

# A Balls-and-Bins Model of Trade

By ROC ARMENTER AND MIKLÓS KOREN\*

*Many of the facts about the extensive margin of trade—which firms export, and how many products sent to how many destinations—are consistent with a surprisingly large class of trade models because of the sparse nature of trade data. We propose a statistical model to account for sparsity, formalizing the assignment of trade shipments to country, product and firm categories as balls falling into bins. The balls-and-bins model quantitatively reproduces the pattern of zero product- and firm-level trade flows across export destinations, and the frequency of multi-product, multi-destination exporters. In contrast, balls-and-bins overpredicts the fraction of exporting firms.*

International trade has long been concerned with aggregate patterns—what and how much countries trade with each other. The recent availability of finely disaggregated trade data has spurred a fast-growing research that documents the extensive margin in trade—which firms export, and how many products they send to how many destinations. A number of stylized facts have emerged regarding the incidence and pattern of zero trade flows at the product and the firm level, the size and frequency of exporters, and the size and frequency of multi-product and multi-destination exporters.<sup>1</sup>

Are these facts useful for testing new theories of the extensive margin of trade? We argue that some, though not all, of the facts are consistent with a surprisingly large class of trade models. For example, there are many trade models successful at reproducing the distribution of trade across countries—typically relying on a gravity equation—and the sector composition of trade. All of these models, even those without an explicit treatment of the extensive margin, will replicate the

\* *Armenter*: Federal Reserve Bank of Philadelphia (10 Independence Mall, Philadelphia, PA 19106). E-mail: roc.armenter@phil.frb.org. *Koren*: Central European University (Nádor utca 9., Budapest 1051, Hungary), MTA KRTK and CEPR. E-mail: koren@ceu.hu. For useful comments we thank the referees as well as George Alessandria, Arnaud Costinot, Alan Deardorff, Jonathan Eaton, Tom Holmes, László Mátyás, Marc Melitz, Virgiliu Midrigan, Esteban Rossi-Hansberg, Peter Schott, Adam Szeidl, Ayşegül Şahin, and seminar participants at the Federal Reserve Bank of New York, the Institute for Advanced Studies in Vienna, Central European University, UC San Diego, Princeton, the SED, the NBER Summer Institute, Michigan, Stanford, MIT, the University of Zurich, CEPR ERWIT, and the LSE. We also thank Jennifer Peck for excellent research assistance. Much of this research was carried out while Koren was visiting the International Economics Section of Princeton University, and he gratefully acknowledges their hospitality. The views expressed here do not necessarily reflect the views of the Federal Reserve Bank of Philadelphia or the Federal Reserve System.

<sup>1</sup>The following is a necessarily incomplete list of references: Helpman, Melitz and Rubinstein (2008), Baldwin and Harrigan (2011), Haveman and Hummels (2004) and Hummels and Klenow (2001, 2005) on zero trade flows; Bernard and Jensen (1999), Bernard, Eaton, Jensen and Kortum (2003), Bernard, Jensen and Schott (2007), Bernard, Jensen, Redding and Schott (2007) and Eaton, Kortum and Kramarz (2004, 2007) for firm-level facts. See the main text and the web Appendix A for further discussion.

prevalence and pattern of zero trade flows as well. Similarly, provided a model matches the observed skewness in export sales, it will predict that most exports are done by relatively few, but large, multi-product, multi-destination exporters. The reason lies in the nature of trade data.

Our first observation is that disaggregate trade data are categorical by construction. We observe a finite number of shipments which constitute the basic units of observation.<sup>2</sup> Each shipment is then assigned a unique category in each classification, that is, the shipment belongs to one of many product codes, goes to one of many destination countries, and is sold by one of many firms in the economy.

The categorical nature of trade data is important, because these data are typically *sparse*. That is, the number of observations is low relative to the number of possible classifications. There were about 22 million export shipments originating in the U.S. in 2005—and thus the same number of observations. At the same time, there are 229 countries and 8,867 product codes with active trade, so a shipment can have more than 2 million possible country-product classifications. More than 40 percent of the traded country-product pairs had only 1 or 2 shipments during the year, a clear sign that the data are sparse.<sup>3</sup>

We propose a parsimonious statistical benchmark for sparse categorical datasets to discern which data moments are informative about the correct model of the extensive margin. Our starting point is a simple random-assignment model which only requires information on the number of observations and the marginal distributions across categories, say, the distribution of trade across countries or across products. Importantly, we do not need to know which product is exported to which country. We instead keep our benchmark as atheoretical as possible by assuming no systematic relationship between categories when constructing the joint distribution across all classifications, e.g., each country is expected to receive the same distribution of products and so on.

We formalize the assignment of observations to categories as balls falling into bins. Each observation constitutes a discrete unit (the ball), which, in turn, is allocated into mutually exclusive categories (the bins). Because we want a parsimonious benchmark, the model assigns balls to bins at random. That is, a ball falling in a particular bin is an independent and identically distributed random event whose probability distribution is determined solely by the distribution of bin sizes.

In spite of its simplicity, the balls-and-bins model makes a rich set of predictions. After a number of balls, some bins will end up empty and some will not. Among the latter, some will contain a large number of balls, some few. These are taken

<sup>2</sup>Trade data are collected through customs forms, one for each export shipment. See Appendix A for details.

<sup>3</sup>Statistical inference on categorical datasets requires a much larger sample size. Indeed, the sample size needed grows very fast with the number of categories  $K$ . For example, the number of observations must be of order  $O(K \log(K))$  for maximum likelihood estimates to exist. See Section 9.8 of Agresti (2002) for a summary discussion of statistical inference in sparse categorical data.

to be the model’s predictions for the extensive and intensive margin, respectively. We characterize the prevalence of zeros and how it varies with the number of balls and the bin-size distribution. These are indeed all the model’s systematic relationships between export flows and the extensive margin: the *assignment* of balls to bins is random.

The number of balls and the distribution of bin sizes are treated as parameters in the model. These can be calibrated directly from the data, i.e., the actual number of observations, the observed distribution of trade across countries; or from the corresponding predictions of a structural model or a reduced form specification like the gravity equation. Our baseline calibration assumes no systematic relationship between firms, countries and products: for example, the probability a shipment is classified as a certain product is independent of the shipment’s destination country. We chose this approach to avoid imposing any structure regarding the joint distribution of trade across products, countries and firms.

We contrast the balls-and-bins model against the data for several statistics about the extensive margin typically reported in the literature (see Table 1). Statistics that are matched by the balls-and-bins model are not helpful in differentiating models of the extensive margin of trade: these statistics will be consistent with a very large class of models, namely, any model that is successful at matching the marginal distributions across categories—the distribution of trade across countries and products. Many trade models, including several with no explicit treatment of the extensive margin, can reproduce accurately how aggregate trade flows vary across countries (the gravity equation) and are capable of encompassing the necessary heterogeneity at the product and firm level.

If, on the other hand, the balls-and-bins model does not match a statistic, then a trade model has a chance to distinguish itself from the competition by positing the correct structure on the joint distribution of trade across countries, products and firms. Indeed, we envision the balls-and-bins model establishing a benchmark for the quantitative evaluation of structural models in sparse datasets.

Table 1 summarizes our findings. For twelve statistics we report the data and the corresponding prediction by the model—the details on both are in the main text. Zero product-level trade flows are as prevalent in the model as in the data; the pattern of zeros across export destinations is also the same. Indeed, we replicate facts regarding zeros as long as trade flows across countries follow a gravity specification, and the trade shares across product categories are skewed. Trade with most of the countries is then very small and most of the traded product categories are tiny. It is exactly for these country-product pairs that the trade flows are missing in the data. They go missing in the model as well: few balls and tiny bins make for many empty bins.

The balls-and-bins benchmark also matches several of the firm-level facts: in the model, as in the data, most firms export a single product to a single country, but these firms represent a very small fraction of total exports. It is the left tail of the distribution of exports across firms which proves key to reproducing

Description	Data	Balls-and-bins
HS10-level product $\times$ country U.S. export flows		
Share of zeros	82%	72%
OLS coefficient of non-zero flow on GDP	0.08	0.10
Firm $\times$ country U.S. export flows		
Share of zeros	98%	96%
Gravity equation for firms, GDP OLS coefficient	0.71	0.60
Single-product exporters		
Fraction of total exporters	42%	43%
Share of total exports	0.4%	0.3%
Single-destination exporters		
Fraction of total exporters	64%	45%
Share of total exports	3.3%	0.3%
Single-destination, single-product exporters		
Fraction of total exporters	40%	43%
Share of total exports	0.2%	0.3%
Exporters in U.S. manufacturing		
Fraction of total firms	18%	74%
Size-premium of exporters	4.4	34

TABLE 1—SUMMARY OF FINDINGS

*Details on sources, data and model are in the main text and in the web Appendix A.*

these firm-level facts. Most exporters are tiny and are hence assigned only one ball in the model. They are thus predicted to be single-product, single-country exporters. Several models in the literature are able to reproduce the skewness in sales; the incidence and relative size of single- and multi-product exporters follows for all of them.

