



Discussion Papers in Economics

TRADE, GRAVITY AND AGGREGATION

By

Holger Breinlich

(University of Surrey),

Dennis Novy

(University of Warwick)

&

Joao M.C. Santos Silva

(University of Surrey).

DP 07/21

School of Economics

University of Surrey

Guildford

Surrey GU2 7XH, UK

Telephone +44 (0)1483 689380

Facsimile +44 (0)1483 689548

Web <https://www.surrey.ac.uk/school-economics>

ISSN: 1749-5075

Trade, Gravity and Aggregation*

Holger Breinlich[†]
University of Surrey

Dennis Novy[‡]
University of Warwick

J.M.C. Santos Silva[§]
University of Surrey

13 September 2021

Abstract

Gravity regressions are a common tool in the empirical international trade literature and serve an important function for many policy purposes. We study to what extent micro-level parameters can be recovered from gravity regressions estimated with aggregate data. We show that estimation of gravity equations in their original multiplicative form via Poisson pseudo maximum likelihood (PPML) is more robust to aggregation than estimation of log-linearized gravity equations via ordinary least squares (OLS). In the leading case where regressors do not vary at the micro level, PPML estimates obtained with aggregate data have a clear interpretation as trade-weighted averages of micro-level parameters that is not shared by OLS estimates. However, when regressors vary at the micro level, using disaggregated data is essential because in this case not even PPML can recover parameters of interest. We illustrate our results with an application to Baier and Bergstrand's (2007) influential study of the effects of trade agreements on trade flows. We examine how their findings change when estimation is performed at different levels of aggregation, and explore the consequences of aggregation for predicting the effects of trade agreements.

JEL classification: C23, C43, F14, F15, F17

Keywords: Free trade agreements, gravity equation, OLS, PPML, trade costs

*We gratefully acknowledge research support from the Economic and Social Research Council (ESRC grant ES/P00766X/1). We thank conference and seminar participants at the European Trade Study Group conference 2019 and the London School of Economics for helpful comments.

[†]Holger Breinlich, School of Economics, University of Surrey, Guildford GU2 7XH, UK, CEP/LSE and CEPR. Email: h.breinlich@surrey.ac.uk.

[‡]Dennis Novy, Department of Economics, University of Warwick, Coventry CV4 7AL, UK, CEP/LSE, CEPR and CESifo. Email: d.novy@warwick.ac.uk.

[§]João Santos Silva, School of Economics, University of Surrey, Guildford GU2 7XH, UK. Email: jmcss@surrey.ac.uk.

1 Introduction

Gravity equations are the workhorse model in international trade to estimate trade cost parameters and to evaluate the effects of policy changes. Gravity equations have been used to estimate the trade effects of free trade agreements, currency unions, WTO membership and colonial history, amongst other institutional features (see Anderson, 2011, and Head and Mayer, 2014). When trade costs change, the impact typically materializes at the level of individual agents such as firms and customers. That is, the impact is governed by parameters at the micro level, for example product-level demand elasticities. However, due to data constraints, gravity equations are routinely estimated at the aggregate level using country-level data.¹ The resulting estimates are often assumed, explicitly or implicitly, to be informative about the more fundamental micro-level parameters.

In this paper, we investigate to what extent this practice is justified. Specifically, we ask two related questions. First, can we infer micro-level elasticities and other parameters from aggregate-level gravity regressions and, if yes, under what conditions? Second, if the elasticities estimated with aggregate data differ from the micro elasticities, what are the implications of aggregation for the use of gravity equations in evaluating policy changes?²

Regarding the first question, we show that if there is no parameter heterogeneity at the micro level and if regressors do not vary across micro units within more aggregated macro units, we can recover micro-level parameters from estimation based on aggregate data. For instance, this scenario applies when the regressors are bilateral distance and a bilateral free trade agreement (FTA) dummy whose elasticities do not differ at the micro level, and the fixed effects included in the regressions also do not vary at the micro level. Although ordinary least squares (OLS) estimation with aggregate data is in principle able to recover the micro-level elasticities, this is only possible under very restrictive assumptions that are unlikely to hold in practice, for example because errors are heteroskedastic, leading to aggregation bias. In contrast, we show that Poisson pseudo maximum likelihood (PPML)

¹Throughout the paper we use the terms “micro”, “product” and “sector” level interchangeably. This is in contrast to the aggregate country pair-level analysis of bilateral trade.

²Although in this paper we only explicitly consider aggregation across sectors, aggregation over time raises similar issues and our results can easily be applied to that problem.

estimation recovers the micro-level elasticities under more realistic assumptions allowing for heteroskedastic errors. This invariance result extends the well-known finding by Santos Silva and Tenreyro (2006) to the problem of aggregation.

For the more general case where parameters vary at the micro level, naturally it is not possible to recover the micro parameters from aggregate data. However, we show that results from aggregate regressions can to some extent still be informative about micro-level parameters. Specifically, if there is micro-level parameter heterogeneity but regressors do not vary at the micro level, PPML estimation using aggregate data will recover a trade-weighted average of micro-level parameters. By contrast, the more traditional OLS estimation of log-linearized gravity equations leads to results that are not interpretable because they combine the aggregation bias with the bias resulting from log-linearization (see Santos Silva and Tenreyro, 2006). Moreover, this second bias also varies with the level of aggregation, making the OLS estimates very unstable.

If the regressors vary at the micro level, estimates obtained with aggregate data will generally yield biased estimates of the underlying micro-level parameters. We show that it is generally not possible to determine the sign or size of this bias, unless one is willing to make very specific assumptions about the underlying data generating process. Therefore, if the model contains regressors that vary across products, there is effectively no alternative but to use disaggregate data because in this case even the PPML estimates are difficult to interpret.

Having established these theoretical results, we investigate the implications of aggregation for the evaluation of trade policy, focusing on the classic question of the impact of free trade agreements. Specifically, we ask whether and how the level of aggregation is important for predictions regarding the impact of trade agreements on trade flows. Consistent with our theoretical results, we demonstrate that if the conditions required for invariance hold and the corresponding gravity equations are estimated with PPML, then the predicted trade flow increase does not vary with the level of aggregation at which gravity equations are estimated. However, if we use OLS estimation or if the invari-

ance conditions are violated, the predicted increase in trade flows can vary substantially depending on whether we use micro-level or macro-level data.

