. 810–22.
- García-Santana, Manuel, and Josep Pijoan-Mas, 2010, "Small Scale Reservation Laws And The Misallocation Of Talent," Working papers, CEMFI.
- Garicano, Luis, Claire LeLarge, and John Van Reenen, 2013, "Firm Size Distortions and the Productivity Distribution: Evidence from France," Working Paper 18841, National Bureau of Economic Research.
- Ghani, Ejaz, Arti Grover Goswami, and William R. Kerr, 2012, "Is India's manufacturing sector moving away from cities?" Policy Research Working Paper Series 6271, The World Bank.
- Gollin, Douglas, 1995, "Do Taxes on Large Firms Impede Growth? Evidence from Ghana," Bulletins 7488, University of Minnesota, Economic Development Center.
- Guner, Nezih, Gustavo Ventura, and Xu Yi, 2008, "Macroeconomic Implications of Size-Dependent Policies," *Review of Economic Dynamics*, Vol. 11, No. 4, pp. 721–744.
- Hallak, Juan Carlos, 2006, "Product quality and the direction of trade," *Journal of International Economics*, Vol. 68, No. 1, pp. 238–265.
- Hallak, Juan Carlos, and Jagadeesh Sivadasan, 2011, "Firms' Exporting Behavior under Quality Constraints," Working Papers 628, Research Seminar in International Economics, University of Michigan.
- Hasan, Rana, and Karl Robert L Jandoc, 2010, "The Distribution of Firm Size in India: What Can Survey Data Tell Us?" Techn. rep., Citeseer.
- Hillberry, Russell, and David Hummels, 2008, "Trade responses to geographic frictions: A decomposition using micro-data," *European Economic Review*, Vol. 52, No. 3, pp. 527–550.
- Hsieh, Chang-Tai, and Peter J. Klenow, 2012, "The Life Cycle of Plants in India and Mexico," Working Paper 18133, National Bureau of Economic Research.
- Hsieh, Chang-Tai, and Benjamin A. Olken, 2014, "The Missing & Missing Middle," *Journal of Economic Perspectives*, Vol. 28, No. 3, pp. 89–108.
- Hummels, David, and Peter J. Klenow, 2005, "The Variety and Quality of a Nation's Exports," *American Economic Review*, Vol. 95, No. 3, pp. 704–723.
- Iacovone, Leonardo, and Beata Javorcik, 2012, "Getting Ready: Preparation for Exporting," CEPR Discussion Papers 8926, C.E.P.R. Discussion Papers.
- Kugler, Maurice, and Eric Verhoogen, 2012, "Prices, Plant Size, and Product Quality," *Review of Economic Studies*, Vol. 79, No. 1, pp. 307–339.
- La Porta, Rafael, and Andrei Shleifer, 2008, "The Unofficial Economy and Economic Development," NBER Working Papers 14520, National Bureau of Economic Research, Inc.
- Little, Ian, Dipak Mazumdar, and John M. Page Jr, 1987, *Small Manufacturing Enterprises: A Comparative Analysis of India and Other Economies* (NY: Oxford U. Press).

- Loayza, Norman V., 1996, "The Economics of the Informal Sector: A Simple Model and Some Empirical Evidence from Latin America," *Carnegie-Rochester Conference Series on Public Policy*, Vol. 45, No. 0, pp. 129 – 162.
- Loayza, Norman V., Ana Maria Oviedo, and Luis Serven, 2005, "The impact of regulation on growth and informality - cross-country evidence," Policy Research Working Paper Series 3623, The World Bank.
- Loayza, Norman V., Luis Serven, and Naotaka Sugawara, 2009, "Informality in Latin America and the Caribbean," Policy Research Working Paper Series 4888, The World Bank.
- Mandel, Benjamin R., 2010, "Heterogeneous firms and import quality: evidence from transaction-level prices," International Finance Discussion Papers 991, Board of Governors of the Federal Reserve System (U.S.).
- Manova, Kalina, and Zhiwei Zhang, 2012, "Export Prices Across Firms and Destinations," *The Quarterly Journal of Economics*, Vol. 127, No. 1, pp. 379–436.
- McFadden, Daniel F., 1974, *Conditional Logit Analysis of Qualitative Choice Behavior*, pp. 105–142 (Academic Press: New York).
- Mitra, Devashish, and Vitor Trindade, 2005, "Inequality and trade," *Canadian Journal of Economics*, Vol. 38, No. 4, pp. 1253–1271.
- Nataraj, Shanthi, 2011, "The impact of trade liberalization on productivity: Evidence from India's formal and informal manufacturing sectors," *Journal of International Economics*, Vol. 85, No. 2, pp. 292–301.
- Restuccia, Diego, and Richard Rogerson, 2013, "Misallocation and productivity," *Review of Economic Dynamics*, Vol. 16, No. 1, pp. 1–10.
- Schott, Peter K., 2004, "Across-product Versus Within-product Specialization in International Trade," *The Quarterly Journal of Economics*, Vol. 119, No. 2, pp. 646–677.
- Train, Kenneth, 2009, *Discrete Choice Methods with Simulation* (Cambridge University Press).
- Tybout, James R., 2000, "Manufacturing Firms in Developing Countries: How Well Do They Do, and Why?" *Journal of Economic Literature*, Vol. 38, No. 1, pp. 11–44.