Finally, we attempt to predict the share of exporters among manufacturing firms. According to the balls-and-bins model, 74 percent of firms should export — in contrast with 18 percent in the data. The balls-and-bins benchmark is clearly nowhere close to the data, signaling that the relationship between firms and foreign market access is not driven by sparsity alone. Thus the split between exporters and non-exporters is a useful statistic to discern among structural models: matching the data requires imposing the correct *systematic* relationship between firms and export flows.

We must emphasize that in a dense dataset—i.e., with many observations relative to the number of categories—the balls-and-bins model would be unable to match *any* fact on the extensive margin. Indeed, all bins would be non-empty and the predictions for the extensive margin would be trivial.

In Section V, we illustrate why identifying theories can be difficult when data

are sparse and how the balls-and-bins benchmark can identify which statistics remain informative nonetheless. We show how to take structural trade models to sparse categorical data and under what conditions will the predictions of these models be similar to balls and bins. In other words, when does a richer structural model nests balls and bins as a special case? The predictions of any two such models will be close to balls and bins, and hence, to one another, whenever balls and bins matches the data. If the balls-and-bins prediction is far from the data, however, we can identify the relevant model more precisely. We illustrate each case with the fraction of product-country zeros and the firm's export participation margin, respectively.<sup>4</sup>

A paper close to us in spirit is Ellison and Glaeser (1997). They ask whether the observed levels of geographic concentration of industries are greater than would be expected to arise randomly. To this end they introduce a "dartboard" model of firm location. In contrast with our results, the dartboard model reaffirms the previous results on geographic concentration. Ellison and Glaeser (1997) are also able to provide a new index for geographic concentration which takes a value of zero under the dartboard model and thus controls for the mechanical degree of concentration arising from randomness. Such an index is more difficult for trade facts which do not focus on a particular dimension.

The questions sparsity raises are similar to the debate about the theoretical content of the gravity equation for bilateral trade flows. The gravity equation is hugely successful in predicting trade flows, yet it may be of limited use in distinguishing trade theories. Deardorff (1998) argues that "just about any plausible model of trade would yield something very like the gravity equation," hence the gravity equation should not be the basis for favoring one theory over another. Evenett and Keller (2002) and Haveman and Hummels (2004) also show that the gravity equation is consistent with both complete and incomplete specialization models.

Our paper is also related to a large literature that tests the robustness of empirical findings through Monte Carlo techniques or sensitivity analysis. To our knowledge these tests have not been commonplace in international trade. An early exception is the analysis on trade-related international R&D spillovers in Keller (1998). There has also been some work on the robustness of gravity equation models. Ghosh and Yamarik (2004) use Leamer extreme bounds analysis to construct a rigorous test of specification uncertainty and find that the trade creation effect associated with regional trading arrangements is fragile. Schaefer, Anderson and Ferrantino (2008) use Monte Carlo experiments to explore alternative specifications of the gravity model and find coefficient bias to be pervasive.

<sup>4</sup>In Appendix D we show how to nest the balls-and-bins framework and generate sparse data predictions from a standard trade model featuring comparative advantage and economies of scale. We also show that the model encompasses our baseline calibration as a special case. In a numerical exercise we have no problem pinning down the firm's fixed cost of exporting, a parameter closely tied to the share of exporters; however, the fraction of zeros in trade flows fails to identify the relevant parameters.

## I. Balls and bins

We model the assignment of export shipments to categories as balls falling into bins. The balls-and-bins model reproduces the categorical structure inherent in disaggregate trade data. A trade flow (such as total exports from the U.S. to Argentina, or total exports of a given firm) is composed of a finite number of shipments, each of them a discrete unit of observation (the balls). Every shipment has been classified into mutually exclusive categories, for example, into one of the 10-digit Harmonized System product classifications (the bins).

Formally, let  $n \in \mathbb{N}$  be the number of balls (observations). Let  $K \in \mathbb{N}$  be the number of bins (categories), each of them indexed by subscript  $i \in \{1, 2, \dots, K\}$ . The probability that any given ball lands in bin  $i$  is given by the bin size  $s_i$ , with  $0 < s_i \leq 1$  and  $\sum_{i=1}^K s_i = 1$ . Thus where a ball lands is independent of the number and location of the other balls.

The state of the system is given by the full distribution of balls across bins,  $\{x_1, x_2, \dots, x_K\}$ . Clearly, this distribution is a random variable. Since we are primarily interested in the “extensive margin,” that is, the split between empty and non-empty bins, we define  $d_i$  to be an indicator variable that takes the value of 1 if bin  $i$  is non-empty,  $x_i > 0$ , and 0 otherwise. The “intensive margin” will be given by the number of balls per non-empty bin.

Figure 1 shows that the balls-and-bins model looks as simple as it sounds. Figure 1A depicts five bins, ordered by size. Figure 1B shows a particular realization after throwing seven balls. Bins 3 and 5 are empty and thus we have  $d_3 = d_5 = 0$ .



Fig. 1A

Fig. 1B

FIGURE 1. BALLS AND BINS

We can derive the key moments of the model analytically. For given bin sizes  $\{s_1, s_2, \dots, s_K\}$ , the joint probability of a number of balls  $\{x_1, x_2, \dots, x_K\}$ , is given by the multinomial distribution,

$$\Pr(x_1, x_2, \dots, x_K) = \frac{n!}{x_1! x_2! \dots x_K!} s_1^{x_1} s_2^{x_2} \dots s_K^{x_K},$$

where  $n = \sum_{i=1}^K x_i$ . Note that, given a total number of balls  $n$ , the particular number of balls in two given bins,  $x_i$  and  $x_j$ , are not independent random vari-

ables. A ball falling in bin  $i$  is a ball less falling elsewhere, so it reduces the expected number of balls in bin  $j$ .

The model has a known probability distribution for the extensive margin. After dropping  $n$  balls the expected value of  $d_i$  is the probability that bin  $i$  receives at least one of those:

$$E(d_i|n) = 1 - \Pr(x_i = 0|n) = 1 - (1 - s_i)^n.$$

Each ball has a  $(1 - s_i)$  probability of landing elsewhere. Where a ball lands is an independent event, therefore the probability that none of  $n$  balls falls in a given bin  $i$  is  $(1 - s_i)^n$ . Obviously, as the number of balls increases, it is less and less likely that any given bin remains empty. In the limit, as  $n \rightarrow \infty$ , the probability  $\Pr(x_i = 0|n)$  is zero for all bins  $i \in K$ .

We denote the total number of non-empty bins by  $k$ ,

$$k = \sum_{i=1}^K d_i.$$

Clearly,  $k$  is a random variable itself with  $k \in \{1, 2, \dots, K\}$ . Since the number of non-empty bins is a sum of random variables, we easily obtain that

$$(1) \quad E(k|n) = \sum_{i=1}^K [1 - (1 - s_i)^n].$$

This is our key statistic out of the balls-and-bins model. We will use it to derive many of the facts on the extensive margin, both at the country and at the firm level.

The comparative statics with respect to the number of balls are as one would expect: more shipments increase the expected number of non-empty bins. Perhaps more subtly, the relationship is not linear. The first few balls fall into distinct bins almost surely. Because of that, as long as balls are few, the number of filled bins is close to the number of balls and the relationship is essentially linear. In other words, most adjustment is on the “extensive margin.” As the number of balls increases, it is more and more likely that balls fall in non-empty bins, and the number of filled bins trails the number of balls.<sup>5</sup> Hence, the relationship flattens out and the number of filled bins increases slowly. The remaining balls can only add to the “intensive margin.” Note that the model is very stark in its limiting predictions as the number of shipments grows large: the number of empty bins converges almost surely to zero.