Combining the insights obtained from this empirical exercise with our theoretical results, we then formulate a set of recommendations for applied researchers on how to interpret and use estimates obtained from gravity equations. In essence, we recommend that PPML and micro-level data should be employed whenever possible. If the model contains regressors such as tariffs that vary by product, there may be no alternative but to use disaggregate data. However, in the leading case where the regressors do not vary by sector, the PPML estimates obtained with aggregate data still have a clear and interesting interpretation as trade-weighted averages of micro-level parameters.

Our work contributes to several related strands in the literature. First, we contribute to the econometrics literature on cross-sectional aggregation of non-linear economic models. For example, Lewbel (1992) and van Garderen, Lee and Pesaran (2000) study the consequences of aggregation in the context of log-linearized constant-elasticity models. They conclude that the least squares estimator of aggregated log-linearized constant-elasticity models is consistent only under very strong conditions that are unlikely to hold in many applications. These conditions are even stronger than the ones needed for the OLS estimator of the log-linear model to be valid (see Santos Silva and Tenreyro, 2006). By contrast, we consider the effects of aggregation when constant-elasticity models are estimated in their exponential form by PPML and find that aggregate parameter estimates are often informative.

Second, our work is related to papers in the trade and international macro literature concerned with learning about micro-level parameters when using more aggregate trade data. For example, Helpman, Melitz, and Rubinstein (2008) show how to account for the self-selection of firms into export markets when estimating aggregate gravity equations. The procedure they propose yields consistent estimates of important micro-level parameters even when only aggregate data is available. However, their procedure makes strong assumptions about the data generating process underlying the observed aggregate trade flows such as CES demand, monopolistic competition and, for some of their results,

Pareto-distributed firm-level productivities.³ By contrast, we provide results for the wider class of log-linear models that nest the CES monopolistic competition model considered by Helpman, Melitz, and Rubinstein (2008) as a special case.

Our work is also related to Imbs and Mejean (2015) who argue that smaller trade elasticity estimates at the aggregate level are an artefact of aggregation, driven by heterogeneity bias as sectoral elasticities are constrained to be homogeneous.⁴ We show that aggregation can lead to different elasticity estimates even in the absence of heterogeneous elasticities at lower levels of aggregation.⁵ Redding and Weinstein (2019a and 2019b) show that the theoretical aggregation of gravity equations is not straightforward. The key tension is that the typical gravity equation is log-linear while aggregation of trade flows implies summation in levels. We focus on the empirical estimation and the underlying econometric theory.

In independent work, French (2017) also considers the aggregate estimation of bilateral gravity equations in the presence of micro-level heterogeneity. He approaches aggregate estimation as a problem of omitted variable bias if micro-level heterogeneity is not properly accounted for. This view is consistent with our results, but we provide a different analytical framework that provides an interpretation of aggregation effects that is arguably more practical and intuitive. In particular, we show that in leading cases, aggregate PPML estimates of gravity coefficients can be seen as weighted averages of product-specific parameters. This interpretation emphasizes the relation between the estimated parameters and the parameters of interest.

The paper is structured as follows. In Section 2 we present a simple international trade model that delivers gravity equations at two different levels of aggregation. This framework provides theoretical guidance for our approach and helps to clarify the link

³See also Santos Silva and Tenreyro (2015).

⁴See also Pesaran and Smith (1995).

⁵Imbs and Mejean (2015) show that the bias typically pushes estimates towards zero in an international trade context. All else equal, trade costs are a lesser impediment to less elastic trade flows characterized by lower elasticities, and therefore a larger weight is placed on less elastic products. Feenstra, Luck, Obstfeld and Russ (2018) specify a monopolistic competition model with a separate ‘macro’ elasticity between home and foreign varieties and a ‘micro’ elasticity between different foreign varieties. We do not make such a distinction but rather focus on the variation of elasticities across products.

between parameter estimates and the underlying theoretical parameters common in international trade models. In Section 3 we present initial motivating evidence regarding the effects of aggregation in gravity estimation. In Section 4 we explain these findings by deriving a number of theoretical results on the aggregation of constant-elasticity models under different assumptions. In Section 5 we apply these results to the question of how aggregation bias can matter for predicting the effects of policy interventions, with a focus on the effects of free trade agreements. Drawing on the previous sections, Section 6 makes a set of recommendations for applied researchers on how to interpret estimates obtained from gravity equations at different levels of aggregation. Section 7 concludes.

2 Gravity at different levels of aggregation

We sketch a theoretical framework that yields gravity equations at different levels of aggregation. It is based on a simple model of international trade with a two-tier nested constant elasticity of substitution (CES) demand system. The upper tier represents the aggregate level of the economy, and the lower tier the disaggregated (sector/industry/product) level. Varieties in each sector are differentiated by origin according to the Armington assumption.

Aggregate consumption at the upper tier is given by

$$C_j = \left(\sum_s (c_{js})^{\frac{\nu-1}{\nu}} \right)^{\frac{\nu}{\nu-1}},$$

where c_{js} is real consumption by country j of sector s aggregates, and ν is the elasticity of substitution between sectors. The lower-tier aggregator is given by

$$c_{js} = \left(\sum_i (\theta_{ijs} c_{ijs})^{\frac{\sigma_s-1}{\sigma_s}} \right)^{\frac{\sigma_s}{\sigma_s-1}},$$

where c_{ijs} is real consumption by country j of sector s varieties originating from country i , σ_s is the elasticity of substitution across sector s varieties, and $\theta_{ijs} \geq 0$ is a taste

parameter that implies zero trade flows between countries i and j in sector s if $\theta_{ijs} = 0$ (see Redding and Weinstein, 2019b).