A. Appendix

This paper uses data from the following surveys from India:

1. Annual Survey of Industries of 2005-06, 1989-90, 1994-95, 2000-01, and 2009-10
2. Survey of Unorganized Manufacturing of 2005-06, 1989-90, 1994-95, 2000-01, and 2010-11
3. Consumer Expenditure Survey of India of 2003 and 2004-05
4. Employment-Unemployment Survey of India of 2004-05

This section provides some more details regarding these surveys. It also provides a brief description of the County Business Database of the US.

A.1. Annual Survey of Industries

The Annual Survey of Industries (ASI) is conducted by the Central Statistics Office of the Government of India every year. It covers all factories registered under Sections 2m(i) and 2m(ii) of the Factories Act, 1948 i.e. those factories employing ten or more workers using power, and those employing twenty or more workers without using power.

The paper primarily uses data from the 2005-06 ASI (as the SUM was also conducted in 2005-06) which reports data for the financial year ending March 2006. The geographical coverage of the 2005-06 ASI was all of India except the states of Arunachal Pradesh, Mizoram, and Sikkim and the Union Territory of Lakshadweep.

ASI 2005-06 uses the National Industrial Classification (NIC) 2004 (which is closely based on International Standard of Industrial Classification (ISIC) Rev 3.1) to classify economic activity. For all the analysis done in the paper, I restrict the sample to plants which report a 2-digit NIC between 15 to 36 as this constitutes the manufacturing sector and matches the coverage of the SUM. Furthermore, for some of the figures, attention is restricted to 15 large Indian states.⁴⁴

The main variables used in the paper are total employment level of the plant, and details regarding the products produced and inputs used (quantities and rupee values) by each plant.

⁴⁴The main states included are: Andhra Pradesh, Bihar, Gujarat, Haryana, Himachal Pradesh, Karnataka, Kerala, Madhya Pradesh, Maharashtra, Orissa, Punjab, Rajasthan, Tamil Nadu, Uttar Pradesh, and West Bengal. Three Indian states were split into two in 2000. In order to maintain comparability with some of the time-series results in Section 5.2 and 6, the pre-split definition of states is used throughout the paper.

Plants report the average number of employees working in the plant for seven different categories, namely: male workers employed directly, female workers employed directly, child workers employed directly, workers employed through contractors, supervisory and managerial staff, other employees, and unpaid family workers. The size of the plant (total employment) is defined as the sum across all these categories.

All plants report the output they produce using a standardized classification of products called the ASICC product classification. The ASICC has about 5,500 product categories. Plants can report up to ten main products produced in terms of this ASICC classification. Each product category has an associated standardized unit (kilograms, tonnes, numbers, etc) in terms of which the quantity produced is to be reported. Plants also report the total value of production before taxes and distribution costs for each product which can be combined with the information on quantity produced to infer per-unit prices. As all plants are supposed to report quantities in standardized units, the prices inferred should be comparable for all plants producing the same product. However, there seems to be some misreporting in units and this issue is discussed further in Section C. The same commodity classification and units are used to report the quantity and value of materials inputs used.

In addition to the 2005-06 ASI, Section 5.2 also uses data on level of employment of each plant from four other years of the ASI, namely 1989-90, 1994-95, 2000-01, and 2009-10. As with the 2005-06 survey, the broadest definition of employment was used for all years which included part-time workers and unpaid workers. Arunachal Pradesh, Mizoram, Sikkim, and Lakshadweep were excluded from the sample for all years as these states were not covered in the ASI for many of the waves. Different years of the survey used different industrial classifications (NIC 1987, NIC 1998, NIC 2004, and NIC 2008). I created a concordance across these different classifications and only industries which corresponded to 2-digit NIC 2004 between 15 and 36 were included in the sample.

Table A.1 reports the number of observations, estimated number of establishments (using sampling weights provided by the ASI), and the estimated total number of workers employed based on the ASI for all five years that are used in the paper.

More details about the ASI can be found on the website of the Ministry of Statistics and Programme Implementation, Government of India (<http://mospi.nic.in/>).