The expected number of non-empty bins also depends on the distribution of

<sup>5</sup>The first ball falling into a non-empty bin comes very early, roughly in proportion to the square root of the number of bins,  $\sqrt{K}$ . This is sometimes known as the “birthday paradox:” it takes only 23 balls before any one of 365 equal-sized bins will contain two or more balls with probability 1/2.

bin sizes. Two bins of equal size fill up very fast: toss a coin ten times and, with almost absolute certainty, the coin will have turned heads some times and tails some others. But if a bin is, say, 10 times the size of the other, then a lot of balls may be needed to hit the small bin. This property of the model will play an important role later, as in many of the quantitative exercises the distribution of bin sizes is particularly skewed.

Formally, the expected number of non-empty bins (1) is convex in  $s_i$  for all  $n \geq 2$ . This implies that as we even out a bin-size distribution, the expected number of non-empty bins increases.

**Proposition 1.** *Let  $\{s_i\}$  be a bin-size distribution and let*

$$(2) \quad \{\tilde{s}_i\} = \alpha\{s_i\} + (1 - \alpha)1/K$$

*for  $\alpha \in [0, 1]$ . Then for all  $n \geq 2$  the expected number of non-empty bins under  $\{\tilde{s}_i\}$  is not less than under  $\{s_i\}$ .*

In some occasions we will focus not on the extensive margin but on zeros, that is, the number of empty bins. It is, of course, trivial to derive the corresponding statistic:

$$K - E(k|n) = \sum_{i=1}^K (1 - s_i)^n.$$

This is clearly decreasing in the number of balls,  $n$ .

We are also interested in the proportion of firms that sell only one product or serve only one country. To this end we derive the probability that a single bin contains all the balls or, equivalently, that exactly one bin is non-empty. Each ball had  $s_i$  probability of falling into bin  $i$ , so with probability  $s_i^n$  all balls fell in bin  $i$ . Of course, this could happen to any of the  $K$  bins, but they are mutually exclusive events. Hence,

$$(3) \quad \Pr(k = 1|n) = \sum_{i=1}^K s_i^n.$$

The probability of a single non-empty bin decreases with the number of balls,  $n$ , and increases with the dispersion of bin sizes. Again, the model becomes degenerate as the number of balls grows: the probability of a single non-empty bin converges to zero.

#### A. Multiple classification systems

So far we have derived the relevant moments for a single trade flow. Often, however, we will be interested in aggregate statistics that involve many trade flows. For example, we will look at the fraction of empty product categories for total U.S. exports as well as how this fraction varies across destinations.



In order to derive aggregate statistics we need to work with the dataset as a whole. The key difference is that each shipment is now classified along many dimensions. For example, in a dataset containing all U.S. export each shipment is given one HS code as well as one export destination out of many different countries.

We introduce a two-dimensional version of the balls-and-bins model, where each shipment is randomly assigned a classification in two systems, with  $T$  and  $K$  categories, respectively.<sup>6</sup> There is, conceptually, nothing different from the previous case: we can always re-arrange the classification system into a row of bins of length  $TK$ , so that  $s_{ij}$  denotes the likelihood of the ball falling into category  $i$  along the first dimension and category  $j$  along the second. Using (1) we can derive the expected number of zeros,

$$(4) \quad E(k|n) = \sum_{j=1}^T \sum_{i=1}^K [1 - (1 - s_{ij})^n],$$

and similarly for the remaining predictions.

In order to keep the benchmark as parsimonious as possible, we assume each ball is randomly and *independently* allocated across classification systems. There is thus no systematic relationship across categories. For example, destination countries would buy the same basket of products in exactly the same proportions; or all exporters are equally likely to sell a given product to a given market. The independence across classifications suits our pursuit of a parsimonious benchmark as we would only require information on the marginal distribution across categories and not on the joint distribution.<sup>7</sup>

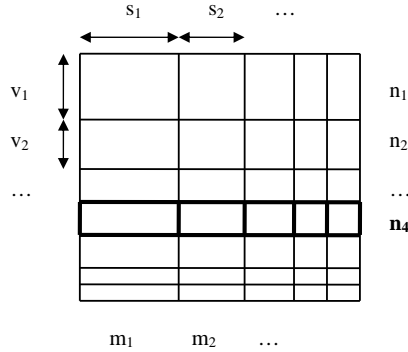


FIGURE 2. BALLS AND BINS:  $T$  BY  $K$  CASE

<sup>6</sup>It is also easy to extend the model to higher-dimensional classification systems.

<sup>7</sup>The model remains tractable if we want to introduce a richer specification: see Section 5 and Appendix D on how to generate sparse data predictions from structural models.

Visually, one can think of throwing balls over a  $T$  by  $K$  grid of bins as in Figure 2. Each classification system comes with its size distribution,  $v_1, v_2, \dots, v_T$  and  $s_1, s_2, \dots, s_K$ , which in Figure 2 pins down the size of rows and columns, respectively. The probability of a given ball falling in the bin  $(i, j)$  is then  $v_i s_j$ .

An additional advantage of approaching the dataset as a whole is that it allows working with *conditional* moments, for example, the number of empty product bins for a given country. For each realization of ball throws there will be a number of balls in each row and in each column, denoted  $n_1, n_2, \dots, n_T$  and  $m_1, m_2, \dots, m_K$ , respectively. (Note that  $n_i$  or  $m_j$  may be zero.) We can then ask the distribution of balls across columns  $1, 2, \dots, K$  within a given row with  $n_j$  balls.

More interestingly, we can compute the statistics of interest given a distribution of balls  $n_1, n_2, \dots, n_T$  across rows. This will allow us, for example, to derive how the fraction of zero product-level bilateral flows varies across U.S. export destinations using the actual aggregate export flows. Since the classification in each system is independent, the conditional statistics for any given row are as in the first version of the model. Let  $k_t$  denote the number of non-empty bins in row  $t$ . We can thus easily construct the distribution of the expected number of non-empty bins per category  $t \in T$  using (1):

$$(5) \quad E(k_t | n_t) = \sum_{i=1}^K [1 - (1 - s_i)^{n_t}],$$

for  $n_t \in \{n_1, n_2, \dots, n_T\}$ . The expected total number of non-empty bins given  $\{n_1, n_2, \dots, n_T\}$  is thus

$$(6) \quad E(k | n_1, n_2, \dots, n_T) = \sum_{j=1}^T \sum_{i=1}^K [1 - (1 - s_i)^{n_j}].$$

It is important to note that, since  $\{n_1, n_2, \dots, n_T\}$  is a random variable, conditional aggregate statistics will not coincide with the corresponding unconditional expectation  $E(k | n)$  with  $n = \sum_{j=1}^T n_j$ .

Similarly, we can compute the probability of a single non-empty bin for each row using (3). We can then derive the proportion of rows which are expected to contain a single non-empty bin. This will allow us, for example, to derive how the fraction of single-product exporters varies across U.S. export destinations using the actual aggregate export flows. As discussed above, the conditional statistics for any given row are as in the first version of the model.

$$\Pr(k_t = 1 | n_1, n_2, \dots, n_T) = \frac{1}{T} \sum_{j=1}^T \sum_{i=1}^K s_i^{n_j}.$$

In practice, we will sometimes approximate the distribution of balls across rows with some parametric distribution. The web Appendix C shows how to compute

aggregate statistics in this case. The Appendix C also describes how to compute the fraction of balls that are expected to fall into single non-empty bin rows: this is useful when we want to derive the fraction of exports originated in single-product or single-destination exporters.

## II. Zeros in trade flows

### A. Product-level zeros

The first data pattern we explore is the prevalence of product-level zeros (i.e., missing trade flows) in country-level exports. In other words, we look at the extensive margin of products when the units of observation are countries. We later discuss firm-level evidence.

We also take the opportunity to carefully describe how we map the data to the balls-and-bins model and back. The methodology is essentially the same for every exercise in the paper.

### B. The facts

Baldwin and Harrigan (2011) recently reported that most potential destination-country product combinations are missing in U.S. exports. In 2005, the U.S. exported 8,867 different 10-digit Harmonized System categories to 229 different countries. Of these 2,030,543 potential trade flows, 1,666,046 (or 82 percent) were missing.<sup>8</sup> In other words, the average country only bought 18 percent of the 8,877 products the U.S. exports. Helpman, Melitz and Rubinstein (2008) look at the country-level zeros in the gravity equation. Of all potential country pairs, only about 50 percent have positive trade in either direction.<sup>9</sup>

**Empirical regularity 1.** *Most of the potential product-country export flows are zero — 82 percent of them in the U.S.*

Other levels of aggregation lead to a similar incidence of zeros. Table 2 reports the incidence of zeros for four classification levels. Zeros only stop being prevalent at the very broad, 2-digit level.