The CES demand relationship at the lower tier follows as

$$x_{ijs} = \left(\frac{p_{ijs}}{\theta_{ijs} P_{js}} \right)^{1-\sigma_s} E_{js}, \quad (1)$$

where x_{ijs} denotes nominal trade flows from country i to country j in sector s , and p_{ijs} denotes their unit price. P_{js} is the sectoral CES price index in country j , and E_{js} is the corresponding sectoral expenditure. We assume trade costs are of the iceberg type such that

$$p_{ijs} = \tau_{ijs} p_{is}, \quad (2)$$

where p_{is} denotes the price (or unit cost) at origin i . We assume a standard log-linear specification of the trade cost function

$$\ln \tau_{ijs} = \rho_s \ln dist_{ij}, \quad (3)$$

where for simplicity we use bilateral distance $dist_{ij}$ as the sole trade cost component with an elasticity ρ_s that can vary by sector.⁶

Combining equations (1) and (2), we can write the micro-level gravity equation in log-linearized form as

$$\ln x_{ijs} = \phi_{is} + \xi_{js} - (\sigma_s - 1) \ln \tau_{ijs} + \ln \eta_{ijs}, \quad (4)$$

where the sector-origin fixed effect ϕ_{is} captures the origin price p_{is} , and the sector-destination fixed effect ξ_{js} captures the price index P_{js} and expenditure E_{js} . The error term $\ln \eta_{ijs}$ absorbs the idiosyncratic taste parameter θ_{ijs} and is traditionally assumed to be independent of trade costs τ_{ijs} . If trade costs are not fully observed, then $\ln \eta_{ijs}$ could also be seen as measurement error in trade costs.

⁶This trade cost function can be extended to a bilateral FTA_{ij} dummy and also to sector-specific components such as tariffs $tariff_{ijs}$.

Aggregate bilateral trade is defined as the sum of bilateral trade flows at the sector level

$$x_{ij} \equiv \sum_s x_{ijs} \quad (5)$$

with $i \neq j$. We proceed to show that based on the micro-level framework in equations (1)-(4), an aggregate gravity equation can be constructed but only with non-standard properties. For this purpose, we substitute the demand function (1) into the definition of aggregate bilateral trade (5) using equation (2):

$$\begin{aligned} x_{ij} &= \sum_s \left(\frac{\tau_{ijs} p_{is}}{\theta_{ijs} P_{js}} \right)^{1-\sigma_s} E_{js} \\ &= \left(\frac{\tau_{ij} p_i}{P_j} \right)^{1-\sigma} E_j \exp(\varepsilon_{ij}), \end{aligned} \quad (6)$$

where σ denotes the aggregate demand elasticity and

$$\exp(\varepsilon_{ij}) = \sum_s \left(\frac{\tau_{ij} p_i}{P_j} \right)^{\sigma-1} \left(\frac{\tau_{ijs} p_{is}}{\theta_{ijs} P_{js}} \right)^{1-\sigma_s} \frac{E_{js}}{E_j}. \quad (7)$$

Taking logarithms of equation (6) implies

$$\ln x_{ij} = \Phi_i + \Xi_j - (\sigma - 1) \ln \tau_{ij} + \varepsilon_{ij}, \quad (8)$$

with $\Phi_i = (1 - \sigma) \ln p_i$ and $\Xi_j = (\sigma - 1) \ln P_j + \ln E_j$.

Superficially, equation (8) has the same structure as a conventional log-linearized gravity equation. But the key point is that ε_{ij} should not be considered a standard error term because it is by construction a function of bilateral trade costs τ_{ij} .⁷ The exception is the case where θ_{ijs} is the only source of sectoral heterogeneity, and therefore $\sigma_s = \sigma$, $p_{is} = p_i$,

⁷The result in equation (8) resonates with Redding and Weinstein (2019a and 2019b) who also demonstrate that in a nested CES demand system as above, a log-linear gravity equation can be derived at the aggregate level but only with an error term that is not orthogonal to bilateral trade costs.

$P_{js} = P_j$, $E_{js} = E_j$, and $\tau_{ijs} = \tau_{ij} = dist_{ij}^\rho$. In this special case we have

$$\exp(\varepsilon_{ij}) = \sum_s \theta_{ijs}^{\sigma-1},$$

and therefore (8) is a proper log-linearized gravity equation. This is a result we will use later.

3 Motivating evidence

The theoretical framework in Section 2 delivers a gravity equation (4) at the disaggregate (i.e., micro) level as well as an equation (8) at the aggregate level that can be construed as a non-standard gravity equation. We now explore empirically how estimated coefficients on gravity variables behave at different levels of aggregation. As an illustration in this section and later in the paper, we use a replication of results from Baier and Bergstrand's (2007) seminal work on the effects of free trade agreements.

Baier and Bergstrand's empirical framework is based on an OLS regression of logarithmic trade flows on multiple categories of fixed effects and binary dummies for whether two countries have a trade agreement in place. Specifically, they consider models of the form

$$\ln x_{ijt} = \alpha_{it} + \alpha_{jt} + \alpha_{ij} + \beta_1 FTA_{ijt} + \beta_2 FTA_{ijt-1} + \beta_3 FTA_{ijt-2} + \varepsilon_{ijt}, \quad (9)$$

where x_{ijt} are imports of country j from country i in period t , the FTA dummies (which include a contemporaneous term as well as two lags to allow for phasing-in effects) are the regressors of interest, α_{it} and α_{jt} denote exporter-year and importer-year fixed effects that control for price index and expenditure terms, α_{ij} are bilateral fixed effects introduced by Baier and Bergstrand to help address the potential endogeneity of free trade agreements, and ε_{ijt} is the error term.

Santos Silva and Tenreyro (2006) show that, in general, OLS estimation of log-linear gravity models such as specification (9) leads to biased estimates of the elasticities, and

that estimation by PPML (see Gourieroux, Monfort and Trognon, 1984) using trade flows in levels solves this problem. Therefore, we also estimate by PPML models of the form

$$x_{ijt} = \exp(\alpha_{it} + \alpha_{jt} + \alpha_{ij} + \beta_1 FTA_{ijt} + \beta_2 FTA_{ijt-1} + \beta_3 FTA_{ijt-2}) \eta_{ijt}, \quad (10)$$

where η_{ijt} is the multiplicative error term.