A.2. Survey of Unorganized Manufacturing

The Survey of Unorganized Manufacturing (SUM) is conducted by the National Sample Survey Office (NSS) of India. The coverage of the survey includes all manufacturing enterprises not registered under Sections 2m(i) and 2m(ii) of the Factories Act, 1948. The SUM is usually conducted every five years. The

last five waves were done in 1989-90, 1994-95, 2000-01, 2005-06, and 2010-11.

The paper primarily uses data from the 2005-06 SUM (62nd Round of the NSS). The survey period was from July 2005 to June 2006.⁴⁵

The geographical coverage of the survey was comprehensive and included all States and Union-Territories of India, with only Leh and Kargil districts of Jammu and Kashmir and a few remote villages in Nagaland and Andaman and Nicobar Islands being excluded. The states of Arunachal Pradesh, Mizoram, and Sikkim and Union Territory of Lakshadweep were dropped to maintain comparability with the coverage of the 2005-06 ASI.

Like the ASI, the SUM 2005-06 uses the National Industrial Classification (NIC) 2004 to classify economic activity. For all the analysis done in the paper, I restrict the sample to plants which report a 2-digit NIC between 15 to 36. Furthermore, for some of the figures, attention is restricted to 15 large Indian states.

The main variables used in the paper are the total employment, and details regarding the products produced and inputs used (quantities and rupee values) by each plant.

Plants report the average number of employees working in the plant for the reference period for which the data is collected (for most plants this was one month). The plants reported the average number of hired workers, working owners, and other workers that they employed on a part-time and full-time basis. Like the ASI, the broadest definition of employment is used with the size of the plant (total employment) being defined as the sum across all these categories.

All plants report the output they produce and material inputs consumed using the same standardized classification of products as is used by the ASI plants. Plants can report up to five main products produced in terms of this product classification. However, unlike the ASI, SUM plants can choose the units in which they are reporting quantities and prices. For example, all ASI plants which produce matchsticks must report quantities in kilograms. However, different SUM plants report quantities and prices of matchsticks in different units including kilograms, tonnes, and numbers (number of matchsticks). I concord units across the two surveys when combining the two surveys. If the same product is being reported in different units which are simple scalar multiples of each other (kilograms and tonnes for example), then I convert the units so that all quantities and prices are being measured in the same unit i.e., divide quantities and prices of all SUM units which report quantities of matchsticks in tonnes by 1000, to get per kilogram prices which are comparable to ASI prices. However, if a SUM plant is reporting the output of matchsticks in numbers, then it is not possible to make this comparable to the the ASI plants which are reporting in kilograms. In such cases, the SUM products are treated as a separate product category.

⁴⁵Note that there is a three month difference in coverage period between the ASI and SUM.

In addition to the 2005-06 SUM, Section 5.2 also uses data on level of employment of each plant from four other years of the SUM, namely 1989-90, 1994-95, 2000-01, 2005-06, and 2010-11. As with the 2005-06 survey, the broadest definition of employment was used for all years which included part-time workers and unpaid workers. The same sampling on states and industries was done as in the ASI.

Table A.1 reports the sample size, number of establishments (using sampling weights provided by the SUM), and the total number of workers employed based on the SUM for all five years that are used in the paper.

More details about the SUM can be found on the website of the Ministry of Statistics and Programme Implementation, Government of India (<http://mospi.nic.in/>).

A.3. Consumer Expenditure Surveys

The National Sample Survey Office of India (NSS) conducts an annual Consumer Expenditure Surveys (Schedule 1.0) in India. From 1972-73, the NSS started a quinquennial series in which every five years, it conducts a survey with a sample size which is about four times larger than the annual survey.

The paper uses data mainly from the 2004-05 (61st Round of the NSS) Consumer Expenditure Survey which was part of the quinquennial series and interviewed about 125,000 households. The geographical coverage of the survey was comprehensive and included all States and Union-Territories of India, with only Leh and Kargil districts of Jammu and Kashmir and a few remote villages in Nagaland and Andaman and Nicobar Islands being excluded. The survey period was from July 2004 to June 2005.

The Consumer Expenditure Surveys of 2004-05 asks households to report the value of consumption for 339 different goods. Households report quantities and rupee values separately for 209 goods, which can be used to compute prices for these goods. 156 of these 209 goods are food items, 10 fall under the “fuel and light” category, another 24 are clothing and footwear, while the remaining are durables.

For food items, households report consumption out of home production (quantities and imputed rupee values) and total consumption (which includes home production and market purchases). The price computed divides total value of consumption by total quantity consumed, thus averaging across home and market consumption.