Baldwin and Harrigan (2011) then report how the incidence of zeros relate to the size of the importer and its distance to the U.S. Larger countries that are closer buy a larger variety of products. Larger countries are more likely to import any given product. The same is true for richer countries. The incidence of non-zero flows decreases with distance: closer countries have more non-zero flows than farther countries (the omitted category is the intermediate distance).

**Empirical regularity 2.** *The incidence of non-zero product exports increases with destination-country size and decreases with distance.*

<sup>8</sup>Haveman and Hummels (2004) report a similar incidence of zeros for imports.

<sup>9</sup>Hummels and Klenow (2005) also look at the product-margin of aggregate exports. They have a different measure of the extensive margin.

Classification	Number of bins	Incidence of zeros
10-digit	8,877	82%
6-digit	5,182	79%
4-digit	1,244	66%
2-digit	97	36%
Section	21	16%

TABLE 2—THE INCIDENCE OF ZEROS UNDER DIFFERENT CLASSIFICATIONS

### C. From the data to the model

In order to map the balls-and-bins model to the data, we proceed as follows. The trade flow of interest is the total U.S. exports to a given country, that is, we will have as many trade flows as destination countries (229). We measure the number of shipments going to a country to calibrate the number of balls. For example, Canada (the biggest importer) received 7.3 million shipments in 2005. Equatorial Guinea, the median buyer of U.S. exports, had 2,641 shipments.<sup>10</sup>

The bins correspond to the 8,867 10-digit HS categories in which the U.S. exports at all.<sup>11</sup> The size of each bin ( $s_i$ ) is the share of each HS code in *total* U.S. exports in 2005. That is, we divide the number of export shipments in a given HS code with the total number of shipments (21.6 million).

We then calculate the expected number of non-empty bins for each country using the previous formula (1),

$$E(k_c|n_c) = \sum_{i=1}^{8867} [1 - (1 - s_i)^{n_c}],$$

where  $n_c$  is the number of balls for country  $c$  and  $k_c$  is the number of non-empty HS categories in exports to country  $c$ . The expected number of non-empty bins overall is then

$$E(k|n_1, n_2, \dots, n_{229}) = \sum_{c=1}^{229} k_c.$$

Note that we are computing the expectation conditional on the number of export shipments from the U.S. to each country. To retrieve the incidence of zeros we only need to subtract from and divide by the appropriate number of categories; 8,867 if we are looking at the zeros for a particular trade flow, or  $229 \times 8,867$  for

<sup>10</sup>Web Appendix A discusses the definition of shipments and the collection of trade data in more detail. A natural question is why there are no more trade shipments in a year. Hornok and Koren (2013) document how the number and size of shipments vary across products, mode of transportation and with the administrative barriers set by trading countries.

<sup>11</sup>We ignore the 121 HS codes for which we did not observe any shipments in 2005. It is possible to account for the missing bins with a simple specification: if anything, ignoring the missing bins reduces the expected fraction of zeros in the model.

overall U.S. exports.

The assumption underlying this calibration is that each destination country would buy the same basket of American products in exactly the same proportions, that is, we have matched the marginal distribution across products and assumed that the product bin a ball falls is independent of the country of destination. The only difference across countries is that smaller countries (such as Equatorial Guinea) have a smaller sample of shipments—drawn from the same distribution—than larger ones (such as Canada). Most trade theories are concerned with the joint distribution of trade across countries and products: our calibration provides a neutral, atheoretical benchmark.

#### *D. The model's predictions*

We find that indeed most of the potential product-level bilateral flows are zero in the model. The expected share of zeros is 72 percent, surprisingly close to the data (82 percent). That is, seven out of every eight zeros are to be expected given the sparsity of the data. Table 3 reports the predicted fraction of zeros for other levels of sectoral aggregation. The model's predictions track the observed incidence of zeros rather well at all levels.

Classification	Number of bins	Data	Balls and bins
10-digit	8,867	82%	72%
6-digit	5,182	79%	68%
4-digit	1,244	66%	52%
2-digit	97	36%	23%
Section	21	16%	10%

TABLE 3—THE INCIDENCE OF ZEROS UNDER DIFFERENT CLASSIFICATIONS

Moreover, the model matches quantitatively the pattern of zeros across flows in the data. To show this, we plot the number of exported products for each destination country against the total number of export shipments to that country in Figure 3. The dots represent the actual number of products in the data, the line is the predicted number of non-empty bins for each country. We already know that the balls-and-bins model somewhat underpredicts zeros, hence overpredicts the number of exported products, but the shape of the relationship to total exports is strikingly similar.

Zeros are more likely to occur in small export flows (those with few balls). This already suggests that non-zero flows may follow a gravity equation, as total export flows are well known to adhere to gravity. We then try to replicate the gravity specification in Baldwin and Harrigan (2011). We take the predicted probability of a non-zero flow  $(1 - (1 - s_i)^{n_c})$  and regress it on the gravity variables such as country size and distance. We emphasize that the balls-and-bins model has

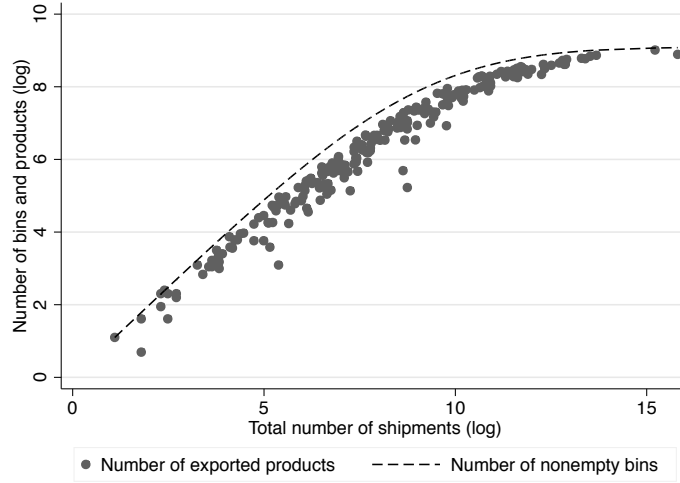


FIGURE 3. THE NUMBER OF SHIPMENTS AND THE NUMBER OF PRODUCTS

nothing to say about gravity, but given that the total number of balls ( $n_c$ ) is highly correlated with the gravity variables, we may find some significant correlations.

The second column of Table 4 reports the results for the balls-and-bins model. For convenience, we replicate the regression in Baldwin and Harrigan (2011) and report the resulting coefficients in the first column.<sup>12</sup> Bigger and closer countries are more likely to have a non-zero flow under the balls-and-bins model, just as in the data. Moreover, the magnitudes of the coefficients are surprisingly similar. The only exception are the two countries bordering the U.S. (“distance= 0”), Canada and Mexico. These seem to import more HS codes in the data than under the balls-and-bins model.

Quantitatively, the dispersion in flow and bin sizes plays an important role. In both cases the distribution is skewed, that is, some product categories and U.S. trade partners are very large, but the vast majority of product categories and trade partners are very small. It is precisely for the combination of latter (small country export for a small product category) that we have the missing trade flows in the data. And it is precisely for smaller bins and fewer balls that the model predicts the most zeros.

Let us start with the distribution of bin sizes. The size of the average bin is  $1/8867 = 1.13 \times 10^{-4}$ . However, the size distribution across bins is rather skewed.

<sup>12</sup>For the top 99 trading partners of the U.S., we regress the incidence of a positive export flow on real GDP of the importer, real GDP per capita, and the distance of the importer from the U.S. using a linear probability model, so coefficients can be understood as marginal effects. Distance is divided in the same categories as in Baldwin and Harrigan (2011). Standard errors are clustered at the country level. These results are comparable to Table 4 of Baldwin and Harrigan (2011). The coefficients are similar, but not identical, potentially due to somewhat different real GDP measures.