Baier and Bergstrand estimate (9) using aggregate (i.e., country-level) bilateral trade data from the IMF’s Direction of Trade Statistics (DOTS). In our context, the key question is how the coefficient estimates on the FTA terms change as we vary the level of aggregation, and how these changes depend on the estimator used. For this purpose we will replicate Baier and Bergstrand’s key results but using data from the UN Comtrade database which provides trade flows at different levels of aggregation.⁸ This allows us to show results from estimating (9) and (10) at three different levels of aggregation: aggregate bilateral imports, and imports at the 2-digit and 4-digit SITC levels.⁹

Baier and Bergstrand only estimate models using aggregate data and therefore do not have a sector dimension to their fixed effects. However, when estimating (9) and (10) at different levels of aggregation we need to decide whether the fixed effects are allowed to vary by sector. The inclusion of exporter-year and importer-year fixed effects that vary by sector is necessary for theoretical consistency as price indices and expenditure levels generally vary at the sector level (see equation 4). By contrast, there is a less compelling case to allow the pair fixed effects to vary by sector, but arguably there are pair-specific

⁸Baier and Bergstrand use data from the IMF’s DOTS for the years 1960-2000 at five-year intervals for 96 countries, excluding zero trade flows. To achieve a similar timespan, we rely on data from Comtrade, based on the SITC classification, for the same countries and for the years 1962, 1965, 1970, ..., 2000 (no data is available prior to 1962 so we use 1962 data for 1960). Specifically, we use the value of bilateral imports in current US dollars on a c.i.f. basis. These data are available at five different levels of aggregation, from SITC 4-digit to the country-level bilateral trade flows used by Baier and Bergstrand (SITC 0-digit). Our data on FTAs are the same as in Baier and Bergstrand (2007), based on their Table 3.

⁹For each country pair in the data, we observe trade flows for 61 2-digit SITC sectors and 625 4-digit SITC sectors. However, we drop all 4-digit sectors with fewer than 2,000 observations of positive trade flows. This is done because in the sector-level regressions discussed below, it is not always possible to identify all the parameters of interest when the number of positive observations is small. To keep our sample comparable across the different sections of this paper, we also exclude such observations for the pooled regressions presented here. Dropping these observations reduces the number of 4-digit sectors to 576 and that of 2-digit sectors to 60.

trade cost elements that do not vary over time but vary by sector, justifying the inclusion of time-invariant sectoral importer-exporter fixed effects.

The results obtained when estimating models (9) and (10) at different levels of aggregation are presented in Table 1. Specifically, for each estimator we present results at three levels of aggregation and, when using disaggregated data, we present results for models imposing that the fixed effects are the same across sectors and for models where the fixed effects are allowed to vary by sector.¹⁰ Although we present estimates for each of the three FTA dummies for the cases considered, we will focus our discussion on the total FTA effect, which is reported in the last column of the table and is computed as the sum of the estimated coefficients on the three FTA dummies.

Reassuringly, results for the specification most directly comparable to Baier and Bergstrand’s (the one using OLS estimation with aggregate trade and without sector-level fixed effects) are similar to theirs. We obtain a total FTA effect of 0.714 log points (see the last column of the first line of Table 1) compared to 0.76 log points in the key specification by Baier and Bergstrand (the one reported in their Table 5, column 4). This demonstrates that changing the data source from the IMF DOTS database to UN Comtrade does not in itself change the basic findings in Baier and Bergstrand (2007).¹¹

After this initial check, we now turn to our question of how results change as we change the degree of aggregation in our data. The results in Table 1 show that OLS estimates are sensitive to the level of aggregation, irrespective of the fixed effects we use. For example, when we use 4-digit trade flows, the estimated total effect in the specifications without

¹⁰Specifically, models with sector-level fixed effects include importer-year-sector, exporter-year-sector and exporter-importer-sector fixed effects, whereas models without sector-level fixed effects include only importer-year, exporter-year and importer-exporter fixed effects.

¹¹Note, however, that the two samples are not fully comparable, as Baier and Bergstrand use log exports as their dependent variable and thus have to exclude observations with zero bilateral trade flows from their sample. By contrast, the results in Table 1 are based on a fully rectangularised set of bilateral trade flows following current best practice in applied international trade research (see, e.g., Yotov, Piermartini, Monteiro and Larch, 2016). That is, we fill in all missing country pair-product-year combinations and assign a trade flow value of zero for all such “filled in” observations. While the additional zero observations get dropped when taking logs (as we do for our OLS specifications), rectangularisation also changes the structure of the lags of the FTA regressors included, making the two datasets (Baier and Bergstrand’s and our rectangularised data) incompatible. As an additional comparability check, we have also re-estimated Baier and Bergstrand’s key specification on our sample without zero trade flows, obtaining a total FTA effect of 0.77 log points, which is very similar to the total FTA effect of 0.76 log points estimated by Baier and Bergstrand (results available on request).

sector-level fixed effects is 0.481, which is substantially lower than the aggregate effect of 0.714 obtained with aggregate data (see lines 1 and 3 in the OLS panel in Table 1). Note that this is despite the fact that we impose a common coefficient across sectors for each of the FTA dummy variables, hence assuming that the FTA effect is the same for all sectors, and that the FTA dummy itself does not vary with aggregation, i.e., it is the same for every sector for a given country pair.

Table 1: Regression results at different aggregation levels

Estimator	Aggregation level	Sector-level fixed effects	Regressor coefficients (Standard errors clustered by pair)			
			FTA_t	FTA_{t-1}	FTA_{t-2}	Total
OLS	Aggr. trade	No	0.174 ^{***} (0.0453)	0.379 ^{***} (0.0455)	0.161 ^{***} (0.0510)	0.714 ^{***} (0.0608)
	SITC 2-digit	No	0.355 ^{***} (0.0212)	0.191 ^{***} (0.0206)	-0.004 (0.0246)	0.542 ^{***} (0.0307)
	SITC 4-digit	No	0.285 ^{***} (0.0165)	0.161 ^{***} (0.0140)	0.036 ^{**} (0.0180)	0.481 ^{***} (0.0268)
	SITC 2-digit	Yes	0.334 ^{***} (0.0225)	0.168 ^{***} (0.0209)	0.064 ^{***} (0.0255)	0.566 ^{***} (0.0326)
	SITC 4-digit	Yes	0.326 ^{***} (0.0195)	0.116 ^{***} (0.0159)	0.087 ^{***} (0.0199)	0.529 ^{***} (0.0314)
PPML	Aggr. trade	No	0.278 ^{***} (0.0324)	0.224 ^{***} (0.0242)	0.089 ^{***} (0.0267)	0.591 ^{***} (0.0455)
	SITC 2-digit	No	0.278 ^{***} (0.0324)	0.224 ^{***} (0.0242)	0.089 ^{***} (0.0267)	0.591 ^{***} (0.0455)
	SITC 4-digit	No	0.278 ^{***} (0.0324)	0.224 ^{***} (0.0242)	0.089 ^{***} (0.0267)	0.591 ^{***} (0.0455)
	SITC 2-digit	Yes	0.235 ^{***} (0.0252)	0.179 ^{***} (0.0171)	0.118 ^{***} (0.0217)	0.533 ^{***} (0.0338)
	SITC 4-digit	Yes	0.210 ^{***} (0.0270)	0.172 ^{***} (0.0160)	0.117 ^{***} (0.0201)	0.500 ^{***} (0.0353)