The reference period for consumption of all food items is 30 days, i.e., households report quantity consumed and rupee values for food consumption for the last 30 days. For clothing and footwear categories, households report consumption for a reference period of 30 days as well as 365 days. The 365 day reference period for these categories is used as many households report zero purchases for these items for the 30 day reference period but positive amounts for the 365 day reference period.

Table 2 uses data from the 2003 (59th Round) Consumer Expenditure Survey which was not part of the quinquennial series and interviewed about 41,000 households. The geographical coverage of the 2003 survey was similar to the 2004-05 survey. The survey period was from January 2003 to December 2003. The consumption items recorded across the two surveys were also very similar with only a few minor differences.

Table A.2 reports some summary statistics for the Consumer Expenditure Survey of 2004-05. It reports the number of items and share of expenditure for five broad expenditure heading and also the share of expenditure within the heading for which prices could be computed. The summary statistics for the 2003 survey are very similar and are not reported.

More details about the dataset can be found on the website of the Ministry of Statistics and Programme Implementation, Government of India (<http://mospi.nic.in>).

A.4. Employment-Unemployment Survey

The National Sample Survey Office of India (NSS) conducts an Employment-Unemployment Survey (Schedule 10.0) as part of its quinquennial series. This paper uses the Employment-Unemployment Survey of 2004-05 (61st Round of the NSS). The geographical coverage of the survey was comprehensive and included all States and Union-Territories of India, with only Leh and Kargil districts of Jammu and Kashmir and a few remote villages in Nagaland and Andaman and Nicobar Islands being excluded. The survey period was from July 2004 to June 2005. In this survey it interviews about 125,000 households (about 600,000 individuals).

The survey asks all the individuals in the household to report demographic characteristics like age, education etc. It also asks individuals to report the main industry in which they work, the size of establishment in which they work, and the wage they earned in the last week. To maintain comparability with the production surveys, only individuals who report a 2-digit NIC 2004 between 15 and 36 are used.

The main variables used from this survey are the education level of individuals and the size category of the establishment in which they work.

The survey asks individuals to report the level of general education that they have achieved. The possible responses are: illiterate, literate but not through formal schooling, primary, middle, secondary, higher secondary, diploma/certificate course, graduate, and post graduate or above. For the purpose of the model, a person was defined as skilled if he or she had finished at least secondary education (Grade ten).

Individuals were also asked to report the size category of the establishment in which they worked. They could report one of the following options: establishment of size less than 6, between 6 and 9, between 10 and 19, 20 or greater, and unknown size

The calibration of the wage premium in Section 4.1 also makes use of wage data from this survey. More

details regarding the construction of this variable along with the Mincerian regression results are provided in Section D.

More details about the Employment-Unemployment Survey can be found on the website of the Ministry of Statistics and Programme Implementation, Government of India (<http://mospi.nic.in/>).

A.5. County Business Patterns Database (US)

The County Business Patterns Database maintained by the US Census Bureau provides level of employment for each 6-digit NAICS for each US county. The employment level is as on the week of March 12th of that year.

The paper uses the 2006 release of the data. For many industry-county cells the exact level of employment is not reported. Instead, the dataset reports an employment size class for that cell. In these cases, the employment in the cell is assigned the midpoint of the size class reported. For example, if a NAICS-County cell reports employment in the size class 'B' which represents 20-99 employees, then the cell is assigned an employment level of 60.

The data can be downloaded from <http://www.census.gov/econ/cbp/>.

B. Inter-State Trade: Concordances

B.1. NAICS 2002 to NIC 2004 Concordance for Herfindahl Index

Section 6 uses a Herfindahl Index of employment concentration across US counties as a measure of tradability for industries in India. While the US County Business Patterns Database uses NAICS to classify economic activity, the Herfindahl Index needs to be based on the Indian classification of economic activity (NIC 2004) for it to be applied using Indian data. In order to construct this index, I created a concordance between 6-digit NAICS 2002 and 3-digit ISIC Rev 3.1 (the classification used by the ASI and SUM is the NIC 2004, which is a one to one match to ISIC Rev 3.1 at the 3-digit level).

The concordance between 6-digit NAICS and 3-digit ISIC Rev 3.1 was based on the Census Bureau's concordance file available at <http://www.census.gov/eos/www/naics/concordances/concordances.html>. Although this file gives a many to many concordance, this was reduced to a one to one concordance by taking the 3-digit ISIC which was the closest fit for each 6-digit NAICS.

Of the 59 3-digit ISIC industries in the manufacturing sector, three industries (182, 231, and 233) were not represented in this concordance i.e. none of the 6-digit NAICS industries were mapped into these 3-digit

ISIC industries. These industries employed only 0.16 percent of the total manufacturing workforce in India in 2005. These industries are dropped for all the analysis which uses the Herfindahl Index.