	Non-zero trade flow	B+B model
Real GDP	0.078*** (0.007)	0.099*** (0.008)
Real GDP per capita	0.045** (0.011)	0.063*** (0.012)
Distance = 0	0.338*** (0.060)	0.215*** (0.033)
0 < distance < 4000km	0.235*** (0.024)	0.239*** (0.029)
4000 < distance < 7800	omitted	omitted
7800 < distance < 14000	0.010 (0.035)	-0.013 (0.040)
Distance > 14000	0.034 (0.036)	0.015 (0.047)
Observations	868,966	868,966
Clusters	98	98
$R^2$	0.23	0.44

TABLE 4—NON-ZERO FLOWS AND GRAVITY – *Balls and bins*

The size of the median bin is  $2.2 \times 10^{-5}$ , about five times smaller than the average. For comparison, we find 57 percent zeros if we assume that all 8,867 HS codes have the same size.

What is the source of this skewness across product categories? Category sizes may partly reflect the export specialization of the U.S., as higher exports of a product make that product category bigger. However, they are also affected by the nature of the classification system. As an illustration, we flag all product categories that contain either of the words “parts,” “other,” “e.t.c.” and “n.e.s.o.i.” (for “not elsewhere specified or included”) as *catch-all* categories. These are probably heterogeneous aggregates of various products. Of the 100 biggest categories, 72 are such catch-all. In contrast, only 13 of the 100 smallest categories are catch-all.

It is important to emphasize that it is the dispersion in bin sizes, and not some particular bins being large and other small, that leads balls and bins to predict so many zeros. To check for this, we re-run the model with the bin-size distribution calibrated to the HS shares of U.S. exports to Canada and Mexico only. These two trade flows contain very few zeros and so the size distribution of bins would not be affected by the large incidence of zeros in the data. The predicted fraction of zeros under these bin sizes is 75 percent. We find similar predictions if we use the shares of other countries or some exogenous bin-size distribution with skewness.

The skewness of trade flows is also important. Canada alone accounts for more

than one fifth of total U.S. exports; the top five U.S. trade partners account for more than a half of the total. In order to shut down any shipment size variation across destinations, we also computed the fraction of zeros by dividing export flows (in dollars) by the average shipment value, \$36,000. The fraction and pattern of zeros are virtually unchanged.

We also replace the actual trade flows with the trade flows *predicted* by the gravity equation. We find 67 percent zeros, the number being slightly lower than the baseline result mainly due to the reduced country sample. This exercise also allows us to pin the key determinant of the skewness in trade flows. Assuming distance has no effect on trade flows does not change the number of zeros from 67 percent. In contrast, assuming all countries have identical size brings the fraction of zeros down to 30 percent. Thus it is the skewness in country size, through its impact on export flows, that is most important for the calibration.

On a more positive note, we show how a quantitative evaluation of the models could elicit some identification. We underpredict the fraction of zeros as well as the impact of distance. Both effects are relatively small so we would need trade models capable of matching the data with precision.

#### *E. Firm-level zeros*

Zeros in firm-level trade flows display a remarkably similar pattern in the data. The average exporting firm in 2000 shipped goods to only 3.5 countries from a total of 229.<sup>13</sup> In other words, 98 percent of potential firm–country trade flows are zero. Again, zero trade flows follow a well-defined spatial pattern, with zeros being more frequent for small, distant countries.

We can calibrate the balls-and-bins model to study zeros in firm-level trade. The number of balls per destination country are again taken by counting the shipments going to that country. The key difference is that now we need to create bins for *firms* as opposed to product categories. The total number of bins equals the number of exporting firms, 167,217.<sup>14</sup> The size distribution of firm bins is calibrated as follows. We approximate the distribution of exports with a lognormal distribution with mean  $\mu = 11$  and standard deviation  $\sigma = 3$ . This specification matches the mean exports of \$5.11 million and has a median exports of \$59,300, and does a good job in matching the Lorenz curve reported in Bernard, Jensen and Schott (2007). As it is well known, there is a striking skewness in the distribution of exports across firms.<sup>15</sup>

The balls-and-bins model predicts that 96 percent of the potential firm times country trade flows is going to be zero. This is very close to the 98 percent we see in the data. What about the distribution of firm zeros across destinations?

<sup>13</sup>Bernard, Jensen and Schott (2007), page 11.

<sup>14</sup>Bernard, Jensen and Schott (2007), Table 2.

<sup>15</sup>Note that this is conditional on having positive exports. In other words, we only try to explain the *allocation* of exporting firms across destination markets; we do not analyze the question of which firms export. That is done in Section IV.



For each country, we can calculate the expected number of non-empty firm bins. We can then regress (the log of) this number on GDP and distance.<sup>16</sup> Table 5 presents the results. For convenience, we reproduce the regression estimate by Bernard, Jensen, Redding and Schott (2007) in the first column.

The coefficient estimates in the simulated regression are similar to the ones in the actual data. Just as in the data, bigger, closer countries are served by more exporters: the more balls are thrown, the fewer bins will be left empty. Interestingly, the skewness in firm exports does not play as big a role as it did for product bins: given that there are so many, most firm bins are going to remain empty anyway.<sup>17</sup>

	Log number of exporting firms	Log number of non-empty bins
Log GDP	0.71*** (0.04)	0.60*** (0.03)
Log distance	-1.14*** (0.16)	-0.89*** (0.13)
Observations	175	184
$R^2$	0.74	0.72

TABLE 5—EXPORTING FIRMS AND GRAVITY – *Balls and bins*

### III. Firm-level export patterns

We now turn to evidence on the extensive margin at the level of individual exporting firms. In this section we ask how many products firms export and how many destinations they serve. Note that the universe of interest is the set of *exporting firms*, because the empirical facts are usually reported only for firms that have some exports.<sup>18</sup> This way we can use the balls-and-bins model to understand these moments and abstract, for now, from the split between exporters and non-exporters.

The key facts about the extensive margin at the firm level are that while most firms export a single product to a single country, the bulk of exports is done by multi-product, multi-destination exporters.<sup>19</sup>

<sup>16</sup>We take GDP (in current-price USD) from the World Development Indicators. We take distance from the bilateral distance dataset of CEPII.

<sup>17</sup>We calibrated firm bins to the distribution of overall sales in manufacturing, which resulted in 93 percent of firm-country bins remaining empty and a 0.60 elasticity of the number of firms exporting to a country with respect to country size. When using 167,217 symmetric firm bins, we got 82 percent empty bins and an elasticity of 0.72. We also explored alternative distribution, like a Pareto, with similar results.

<sup>18</sup>Though export datasets can be merged with domestic data such as in Bernard, Jensen, and Schott (2007) and Eaton, Kortum and Kramarz (2004).

<sup>19</sup>The following facts are for U.S. merchandise trade in 2002, reported in Bernard, Jensen, Redding and Schott (2007), Table 4.

To start with, 42 percent of the firms export only a single product, defined by the 10-digit HS code. While being a little less than half of the total firms, they account for a tiny fraction of total exports, 0.4 percent.

**Empirical regularity 3.** *42 percent of firms export a single product (defined as a 10-digit HS code). These firms account for only 0.4 percent of exports.*

A similar pattern exists for firms that export to a single country. These firms account for a little less than two thirds of the total, but still amount to a small fraction of total exports.

**Empirical regularity 4.** *64 percent of firms export to a single country. These firms account for only 3.3 percent of exports.*

But perhaps the most striking fact corresponds to the fraction of firms that export a single product to a single country. These firms represent 40 percent of the total exporters yet account only for a minuscule 0.2 percent of total exports.

**Empirical regularity 5.** *40 percent of firms export a single product to a single country. These firms account for only 0.2 percent of total exports.*

Let us turn now to the calibration of our benchmark. The 10-digit HS codes are calibrated to the aggregate export share of each HS code shipments in total U.S. exports in 2005. The size of each country bin is calibrated to the share of that country shipments in total U.S. export flows.<sup>20</sup> The following table lists the five biggest country bins.

Country	Share
Canada	0.341
Mexico	0.189
Japan	0.041
United Kingdom	0.035
Germany	0.030

TABLE 6—THE FIVE BIGGEST COUNTRY BINS

We assume each firm has a different number of export balls. Because we do not have data on the number of shipments at the firm level, we calibrate the number of balls to the distribution of exports across firms. As we did earlier, we approximate the distribution of exports with a lognormal distribution with  $\mu = 11$  and  $\sigma = 3$ . Corresponding to the average size of export shipments in 2000, we take each \$36,000 of export sales to represent one ball, rounding up. Because of

<sup>20</sup>The assumption here is that the structure of aggregate exports did not change too much between 2002 and 2005. Results are virtually unchanged if we use value shares instead. We maintain the ongoing assumption of no systematic relationship between categories.

the extreme skewness in the distribution of exports by firm, many firms will end up with just one export ball.