Notes: The table presents the results of estimating three-way gravity equations. The dependent variable is logarithmic trade in the OLS panel and the level of trade in the PPML panel. Models with sector-level fixed effects (“Yes”) include importer-year-sector, exporter-year-sector and exporter-importer-sector fixed effects, whereas models without sector-level fixed effects (“No”) include only importer-year, exporter-year and importer-exporter fixed effects; *** $p < 0.01$, ** $p < 0.05$, * $p < 0.1$.

In contrast to the OLS estimates, Table 1 shows that with PPML the estimated coefficients and standard errors are invariant to aggregation if we do not allow the fixed effects to vary by sector. This invariance result is illustrated by the fact that the estimates and standard errors are exactly the same for the first three rows of the PPML panel in Table 1, which of course also implies that the estimated total FTA effect is the same at 0.591.

Moreover, when PPML is used, estimation with aggregate data is equivalent to using disaggregated data to estimate models where the parameters are assumed to be the same for all sectors. This can be seen by noting that the change in the estimates and standard errors resulting from aggregation (e.g., going from the bottom row to the top row in the PPML panel) is the same as the change resulting from imposing coefficient homogeneity (e.g., going from the bottom row to the middle row in the PPML panel). That is, in this context the effect of aggregation is the same as the effect of removing the sector-level dimension of the fixed effects.

Naturally, when we include the sector dimension into the fixed effects, the PPML estimates do depend on the level of aggregation because changing the level of aggregation changes the model specification. However, even in this case the PPML estimates are less sensitive to changes in the level of aggregation, dropping from 0.591 at the aggregate level to 0.500 at the 4-digit level (a drop of 15%). By contrast, the corresponding OLS estimates change from 0.714 at the aggregate level to 0.529 at the 4-digit level (a drop of 26%).

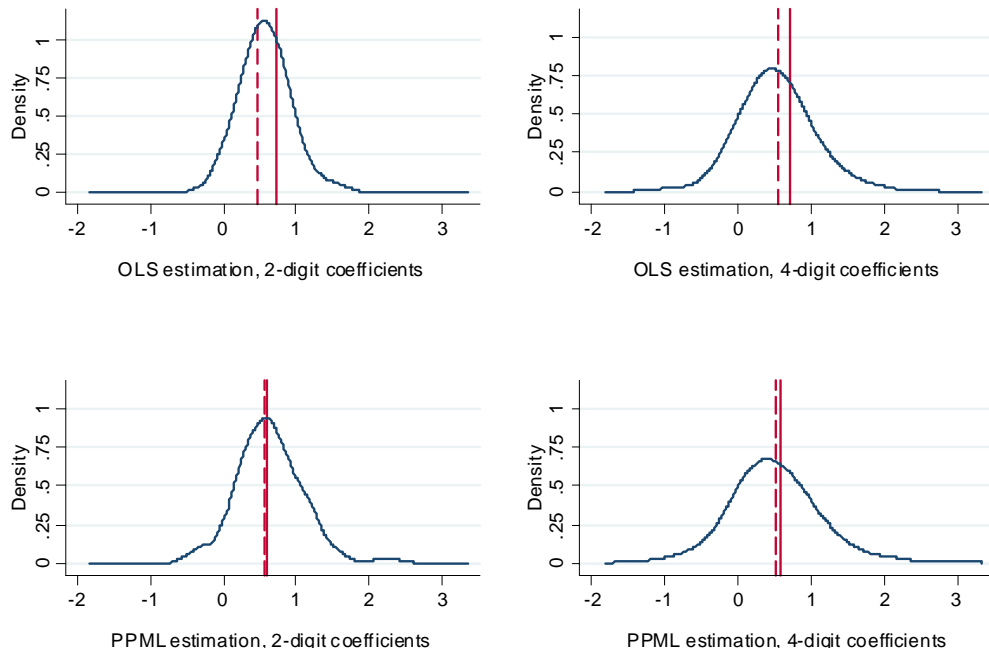
In the next section, we discuss from a theoretical perspective the pattern of results observed in this table. In particular, we will explore the fact that in the absence of sectoral heterogeneity PPML estimates are not affected by aggregation to develop a two-step approach to determine the consequences of aggregation in models with sectoral heterogeneity.

Until now, we have constrained coefficient on the FTA dummies to be the same across sectors. We now relax this restriction and estimate (9) and (10) separately for each SITC sector at both the 2-digit and the 4-digit levels of aggregation. That is, now both the

fixed effects and the coefficients on the FTA dummies are allowed to vary by sector. This yields 60 sets of estimates for each equation at the 2-digit level, and 576 sets at the 4-digit level.

Figure 1 presents kernel density estimates of the total estimated FTA effects obtained with OLS and PPML for each of the sectors at the 2-digit and 4-digit levels of aggregation. In each panel, we add two elements to help interpret the results. The vertical solid line represents the estimated effect obtained by estimating the models with the aggregate data (this corresponds to the results in the first line of each panel in Table 1). We also include a vertical dashed line representing the weighted average of the estimated coefficients using shares of 2-digit or 4-digit sectors in total trade as weights.