B.2. HS Product Classification to NIC 2004 Concordance for Export-Import Index

Section 6 also uses a measure of international trade as a proportion of domestic production in India at the 3-digit level for NIC 2004. The export and import data for India was not at the industry level but rather at the product level using the Harmonized Commodity Description and Coding System (HS). The World Integrated Trade Solution (WITS) provides a one to one concordance from HS 2002 to ISIC Rev 3. WITS is based on a collaboration of the World Bank with UNCTAD, WTO and other international organization associated with international trade data.⁴⁶

Two NIC 04 industries (223 and 273) were not represented in this concordance i.e. none of the HS codes mapped into these industries. These industries employed only 0.37 percent of the total manufacturing workforce in India in 2005. These industries are dropped for all the analysis which uses the Export-Import Index. Furthermore, industry 233 (nuclear fuel) had some imports in the trade data but no local production in India. This industry was also dropped.

B.3. Concordances Across Different NIC Revisions

Different years of the ASI and SUM use different revisions of the NIC. The 1989 and 1994 surveys use NIC87, the 2000 surveys uses NIC98, the 2005 surveys use NIC04 and the 2010 surveys use NIC08. I create a concordance from the different NIC revisions to NIC04 at the 3-digit level as the tradability indexes are constructed for NIC04. The concordances were based on official concordance tables which can be found at <http://mospi.nic.in> under the “Economic and Social Classification Heading”.

The NIC04 industries 341 (Manufacturing of motor vehicles) and 342 (Manufacture of bodies of motor vehicles, trailers, and semi-trailers) cannot be separately identified in NIC87. Hence, these two industries are merged into one industry group for all the tradability regressions.

C. Units Misreporting Problem in the ASI

As mentioned in footnote 14, there seems to be a misreporting of units and quantities in the ASI. I discuss an example here to clarify the problem. ASICC code 11401 stands for “milk”. All plants who report that they produce milk are supposed to report the quantity produced in kiloliters (1000 liters) which should mean

⁴⁶The concordance can be found at http://wits.worldbank.org/wits/product_concordance.html.

that when we divide rupee values by the quantity, then it should yield the price of milk that the plant charges in kiloliters. Figure A.2 plots the log of price charged for milk by different plants in the ASI against log of the number of employees in the plant. As can be seen, the log of the price charged by most plants is about ten. However, there is a group of plants who report a price which is seven log points lower or about 1000 times lower ($\exp(7) = 1096$). This is clearly a case of some plants reporting quantities in liters instead of kiloliters which makes the price computed a price per liter.

Such misreporting can potentially bias the results from regressions of price on size if larger plants are more likely to misreport quantities in terms of larger units. To account for this problem, I manually go through about a 1000 product categories to see which product categories suffer from this problem. I split products which suffer from this problem into two separate product categories based on a sensible price cutoff (for the milk example, all plants charging a log price greater than six were placed in a different product category).⁴⁷ As a different product fixed effect is allowed for this new product category, the regressions control for the price level differences arising from misreporting of quantity units. However, the clustering when computing standard errors does not treat the new product category as a separate category which is why the number of product fixed effects exceed the number of clusters in these regressions. Table A.8 compares the results when the units correction described above is implemented versus when it is not implemented. Column 1 is the same as the first column of Table 3 (it corrects for the units problem). Column 2 repeats the regression but does not split products with the units problem into different categories. As can be seen, the price elasticity with respect to employment is smaller when the units problem is corrected, implying that the misreporting of units is correlated with size.

In addition to the manual correction, I also implement an algorithm which identifies product categories for which units have been potentially misreported. The algorithm consists of the following steps:

1. If the maximum price reported for a product is less than 50 times the minimum price, then the product is classified as one with no units misreporting.
2. I first arrange prices in ascending order within a product category. If there are two consecutive prices which are at least different by a factor of 20, and the average price above the jump is between 500 and 2000 times the average price below the jump, then the product is classified as one with a units misreporting problem and is split into two product categories.
3. I run regressions of log of price on log of employment with a dummy which takes value 1 for all observations below a given price i.e. if a product category has 50 plants producing it, I run 50 separate

⁴⁷While in the milk example presented here, the units problem and the appropriate price cutoff was obvious, for some other products the problem is harder to clearly identify. In these cases I use my judgment to decide on the price cutoff.

regressions - in the first the dummy only takes the value 1 for the lowest price plants, for the next regression, the dummy takes value 1 for the two lowest prices and so on. I then compare the highest R-square that I get with the dummies (within the product category) with the R-square when I run a regressions of log of price on log of employment with no dummy. If the difference in R-square is more than 0.75, and the difference in mean prices above the dummy (for the highest R-square) is at least 300 times higher than the average price below the dummy, then the product is classified as one with units misreporting and is split into two product categories.