The predicted fraction of single-product exporters is 43 percent. This is very close to the actual fraction in the data (42 percent). The predicted fraction of exports coming from single-product producers is 0.3 percent, close to the actual 0.4 percent. Let us see how the balls-and-bins model manages to reproduce the fraction of single-product exporters with such precision. In the model practically all single-product exporters have only one ball. This is because with 8,867 HS codes, the second ball is very likely to fall into an HS category different from the first one. Only 0.3 percent of two-ball exporters are single-product exporters. The key to understanding the incidence of single-product exporters is that there are plenty of very small exporters.

The model underpredicts the data with respect to the fraction of single-country exporters: 45 percent in the model for 64 percent in the data. The reason is that the fraction of single-country exporters falls sharply with firms with the second and third balls. For example, the model predicts that only 11 percent of firms with two shipments export both of them to Canada (and 4 percent to Mexico). In order to match this fact, a structural model will require a mechanism that increases the fraction of relatively large exporters that export only to Canada (and possibly Mexico). We view this an excellent example of how the balls-and-bins benchmark can highlight as informative statistics that were not obviously relevant *ex ante*.

Last but not least, balls and bins is right on the spot with respect to the fraction of single-product, single-country exporters, and the small fraction of exports that they account for. Note that a fraction of 40 percent of single-product, single-country exporters implies that most single-product exporters are also single-country exporters, and vice versa. Is this surprising? The balls-and-bins model makes it clear that the fact follows from the presence of many small exporters. Almost all single-product exporters have only one ball, and these are all going to be single-country exporters. And this exactly what we see in the data. The conditional probability of single-country exporters among single-product exporters is 99.9 percent in the model, close to the 96 percent in the data.

Our results suggest that the skewness of the exporter distribution is key to understanding the split between single-destination, single-product firms and the rest. In particular, the left tail of the export distribution—the small exporters—is what enables the balls-and-bins model to match the data. This property of the distribution is not specific to exporters. For example, our results do not change when we calibrate the model to match the observed skewness in *domestic* sales for the U.S.<sup>21</sup> In contrast, the balls-and-bins model underpredicts the data once we censor the left tail. Interestingly, the right tail properties of the distribution have little bearing on the results as virtually all firms selling more than \$100,000 are

<sup>21</sup>The web Appendix B contains several alternative calibrations of firm size skewness.

predicted to be multi-country, multi-product exporters. We thus conclude that trade models capable of matching the fraction of small exporters in the data will also be able to reproduce the firm-level export patterns discussed here.<sup>22</sup> As we shall see in the next Section, the split between exporters and non-exporters is not due to sparsity. There are thus strong economic forces, yet to be fully understood, shaping the distribution of exporters and the firm-level facts discussed here.

#### IV. Exporting firms

We now move on to the differences between exporting and non-exporting firms. It is a well-established fact that exporters are few in number and they are significantly larger than non-exporting firms.

According to Bernard, Jensen, Redding and Schott (2007), only 18 percent of manufacturing firms export at all. The fraction drops to about 3 percent when all firms outside manufacturing are included.<sup>23</sup> Other studies have likewise confirmed the scarcity of exporters. Plant-level statistics also fall in the same pattern. For the quantitative exercise, we stay with the fraction of exporters among U.S. manufacturing firms.

**Empirical regularity 6.** *Exporters are few — only 18 percent of manufacturing firms export in the U.S.*

The second fact is that exporters sell significantly more than non-exporters — about 4.4 times more, according to Bernard, Jensen, Redding and Schott (2007). Again, firms outside manufacturing and plant-level evidence reveal similar patterns.

**Empirical regularity 7.** *Exporters are large — among U.S. manufacturing firms, exporters sell 4.4 times more than non-exporters.*

That exporters are few and they are larger than non-exporters have been confirmed in other datasets, in other settings, and with other measures of size.

We follow essentially the same steps as before to map the model to the data. The key difference is that now the output flow will include total sales, not only exports. We thus need data on total sales per firm in order to construct the distribution of balls ( $\pi_n$ ). Using publicly available data from the Statistics of U.S. Businesses of the Census for year 2002, we approximate the distribution of firm sales by a lognormal distribution with  $\mu = 13.2$  and  $\sigma = 2.66$ . This corresponds to median sales of \$680,000 and average sales of \$13.2 million. As it is well known, there is enormous skewness in the size distribution of firms. Whereas 59 percent of firms sell less than \$1 million, the average firm sells \$13.2

<sup>22</sup>Most of the literature has not paid much attention to small exporters, with the exception of Arkolakis (2010).

<sup>23</sup>See Table 2 in Bernard, Jensen, Redding and Schott (2007). The data is from the 2002 Economic Census.

million.<sup>24</sup> In the 2002 Economic Census, there were 297,873 manufacturing firms. As before, we obtain the number of balls  $n$  per firm by dividing its total sales by \$36,000 and rounding up.<sup>25</sup>

To distinguish between exporters and non-exporters we only need two bins: one for domestic sales, the other for foreign sales. Total receipts amounted to \$3.94 trillion for manufacturing firms in the 2002 Economic Census. Exports of manufactured goods amounted to \$545 billion in 2002.<sup>26</sup> That is, 13.9 percent of manufacturing receipts came from exports. This pins down the size of the domestic bin at 0.861 and the size of the export bin at 0.139.

We find that exporters are much less common in the data than in the model: 74 percent of the manufacturing firms should be exporting according to the balls-and-bins model, compared to 18 percent in the data. Clearly, the scarcity of exporters is not a byproduct of sparsity and is thus vindicated in its status as a key stylized fact in the trade literature.

It is easy to see why the model overpredicts the fraction of exporters. The probability that a firm with  $n$  balls of total sales does not export is

$$(1 - s)^n = 0.86^n.$$

Among the smallest firms, that is, with one ball, 14 percent of them export. This is already a very high number given that only 18 percent of total manufacturing firms export. It obviously gets worse. Since where each ball ends up is independent of the distribution of existing balls, each \$36,000 has quite a high chance of ending up going to a foreign market. Almost half of the firms with a paltry \$100,000 of total sales should export. A median firm has a 95 percent chance to export. It is clear that this is not the case in the data: exporting is a more unlikely event than the balls-and-bins model would indicate.

The unconditional probability of exporting is convex in the fraction of exports,  $s$ , so if there is heterogeneity across industries, the aggregate economy will contain fewer exporters than predicted by the average  $s$ . However, at the 3-digit level, this heterogeneity is rather small, and does not change the exporting probability substantially.

The model's prediction for the exporter's size premium is also off. Surprisingly, though, the model *overpredicts* the size of exporters. That is, despite exporters being four fifths of total firms in the model for one fifth in the data, the model predicts that exporters are 34 times larger than non-exporters on average, while in the data they are "only" 4.4 times larger. In terms of the exporter size premium, in log sales, the difference in the model is 3.53, for 1.48 in the data.<sup>27</sup>

<sup>24</sup>We also experimented with fitting a Pareto distribution with similar results.

<sup>25</sup>In the previous section we used evidence on the average shipment value to pin down the "ball size." We have no direct equivalent for total sales.

<sup>26</sup>Bureau of the Census, FT-900, "International Trade in Goods and Services." We converted all figures to 2000 dollars.

<sup>27</sup>In the web Appendix C we formally derive the exporter's size premium and include a parametric example.

To understand why exporters are larger under balls-and-bins than in the data, note that balls-and-bins implies that the largest firms export with a probability close to one. Even the median firm that has \$660,000 dollars in sales, corresponding to 18 balls, exports with probability 0.93. The skewness of the firm sales distribution then implies that the average firm in the top half of the distribution is much larger than any of the non-exporters, who mainly come from the bottom half. The fact that the size premium is smaller in the data suggests the data has a weak sorting of exporters by size: exporters are smaller, not larger, than expected. In other words, there has to be a substantial fraction of very large firms that do not export – in contrast with the model.

## V. The implications of balls and bins for trade theory

In this Section we back up the claim that if a fact is matched by the balls-and-bins benchmark, then there is a wide range of theoretical models which are also consistent with that fact. First we show how to generate predictions for finite-sample data from any structural model. We also show that the baseline calibration for the balls-and-bins exercise is nested as a special case in most structural trade models—and thus they are all capable of matching a fact if it is matched by the balls-and-bins model. Identification problems, though, can also arise within a single structural model class: we briefly discuss the case of models with economies of scale. We refer the reader to Appendix D for a step-by-step description on how to nest the balls-and-bins within an extended version of Helpman, Melitz, and Rubinstein (2008).