Figure 1: Kernel density plot of the estimated FTA effects at sectoral level



Notes: The dashed line is the trade-weighted average of the estimated sectoral effects, and the solid line is the effect estimated with aggregate data. The left panels show estimates at the 2-digit level, and the right panels show estimates at the 4-digit level. The top panels are estimated with OLS and the bottom panels with PPML.

Figure 1 shows that in all cases the aggregate estimates lie reasonably close to the mode of the distribution of the sector-level estimates. Moreover, when PPML is used, the aggregate estimates are very close to the weighted averages of the disaggregate estimates (i.e., the solid and dashed lines are very close to each other); the same does not necessarily happen when OLS is used. We will also explain these results in the next section.

4 Aggregation of constant-elasticity models

Given the initial motivating evidence from Section 3, we now examine the effects of aggregation from an econometric point of view. We draw a distinction between micro/sector-level *parameters* on the one hand (where coefficients vary at the sector level, for instance, by letting the coefficients on the FTA regressors from Section 3 vary across SITC sectors), and micro/sector-level *regressors* on the other hand (where regressors themselves differ across sectors, for example the fixed effects from Table 1 that have a sector dimension). This analysis provides insights into the coefficient patterns we should expect in the data at different levels of aggregation, hence helping to explain the results from the preceding section.

4.1 Set-up and aggregation with OLS

Let x_{ijs} denote the non-negative outcome of interest and let z_{ijs} be a vector of explanatory variables for observation $s = 1, \dots, l$ in the pair ij . As in equation (5), we define $x_{ij} = \sum_{s=1}^l x_{ijs}$ as the aggregate counterpart of x_{ijs} , but we are agnostic about how z_{ij} , the aggregate value of the vector of regressors, is obtained.

We are interested in estimating gravity equations at either the sector or the aggregate level. Consistent with our earlier theoretical framework (see equation 1), we assume that sector-level trade flows (x_{ijs}) are generated by the following constant-elasticity model

$$x_{ijs} = \exp(z'_{ijs}\beta_s) \eta_{ijs}, \quad (11)$$

where η_{ijs} is a non-negative error term such that $\mathbb{E}(\eta_{ijs}|z_{ijs}) = 1$, β_s is a vector of parameters (including fixed effects) that are potentially allowed to vary with s and in which the slope parameters have the usual interpretation as (semi-) elasticities. The traditional approach to estimating models such as (11) is to take logarithms of both sides and estimate

$$\ln x_{ijs} = z'_{ijs}\beta_s + \ln \eta_{ijs} \quad (12)$$

by least squares, under the assumption that $\mathbb{E}(\eta_{ijs}|z_{ijs})$ is constant.

Lewbel (1992) and van Garderen, Lee and Pesaran (2000) among others have studied the consequences of estimating models such as (12) using aggregate data. The model for the aggregate data is given by

$$x_{ij} = \sum_{s=1}^l x_{ijs} = \sum_s \exp(z'_{ijs}\beta_s) \eta_{ijs}, \quad (13)$$

which in general is not a constant-elasticity model and therefore cannot be log-linearized (see the discussion around equation 8 in Section 2). Therefore, estimation of the aggregate counterpart of the model in (12) will only identify the parameters of interest under very restrictive assumptions, even in the special case where the parameters and regressors do not vary with s . Lewbel (1992) provides details on the conditions needed for the parameters of interest to be identified from aggregated data.

One of the conditions considered by Lewbel (1992) is that the errors of the model in (12) need to be independent of the regressors. As pointed out by Santos Silva and Tenreyro (2006), this condition is very unlikely to hold in trade applications. When this condition is violated, the estimation of log-linearized models will produce biased estimates of the parameters of interest even when disaggregate data are used, as in (12). Because changing the level of aggregation inevitably changes the properties of the error term, the degree to which these conditions on the error terms are violated may depend on the level of aggregation considered. Therefore, in trade applications, aggregation of log-linearized models combines the bias due to aggregation with the bias resulting from log-linearization.

The two biases may (partially) offset or compound each other, making the results very hard to interpret.

To avoid the log-linearization bias, Santos Silva and Tenreyro (2006) recommended that models such as (11) should be estimated in their original multiplicative form using the PPML estimator of Gourieroux, Monfort and Trognon (1984). In the remainder of this section, we study the consequences of using aggregated data when the multiplicative model is estimated by PPML, and we show that this approach lessens or even eliminates the negative consequences of aggregation.

4.2 Aggregation with PPML

Building on our earlier distinction between parameters and/or regressors varying at the sector level s , we consider four particular cases of this problem, summarized in Table 2. Case 1 is the simplest scenario where neither regressors nor parameters vary with s . In Case 2 the parameters vary with s but regressors do not, and the reverse holds in Case 3. Finally, in Case 4 both parameters and regressors vary with s .

Table 2: Four cases of aggregation

		Sectoral parameters	
		No	Yes
Sectoral regressors	No	Case 1	Case 2
	Yes	Case 3	Case 4

4.2.1 Case 1: Parameters and regressors are constant

We start by considering the case where neither the regressors nor the parameters vary with s . We have seen in Section 2 that in this particular case we have proper gravity equations at the disaggregate and aggregate levels, and the same is found in our current set-up. Indeed, in this case equation (11) can be written as

$$x_{ijs} = \exp(z'_{ij}\beta) \eta_{ijs}, \quad (14)$$

and expression (13) becomes

$$x_{ij} = \sum_s \exp(z'_{ij}\beta) \eta_{ijs} = \exp(z'_{ij}\beta) \sum_s \eta_{ijs} = \exp(\ln l + z'_{ij}\beta) \eta_{ij}^*, \quad (15)$$

where $\eta_{ij}^* = l^{-1} \sum_{s=1}^l \eta_{ijs}$ is an error term such that $\mathbb{E}(\eta_{ij}^* | z_{ij}) = 1$. Therefore, as discussed in Section 2, in this particular case both x_{ijs} and x_{ij} are given by stochastic constant-elasticity models.