When implementing the algorithm instead of the manual correction, the elasticity of price to size is 0.1037 in the ASI. The result is not very sensitive to using slightly different thresholds in the three stages of the algorithm. For example, changing the threshold for the R-square in step 3 to 0.8 and 0.7 changes the estimated elasticity to 0.1091 and 0.0986 respectively.

I also implement the algorithm for the input prices regressions and the results are similar to the ones with the manual correction

D. Calibrating Production Parameters - θ_{q_n}

This section provides more details on the calibration of θ_{q_n} , the share of unskilled workers in the intermediate producers production function. As mentioned in the paper, θ_{q_n} is chosen to match the wage premium and the ratio of unskilled to skilled workers for different qualities relative to the lowest quality level.

The target for the wage premium is obtained by running Mincerian type regression using the Employment-Unemployment Survey of 2004-5. Each individual is asked to report the main activities he or she undertook in the last seven days. Individuals can report multiple activities, and report if they were involved in the activity with “full intensity” or “half intensity. The wages earned in the last week are reported for all activities separately (if that activity generated wages). The average wage earned by each individual is computed by dividing the total wage earned for each activity over the last seven days by the number of intensity-days worked (summing across days and treating full intensity as 1 day and half intensity as 0.5 days) in that activity.

The wage premium for skilled workers is computed by running a regression of log of wages on a dummy which takes the value 1 if the worker is skilled (ten or more years of education) controlling for potential experience (age minus years of education minus four) and its square, and dummies for each 4-digit industry, 2-digit occupation, state, sector (urban or rural), and sex. I restrict the sample to workers reporting their industry as manufacturing (2-digit NIC between 15 and 36) and individuals between the age of 15 and 65

only.

Column 1 of Table A.4 reports the results for the regression. The coefficient on the dummy which takes value 1 if a person is classified as skilled is 0.45, implying a wage premium of 56.8 percent which is rounded up to 60 percent when calibrating the model.

Calibrating θ_{q_n} also requires $N - 1$ ratios for equation 9, the unskilled to skilled ratio for different qualities relative to the lowest quality. As mentioned in the paper, size categories reported in the Employment-Unemployment survey are very coarse, and therefore cannot be used to compute eleven ratios for equation 9 for eleven different quality (size) levels. Instead the relation between the size of an establishment and the share of unskilled to skilled workers is extrapolated based on the values reported in Table 6. The table reports that plants of size five or less have a unskilled to skilled ratio of 5.05 while plants of size 5 to 20 have a unskilled to skilled ratio of 2.92. These two points are used to extrapolate the unskilled to skilled ratio for larger sized plants with the ratio taking a minimum value of 0.5 (hire twice as many skilled as compared to unskilled workers). These extrapolated values are used to compute equation 9 for different quality levels given the average size of each quality level.

Table A.1: Summary Statistics: ASI and SUM

	Annual Survey of Industries			Survey of Unorg. Manufacturing		
	Observations	Plants	Employment	Observations	Plants	Employment
1989-90	45	88	6,999	94	13,279	26,968
1994-95	52	105	7,853	156	12,114	29,924
2000-05	30	119	7,762	220	16,994	37,016
2005-06	42	125	8,811	82	17,037	36,376
2009-10	41	144	11,506	98	17,211	34,910

Notes: All numbers are in thousands ('000). The data is from the Annual Survey of Industries and the Survey of Unorganized manufacturing for five different years. The row corresponding to the year 2009 reports results the ASI of 2009-10 but the SUM of 2010-11 . The column "Observations" reports the total number of observations surveyed in the year. The "Plants" and "Employment" columns report the total plants and the total employment in these plants after taking into account the survey weights provided with the surveys. Four states were excluded due to lack of coverage in some years of the ASI. Only plants which reported industries which corresponded to the 2-digit NIC 2004 classification between 15 and 36 were included.

Table A.2: Summary Statistics: Consumer Expenditure Survey

	Items	Share of Expenses	Items with Prices	Share with Prices
Food	161	0.499	156	0.980
Fuel and Light	13	0.094	10	0.932
Clothing and Footwear	27	0.075	24	0.967
Other goods and services	86	0.288	0	0.000
Durables	52	0.044	19	0.449

Notes: The data is from the Consumer Expenditure Survey conducted by the NSS in 2004-05. The rows represent broad expenditure categories. The column "Item" gives the number of distinct goods in the category for which consumption was reported. "Share of Expenses" gives the share of total expenditure that was devoted to the particular expenditure category when summing over all households. "Items with Prices" reports the number of items in the category for which values and quantities were reported, allowing calculation of prices. "Share of Prices" reports the share of expenditure within the category which was devoted to items for which the price could be computed.