### A. Model predictions for finite-sample data

Most trade models take the form of a set of continuous trade flows, that is, if we were to evaluate the models at different frequencies, the predicted export flows would just scale up or down proportionately with the period length, akin to oil flowing through a pipeline at a constant rate.

The data, however, consists of a finite number of observations, corresponding to the transactions in a given time period, usually a year. As several trade models have proven their worth on many dimensions, we aim to generate finite-sample predictions from any model in a parsimonious manner. Our approach is simply to sample the model for  $n \in \mathbb{N}$  observations, ideally replicating the number of transactions or shipments in the data, and re-interpreting the model's predicted market shares as the likelihood that any given transaction belongs to a particular category.

To fix ideas, we focus the discussion on trade flows at the country-product level. Let superscript  $m \in M$  index a model class, say, a Ricardian or a Melitz model, that predicts trade revenues  $R^m(X_j, X_g, \theta)$  for country  $j$  and product  $g$  given some variables  $X_j$  and  $X_g$  as well as a vector of parameters  $\theta$  from a set of admissible values  $\Theta^m$ . The variables  $X_j$  and  $X_g$  may be equilibrium

variables obtained elsewhere in the model or data on distance, factor intensity, other covariates, or perhaps fixed-effects for countries and products.<sup>28</sup> The market share of each country-product pair in model  $m$  is simply

$$s_{jg}^m = \frac{R^m(X_j, X_g, \theta)}{\bar{R}^m}$$

where  $\bar{R}^m$  is the sum of all trade flows.

We let the model's market share for a country-product pair  $j, g$  to pin down, quite naturally, the likelihood that a transaction belongs to such pair. Assuming that transactions are independent of each other, the formulas from Section I carry on with the whole vector of market shares in place of the bin size distribution. For example, the probability that a country-product pair is observed in model  $m$  for a finite sample of size  $n$  is simply  $1 - (1 - s_{jg}^m)^n$ .

Note that the model may predict an empty trade flow and thus a country-product pair may have zero probability. We call this a “fundamental zero.” However, we also know a trade flow may go missing in the sample even if it has positive probability—a “sample zero.” In a dense data set,  $n \rightarrow \infty$ , the realized share of shipments in any category will converge almost surely to the market shares predicted in the model and only fundamental zeros would remain zeroes. In other words, we would recover exactly the underlying continuous flows from the structural model. On the other hand, if the data are sparse, that is, if the number of observations is low relative to the level of detail we wish to analyze, sample zeros will be pervasive and realized market shares will have some sampling variation.

### B. Nesting the balls-and-bins baseline calibration

Throughout the paper we have assumed that there is no systematic relationship across categories in order to construct a parsimonious benchmark. It turns out this assumption is usually nested as a particular parameterization in most standard trade models. Formally, if there exists a parameter vector  $\theta_{bb}^m \in \Theta^m$  such that all trade flows can be written as multiplicatively separable in product and country variables,

$$R^m(X_j, X_g, \theta_{bb}^m) = d_j^m(X_j) d_g^m(X_g),$$

then when  $\theta = \theta_{bb}^m$ , the product and country assignments become independent, that is, the likelihood that a shipment is classified in product  $g$  and country  $j$  is simply  $s_{jg}^m = s_j^m s_g^m$ , where  $s_j^m$  and  $s_g^m$  are the marginal probabilities. The above property also typically allows the model  $m$  to match the distribution of trade flows across products and across countries—and we then retrieve *exactly* the baseline calibration for the balls-and-bins.

<sup>28</sup>The exact variables are expected to vary across models: we omit the explicit dependence on the model  $m$  to keep the notation as concise as possible.

The condition of multiplicatively-separable trade flows is satisfied by the basic trade model featuring an Armington-type, constant-elasticity-of-substitution demand system. The latter is ubiquitous as a building block in richer trade models so it is usually possible to find parameters such that trade flows are multiplicatively separable and thus nests our baseline calibration. For example, a trade model based on comparative advantage will boil down to the multiplicatively-separable specification when the distribution of unit costs across sectors is equated across countries—that is, when there is no systematic relationship between countries and products. There is also no question trade models in use have proven to be very successful at matching quantitatively the marginal distributions across categories, for example, by reproducing the gravity equation.

### C. Identification

We now close our argument. If two classes of models  $m, m'$  allow for trade flows to be multiplicatively separable, then neither of them will do worse than the balls-and-bins model at matching a given data moment in a finite sample. If the balls-and-bins model matches the data moment, then both models will inherit the ability to do so—and hence the data moment is not useful to distinguish the models. Formally, let  $\psi^m(\theta)$  denote the predicted (theoretical moment) from a class  $m$  of models, parametrized by the vector  $\theta$  within the admissible set of parameters  $\Theta^m$ . The corresponding data moment is  $\psi^*$ . If a class of models admits a parametrization  $\theta_{bb}^m$  such that trade flows are multiplicatively separable, then the best estimate in each model  $\theta^m$  cannot be further from the data moment than under  $\theta_{BB}^m$ , that is,

$$|\psi^m(\theta^m) - \psi^*| \leq |\psi^m(\theta_{bb}^m) - \psi^*| = |\psi^{bb} - \psi^*|$$

where  $\psi^{bb}$  is the balls-and-bins model's prediction for that moment. Then, if both  $m, m'$  nest the balls-and-bins baseline calibration, their best fit for moment will be similar, since by the triangle inequality,

$$|\psi^m(\theta^m) - \psi^{m'}(\theta^{m'})| \leq |\psi^m(\theta_{bb}^m) - \psi^*| + |\psi^{m'}(\theta_{bb}^{m'}) - \psi^*| \leq 2|\psi^{bb} - \psi^*|.$$

That is, the difference between the best of the two class of models is smaller than twice the difference between the balls-and-bins predictions and the data. Whenever balls and bins match the data well and  $|\psi^{bb} - \psi^*|$  is small, the two models cannot be distinguished by their ability to reproduce the aforementioned data moment.

Unfortunately, moments in sparse datasets may not be informative even *within* the strict confines of a single model class, that is, model  $\psi^m(\theta)$  and data  $\psi^*$  may be very close for a large range of parametrizations  $\theta \in \Theta^m$ . The next subsection illustrates this problem for trade models with economies of scale.



## D. Zeroes and economies of scale

Consider a model class  $m$  such that trade flows below a certain size, pinned down by the model parameter  $\theta$ , never take place. While admittedly stylized, this specification naturally captures the extensive-margin implications of a wide array of models with economies of scale. Such models are parametrized by a threshold  $t \geq 0$  such that if a bin is smaller than  $t$ , that bin will remain empty.<sup>29</sup> Let  $\theta^m$  denote the number of fundamental zeros; the number of bins that are smaller than  $t$ ,  $\theta^m = \max\{i : s_i \leq t\}$ . Within this class of models the parameter  $\theta^m$  fully describes the model. Clearly,  $\theta_{bb}^m = 0$  retrieves the balls-and-bins benchmark.

The moment we are interested in,  $\psi^m$ , is the number of zero trade flows. The prediction of a model with parameter  $\theta^m$  with respect to the total zeros is given by

$$\psi^m(\theta^m) = \sum_{i=1}^{\theta^m} 1 + \sum_{i=\theta^m+1}^K (1 - s_i)^n,$$

where, without loss of generality, we have sorted the bins in increasing size. The difference between this model prediction and the prediction of balls and bins is

$$\psi^m(\theta^m) - \psi^m(\theta_{bb}^m) = \sum_{i=1}^{\theta^m} [1 - (1 - s_i)^n].$$

The formula has a simple interpretation: it is the number of bins among  $\{1, 2, \dots, \theta^m\}$  that are expected to be non-empty in *the balls-and-bins baseline*. As these are exactly the smallest bins, they are bound to be a sample zero with high probability in the balls-and-bins model. As a result, across a large range of thresholds that possibly close many bins there will be almost no change in total zeros as we are trading fundamental for sample zeros virtually one to one. But the abundance of sample zeros is, in a nutshell, the reason that the balls-and-bins model is capable of generating a large number of total zeros and match the data.

Note that if the data were dense, that is, the number of shipments would be very large, total zeros would increase one-to-one with the number of fundamental zeros, as  $\lim_{n \rightarrow \infty} \psi^m(\theta^m) = \theta^m$ . In this case, total zeros would be perfectly informative about fundamental zeros, and there would be no problem identifying the relevant model of economies of scale.