It is easy to show that the PPML estimates of the slopes in (14) and (15) are identical.¹² To see that this is true, notice that the first order condition of the PPML estimator of β in (14) is (see, e.g., Cameron and Trivedi, 2013)

$$S(\hat{\beta}) = \sum_{ijs} \left(x_{ijs} - \exp(z'_{ij}\hat{\beta}) \right) z_{ij} = 0,$$

where as usual a “hat” is used to denote parameter estimates and \sum_{ijs} is shorthand for $\sum_i \sum_j \sum_s$. This condition can be written as

$$S(\hat{\beta}) = \sum_{ij} \left(x_{ij} - \exp(\ln l + z'_{ij}\hat{\beta}) \right) z_{ij} = 0,$$

which is the first order condition of the PPML estimator of β in the aggregate model defined by (15). Hence, the estimation results are invariant to the level of aggregation of the data (with the exception of the intercept which is adjusted to reflect the number of sectors being aggregated). Moreover, if the dependent variable in the aggregate equation is the mean of x_{ijs} rather than its sum, the estimates are exactly the same at both levels, and the invariance result continues to apply.

It is interesting to note that when clustering is taken into account, the level of aggregation also does not matter for the significance of the estimates. Indeed, we can show that the cluster-robust estimate of the covariance matrix for the estimates from (14) is

¹²This result first appears in the simulation evidence reported by Amrhein and Flowerdew (1992).

identical to the estimate of the robust covariance matrix for the estimates in the aggregate equation when the dependent variable is the average of x_{ijs} over s .¹³

It is important to note that the results above are obtained under the assumption that the number of sectors l is the same for every ij pair. When that is not the case, the same result holds when the models include pair fixed effects that will absorb the differences in the number of sectors by pair. If the disaggregate model does not include pair fixed effects, the aggregate and disaggregate elasticity estimates will not be numerically identical in finite samples, but the two sets of estimates converge to the same limit if the aggregate model includes pair fixed effects. To simplify the exposition, in what follows we continue to assume that the number of sectors l is the same for every pair.¹⁴

In summary, when both the parameters and the regressors are constant across s , both x_{ijs} and x_{ij} are given by constant-elasticity models with the same parameters, and the PPML estimates and standard errors are invariant to the level of aggregation of the data. This contrasts sharply with the results on aggregation of log-linear models where aggregation generally leads to an inconsistent estimator of the parameters of interest, even if the regressors and parameters do not vary with s (see Lewbel, 1992, and van Garderen, Lee and Pesaran, 2000).

Looking back at our results from Section 3, note that the models underlying the estimates in the first three lines of the PPML panel of Table 1 fall into our Case 1 (neither parameters nor regressors vary with s). Thus, the invariance result just outlined explains why the estimates and standard errors obtained with these models are exactly the same.

¹³Details are available on request.

¹⁴Note that this assumption will hold in any fully rectangularised dataset such as the one we are using for our empirical illustrations (see footnote 11).

4.2.2 Case 2: Parameters vary with s but regressors do not

We now consider the case where parameters vary with s but the regressors do not. That is, the relevant model at the disaggregate level is

$$x_{ijs} = \exp(z'_{ij}\beta_s) \eta_{ijs}.$$

Clearly, now it is not possible to recover the individual parameters from aggregate data, but it is interesting to study what we estimate when using aggregate data. We approach this problem in two steps. We first examine the effect of ignoring the parameter heterogeneity with disaggregated data, and we then use the invariance result for Case 1 to find the effect of aggregation.

To see the effect of ignoring the parameter heterogeneity, write the first order conditions for the estimates of β_s with disaggregated data as

$$S_s(\hat{\beta}_s) = \sum_{ij} \left(x_{ijs} - \exp(z'_{ij}\hat{\beta}_s) \right) z_{ij} = 0, \quad s = 1, \dots, l.$$

Since we have that $S_s(\hat{\beta}_s) = 0$ for each s , for the full sample we have $\sum_s S_s(\hat{\beta}_s) = 0$. Imposing homogeneity we estimate a single parameter for all s , say $\hat{\beta}^r$, which by definition will satisfy $S(\hat{\beta}^r) = \sum_s S_s(\hat{\beta}^r) = 0$.¹⁵

To study the relation between $\hat{\beta}^r$ and $\hat{\beta}_s$, $s = 1, \dots, l$, we can use the mean value theorem to write

$$\sum_s S_s(\hat{\beta}_s) = \sum_s S_s(\hat{\beta}^r) - \sum_s H_s(\beta_s^*) (\hat{\beta}_s - \hat{\beta}^r)$$

with $H_s(\beta_s^*) = -\partial S_s(b) / \partial b|_{b=\beta_s^*}$, where β_s^* is a point between $\hat{\beta}_s$ and $\hat{\beta}^r$.

¹⁵But notice that $S_s(\hat{\beta}_r) \neq 0$.

As $\sum_s S_s (\hat{\beta}_s) = \sum_s S_s (\hat{\beta}^r) = 0$, we can write

$$\hat{\beta}^r = \left[\sum_s H_s(\beta_s^*) \right]^{-1} \sum_s H_s(\beta_s^*) \hat{\beta}_s, \quad (16)$$

and therefore $\hat{\beta}^r$ can be interpreted as an average of the estimates of β_s weighted by the matrices $H_s(\beta_s^*)$.¹⁶

Noting that $H_s(\beta_s^*) = \sum_{ij} (\exp(z'_{ij}\beta_s^*) z_{ij} z'_{ij})$, we can see that $H_s(\beta_s^*)$ is itself a weighted sum of $\exp(z'_{ij}\beta_s^*)$, where the weights do not depend on s . Because $\exp(z'_{ij}\beta_s^*)$ is closely related to the expectation of x_{ijs} , heuristically $\hat{\beta}^r$ can be interpreted as a weighted average of the estimates of β_s , giving more weight to the estimates from the subsamples where x_{ijs} tends to be larger.¹⁷

From the invariance result for Case 1, we know that estimating the model that imposes $\beta_s = \beta^r$ with aggregated data will only change the estimate of the intercept, and therefore the parameters estimated with aggregated data can also be interpreted as weighted averages of the estimates of the individual parameters, with weights given by $H_s(\beta_s^*)$. It also follows from these results that in Case 2 the aggregation bias is identical to the bias caused by imposing the restriction that the coefficients do not vary across sectors. In Case 2, the problem is therefore not so much aggregation but the impossibility to account for sector-level parameter heterogeneity when estimating with aggregated data.