Table A.3: Main Activity of Individual - 2003 Consumer Expenditure Survey

Description	Code
Worked in hh enterprise (self-employed): own account worker	11
Worked in hh enterprise (self-employed): employer	12
Worked as helper in hh enterprise	21
Worked as regular salaried/wage employee	31
Worked as casual wage labor: in public works	41
Worked as casual wage labor: in other types of work	51
Did not work but was seeking and/or available for work	81
Attended educational institution	91
Attended domestic duties only	92
Domestic duties and engaged in free collection of goods, sewing, tailoring, etc. for household use	93
Rentiers, pensioner, remittance recipients etc	94
Not able to work due to disability	95
Beggars, prostitutes	96
Others	97

Notes: The 2003 consumer expenditure survey asks each individual in the household to report their main activity during the year. The table lists the different activities which the individuals could report. People who reported codes 92, 93, 94, or 97 were classified as non-workers and households which had at least one person between the age of 15 and 70 who was classified as a non-worker were considered to have low opportunity cost of time.

Table A.4: Wage Premium from Employment-Unemployment Survey

Dependent variable: log(wage)		
	(1)	(2)
skilled	0.450*** (0.0150)	0.445*** (0.0147)
Wage Premium	1.568	1.560
Winsorize 1%		Y
Observations	11,003	11,003

Notes: The data is from the Employment Unemployment Survey of 2004-05. Column 1 reports results for the regression of log of wages earned by an individual on a dummy which takes value 1 if the individual has 10 or more years of education. Column 2 winsorizes 1 percent tails of wages. All regressions include controls for potential experience (age minus years of schooling minus 4) and its square, and dummies for each 4-digit NIC industry, 2-digit occupation, state, sector (urban or rural), and sex. The wage premium implied by the coefficient estimate for skilled is given in the row labeled "Wage Premium". Robust standard errors are reported. ***p<0.01.

Table A.5: Ranking of Industries Based on Tradability Index

	Herfindahl Index	Export-Import Index
Industries Below Median (Non-tradable)	151,152,153,154,155,171,	152,153,154,155,160,171,
	201,202,210,221,222,241,	182,201,202,210,222,231,
	242,251,252,261,269,272,	251,252,269,271,281,293,
	273,281,289,291,292,311	311,313,314,315,341,342,
	312,343,361,369	343,352,359,361
Industries Above Median (Tradable)	160,172,173,181,191,192,	151,172,173,181,191,192,
	223,232,243,271,293,300,	221,232,241,242,243,261,
	313,314,315,319,321,322,	272,289,291,292,300,312,
	323,331,332,333,341,342,	319,321,322,323,331,332,
	351,352,353,359	333,351,353,369

Notes: The table lists the 3-digit industries (NIC04) which fall above and below the median of for the two tradability indexes.

Table A.6: Size Income Relation Across States for Tradables vs. Non-tradables: No Weighting

Dependent Variable: share of employment in <=5 in industry 'i', state 's', time 't'				
	(1)	(2)	(3)	(4)
log(per-capita SNDP)*tradability	0.015 (0.0376)	-0.026 (0.0705)	-0.043 (0.0278)	-0.006 (0.0415)
Index Cutoff	H-index Median	H-index Quartile	Exp-Imp Median	Exp-Imp Quartile
Observations	3,885	1,826	3,899	1,959

Notes: The data is from five rounds of the ASI and SUM. The table reports regression results for the share of employment in plants of size 5 or less in industry 'i' in state 's' at time 't' on log per-capita state NDP interacted with a dummy which takes value 1 if industry 'i' is classified as a tradable industry and 0 if it is classified as non-tradable. Column 1 classifies an industry as tradable if the Herfindahl Index across US counties for the industry was above the median of Herfindahl Indexes, and non-tradable if it was below the median. Column 2 uses top and bottom quartiles of the Herfindahl Index as cutoffs. Column 3 and 4 use the tradability index based on Indian exports and imports and uses the median and the top and bottom quartiles as cutoffs respectively. All regressions include fixed effects for industry interacted with time and state interacted with time. No weights are applied to the observations in the regressions. Standard errors are clustered at the state level.