We briefly report here numerical results for product-level zeros and the prevalence of exporting firms. Let us start with product-level zeros. The balls-and-bins model predicted that 72.1 percent of product-country pairs would be empty, that is,  $\psi^m(0) = 0.721$ . If we set  $\theta^m$  to close *half* of all bins, the prediction for total zeros barely budges to 72.8 percent.<sup>30</sup> That is, even though the model closes more

<sup>29</sup>The threshold is stated directly in terms of bin size, but we can always scale the threshold units to shipments by multiplying  $t$  by the number of shipments  $n$ , and then to dollars by multiplying it with the average shipment value.

<sup>30</sup>As in Section II.A, we condition on the flow of shipments per country. The results are very similar

than a million country-product pairs, the total number of zeros only changes by 0.7 percent. Simply put, the balls-and-bins model matches the data because the vast majority of bins are expected to be empty: even the median bin has less than one chance in a hundred to receive a ball. Closing these bins has virtually no impact on the total number of product-level zeros, and thus the latter cannot be used to identify  $\theta^m$  precisely, or, more broadly, to estimate the fixed costs of exporting with any precision.

The conclusion is drastically different in the exercise regarding exporting firms. As soon as we start closing exporting bins, the share of exporters drops very fast. Recall that under the benchmark balls-and-bins calibrations, 74 percent of firms were exporters, that is, only 26 percent were non-exporters. By shutting down no more than one fifth of the exporting bins the share of exporters drop below 70 percent. From then on, the share of exporters drops virtually one-to-one with the share of fundamental zeros. For example, closing 30 percent more bins we find that the fraction of exporters drops below half. In stark contrast with the previous exercise, the model predictions react sharply to threshold values, so we can take the fraction of exporters as an informative moment of the magnitude of the underlying economies of scale.<sup>31</sup>

## VI. Conclusion

Categorical datasets *do* contain a lot of information, even if they are sparse. We provided a benchmark to discern which statistics are driven by the sparsity—and thus contain little information about the extensive margin—and those that are not—and require a model to posit the correct joint distribution across categories in order to reproduce the fact. We evaluated several stylized facts commonly cited in the trade literature and found a mixed picture. Facts relating the frequency and pattern of zeros in product-country trade flows or to multi-destination, multi-product exporters do not differentiate among any model that can reproduce the marginal distributions, i.e., across products, countries or firms. On the other hand, the split between exporters and non-exporters is clearly not driven by sparsity and thus informative of the underlying mechanism driving the export participation decision.

We hope that our approach can be used in future empirical work using massive micro-level trade datasets. Recent transaction-level datasets are very detailed, and trade flows are typically broken down by firms, 8 or 10-digit product codes, and destination countries.<sup>32</sup> By their very nature, these datasets are *sparse* in

if we do not condition or even after re-calibrating bin sizes to respect the aggregate distribution of sales over products.

<sup>31</sup>In Appendix D we confirm that the export participation margin pins down the firm-level fixed cost of exporting; but the fraction of zeros in country-product flows is not informative even within the strict confines of a structural model.

<sup>32</sup>Bernard, Jensen and Schott (2007) describe the customs dataset of the U.S.; Eaton, Kortum and Kramarz (2004) for France; Mayer and Ottaviano (2007) for Belgium; Damijan, Polanec and Prasnikar (2004) for Slovenia; Halpern, Koren and Szeidl (2011) for Hungary; Eaton, Eslava, Kugler and Tybout (2007) for Colombia.

the sense that the number of observations is low with respect to the number of categories of interest.

## REFERENCES

- Agresti, Alan. 2002. *Categorical Data Analysis*, Second Edition, John Wiley and Sons. Hoboken, NJ.
- Anderson, James E. and Eric van Wincoop. 2003. "Gravity with Gravitas: A Solution to the Border Puzzle", *American Economic Review* 93(1), 170–192, March.
- Arkolakis, Costas. 2010. "Market Penetration Costs and the New Consumers Margin in International Trade", *Journal of Political Economy*, 118(6):1151–1199, December.
- Baldwin, Richard and James Harrigan. 2011. "Zeros, Quality, and Space: Trade Theory and Trade Evidence", *American Economic Journal: Microeconomics*, 3(2):60–88, May.
- Bernard, Andrew B., Jonathan Eaton, J. Bradford Jensen and Samuel Kortum. 2003. "Plants and Productivity in International Trade", *American Economic Review* 93(4):1268–1290, September.
- Bernard, Andrew B. and J. Bradford Jensen. 1999. "Exceptional Exporter Performance: Cause, Effect, or Both?", *Journal of International Economics* 47(1):1–25, February.
- Bernard, Andrew B., J. Bradford Jensen, Stephen J. Redding, and Peter K. Schott. 2007. "Firms in International Trade", *Journal of Economic Perspectives* 21(3):105–130, Summer.
- Bernard, Andrew B., J. Bradford Jensen and Peter K. Schott. 2007. "Importers, Exporters and Multinationals: A Portrait of Firms in the U.S. that Trade Goods", in Timothy Dunne, J. Bradford Jensen and Mark J. Roberts (eds.), *Producer Dynamics: New Evidence from Micro Data*.
- Damijan, Jozse P., Saso Polanec and Crt Kostevc. 2010. "From Innovation to Exporting or Vice Versa?", *World Economy* 33(3):374–398, March.
- Deardorff, Alan V. 1998. "Determinants of Bilateral Trade: Does Gravity Work in a Neoclassical World?," in Jeffrey Frankel (ed.) *The Regionalization of the World Economy*, University of Chicago Press.
- Eaton, Jonathan, Marcela Eslava, Maurica Kugler, and James Tybout. 2007. "Export Dynamics in Colombia: Firm-Level Evidence", NBER Working Paper No. 13531.
- Eaton, Jonathan, Samuel Kortum, and Francis Kramarz. 2004. "Dissecting Trade: Firms, Industries, and Export Destinations", *American Economic Review* 94(2):150–154, May.
- Eaton, Jonathan, Samuel Kortum, and Francis Kramarz. 2011. "An Anatomy of

- International Trade: Evidence from French Firms”, *Econometrica*, 79(5):1453–1498, September.
- Ellison, Glenn and Edward L. Glaeser. 1997. “Geographic Concentration in U.S. Manufacturing Industries: A Dartboard Approach”, *Journal of Political Economy* 105(5):889–927, October.
- Evenett, Simon J. and Wolfgang Keller. 2002. “On Theories Explaining the Success of the Gravity Equation”, *Journal of Political Economy* 110(2):281–316, April.
- Ghosh, Sucharita and Steven Yamarik. 2004. “Are Regional Trading Arrangements Trade Creating? An Application of Extreme Bounds Analysis”, *Journal of International Economics* 63(2):369–395, July.
- Halpern, László, Miklós Koren and Adam Szeidl. 2011. “Imported Inputs and Productivity”, Working Paper.
- Helpman, Elhanan, Marc J. Melitz and Yona Rubinstein. 2008. “Estimating Trade Flows: Trading Partners and Trading Volumes”, *Quarterly Journal of Economics* 123(2):441–487, May.
- Haveman, Jon and David Hummels. 2004. “Alternative hypotheses and the volume of trade: the gravity equation and the extent of specialization”, *Canadian Journal of Economics* 37(1):199–218, April.
- Hornok, Cecília and Miklós Koren. 2013. “Per-Shipment Costs and the Lumpiness of International Trade”, Working paper.
- Hummels, David and Peter J. Klenow. 2005. “The Variety and Quality of a Nation’s Exports”, *American Economic Review* 95(3):704–723, June.
- Keller, Wolfgang. 1998. “Are International R&D Spillovers Trade-Related? Analyzing Spillovers among Randomly Matched Trade Partners”, *European Economic Review* 42(8):1469–1481.
- Mayer, Thierry and Gianmarco Ottaviano. 2007. *The Happy Few: The Internationalization of European Firms*, Bruegel Blueprint Series. Volume III.
- Melitz, Marc J. 2003. The Impact of Trade on Intra-industry Reallocations and Aggregate Industry Productivity, *Econometrica* 71(6), 1695–1725.
- Schaefer, Kurt C., Michael A. Anderson and Michael J. Ferrantino. 2008. “Monte Carlo Appraisals of Gravity Model Specifications”, *Global Economy Journal*, Berkeley Electronic Press, 8(1):1–26, February.