Table A.7: Size Income Relation Across States for Tradables vs. Non-tradables: Exclude “NEC” and “Others”

Dependent Variable: share of employment in <=5 in industry ‘i’, state ‘s’, time ‘t’				
	(1)	(2)	(3)	(4)
log(per-capita SNDP)*tradability	0.050 (0.0399)	0.036 (0.0525)	0.002 (0.0500)	-0.006 (0.0549)
Index	H-index	H-index	Exp-Imp	Exp-Imp
Cutoff	Median	Quartile	Median	Quartile
Observations	3,219	1,531	3,233	1,593

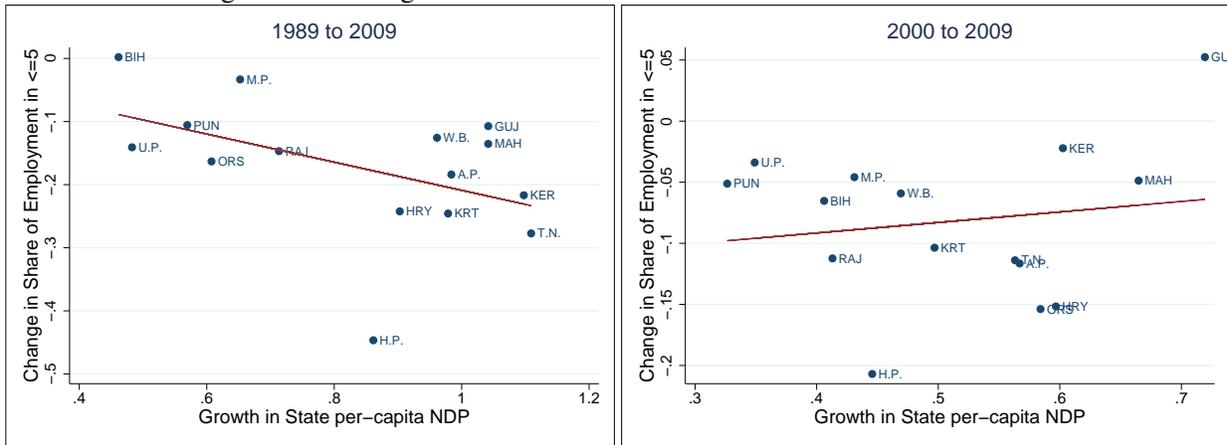
Notes: The data is from five rounds of the ASI and SUM. Residual industries (with words "NEC" or "other") are removed. The table reports regression results for the share of employment in plants of size 5 or less in industry ‘i’ in state ‘s’ at time ‘t’ on log per-capita state NDP interacted with a dummy which takes value 1 if industry ‘i’ is classified as a tradable industry and 0 if it is classified as non-tradable. Column 1 classifies an industry as tradable if the Herfindahl Index across US counties for the industry was above the median of Herfindahl Indexes, and non-tradable if it was below the median. Column 2 uses top and bottom quartiles of the Herfindahl Index as cutoffs. Column 3 and 4 use the tradability index based on Indian exports and imports and uses the median and the top and bottom quartiles as cutoffs respectively. All regressions include fixed effects for industry interacted with time and state interacted with time. Each observation is weighted by the share of observations in the state-industry cell out of the total observations in the ASI and SUM combined for the given year. Standard errors are clustered at the state level.

Table A.8: Units Misreporting Problem in the ASI

Dependent Variable: log(output price)				
	(1)	(2)	(3)	(4)
log(labor)	0.096*** (0.0087)	0.155*** (0.0125)	0.106*** (0.0133)	0.125*** (0.0152)
Units Problem Accounted For	Y	N	Y	N
Sample	ASI	ASI	Both	Both
Observations	46,704	46,704	75,161	75,161
Number of products	1,217	1,077	3,181	3,041
Number of clusters	1,078	1,078	3,042	3,042

Notes: The data is from the ASI and SUM of 2005-06. All columns report results for regressions of log of price charged by plants for their products on log of number of employees hired by the plant. Columns 1 and 2 restrict the sample to the ASI alone while columns 3 and 4 combine the ASI and the SUM. Columns 1 and 3 implement the manual units correction (same as reported in main text) while columns 2 and 4 do not correct for misreporting of units. 1 percent tails of prices (within a product) and plant size are winsorized. All regressions include product fixed effects and state times urban-rural fixed effects. Standard errors are clustered at the product level. ***p<0.01.

Figure A.1: Change in Size Distribution Over Time Across Indian States



Notes: The figure uses data from three waves of the ASI and the SUM (1989, 2000, and 2009). The first figures plot the change in share of employment in plants of size five or less for different states against the change in per-capita NDP in the state between 1989 and 2009. The slope of the linear fitted line is -0.22 with a P-value of 0.077. The second figure takes the change from 2000 to 2009. The slope of the linear fitted line is 0.086 with a P-value of 0.580.

Figure A.2: Price Charged for Milk by Different Plants Against Size