¹⁶It is interesting to note that $\hat{\beta}^r$ can also be seen as a minimum distance estimator obtained as

$$\hat{\beta}^r = \arg \min_b \sum_s (\hat{\beta}_s - b)' H_s(\beta_s^*) (\hat{\beta}_s - b),$$

which is an optimal minimum distance estimator when $H_s(\beta_s^*)$ is proportional to the inverse of the covariance matrix of $\hat{\beta}_s$ as in the Poisson distribution.

¹⁷To illustrate this, consider the case where only the intercept of $\hat{\beta}_s$ varies with s . It is easy to show that in this case

$$\hat{\beta}^r = \frac{\sum_s \exp(\kappa_s) \hat{\beta}_s}{\sum_s \exp(\kappa_s)},$$

where κ_s denotes a point between the intercepts in $\hat{\beta}_s$ and $\hat{\beta}^r$, and therefore in this particular case $\hat{\beta}^r$ is a weighted average of the individual estimates, giving more weight to the estimates from the subsamples where κ_s is larger.

Note, however, that these results do not carry over to OLS estimation because in that case the estimates are not invariant to aggregation, even if the parameters do not vary at the micro level (see Case 1 above). Put differently, we can establish the effects of imposing coefficient homogeneity in the first step, but the resulting estimates will be changed again by aggregation in the second step. Moreover, as discussed in Section 4.1, the bias of the OLS estimator resulting from using logarithmic trade flows as the dependent variable will also vary with the level of aggregation. This bias can partially offset or compound the bias resulting from aggregation, and it is therefore very difficult to meaningfully compare OLS results obtained at different levels of aggregation.

Looking again at our findings from Section 3, the results just outlined help to make sense of the patterns observed in Table 1 and Figure 1. Specifically, a key insight from Case 2 is that with PPML the effect of aggregation can be interpreted as the result of imposing coefficient homogeneity across sectors, and that this result does not carry over to OLS estimates. This explains why in the PPML panel of Table 1 the estimates in the top row are identical to the ones in the middle row, and why the same does not apply to OLS.

Our results also explain why in Figure 1 the trade-weighted average of the sectoral estimates (indicated by dashed vertical line) is always close to the estimate obtained with aggregate data (see the solid vertical line in Figure 1), but only for the PPML estimates. As we have shown, the aggregate estimates obtained by PPML are a weighted average of the underlying sector-level estimates, with subsamples with more trade (larger x_{ijs}) being assigned larger weights. However, this result does not apply to OLS estimation, which is also illustrated by Figure 1.

4.2.3 Case 3: Regressors vary with s but parameters do not

We now consider the case where regressors vary with s but parameters do not. That is, the relevant model is

$$x_{ijs} = \exp(z'_{ijs}\beta) \eta_{ijs}. \quad (17)$$

As in Case 2, we start by considering the effects of ignoring the heterogeneity in the disaggregated data, and we then use the invariance result for Case 1 to find the aggregation effect. That is, we start by considering the effect of estimating

$$x_{ijs} = \exp(z'_{ij}\beta^a) \eta_{ijs}^a, \quad (18)$$

where z_{ij} is obtained by aggregating z_{ijs} , β^a denotes the parameters of the aggregate equation, and η_{ijs}^a is a non-negative error term whose properties are determined by how β^a is defined.

Letting $z_{ijs} = z_{ij} + \varepsilon_{ijs}$, we can write equation (17) as

$$x_{ijs} = \exp(z'_{ij}\beta + \varepsilon'_{ijs}\beta) \eta_{ijs}, \quad (19)$$

and we can then interpret (18) as resulting from omitting $\varepsilon'_{ijs}\beta$ from (19).¹⁸ The effects of omitted variables are well understood in the context of linear regression models, but general results are difficult to obtain for non-linear models (see, e.g., Kiefer and Skoog, 1984, Neuhaus and Jewell, 1993, and Drake and McQuarrie, 1995). To gain some insight into the effect of omitting $\varepsilon'_{ijs}\beta$ we can start by writing

$$\mathbb{E}[x_{ijs}|z_{ijs}] = \mathbb{E}[x_{ijs}|z_{ij}, \varepsilon_{ijs}] = \exp((z'_{ij} + \varepsilon'_{ijs})\beta),$$

from where we obtain

$$\mathbb{E}[x_{ijs}|z_{ij}] = \mathbb{E}_{\varepsilon_{ijs}}[\exp((z'_{ij} + \varepsilon'_{ijs})\beta) | z_{ij}] = \exp(z'_{ij}\beta) \mathbb{E}_{\varepsilon_{ijs}}[\exp(\varepsilon'_{ijs}\beta) | z_{ij}]. \quad (20)$$

Equation (20) makes clear that, as is well known, the PPML estimator of (18) is consistent for the slope parameters in (17) when $\varepsilon'_{ijs}\beta$ is independent of z_{ij} .¹⁹ Unfortunately, this

¹⁸Alternatively, ignoring that the regressors vary with s and estimating specification (18) instead of (17) could be interpreted as estimating a non-linear regression with errors-in-variables. However, this is not a case of classical measurement error, and therefore we cannot use most of the results in the literature on measurement error in non-linear models (see, e.g., Kukush, Schneeweis and Wolf, 2004, and Carroll et al., 2006).

¹⁹See, e.g., Gourieroux, Monfort and Trognon (1984) and Neuhaus and Jewell (1993).

is unlikely to be a realistic scenario, and we therefore need to consider less favorable situations.²⁰

As an illustrative example, it is useful to start by considering the case where, conditional on z_{ij} , $\varepsilon'_{ijs}\beta$ has a normal distribution with mean $z'_{ij}\mu$ and variance $z'_{ij}\omega$ (see Nakamura, 1990, for a related approach). In this case $\exp(\varepsilon'_{ijs}\beta)$ is log-normal with

$$\mathbb{E}_{\varepsilon_{ijs}} [\exp(\varepsilon'_{ijs}\beta) | z_{ij}] = \exp(z'_{ij}\mu + 0.5z'_{ij}\omega),$$

and therefore

$$\mathbb{E}[x_{ijs} | z_{ij}] = \exp(z'_{ij}\beta + z'_{ij}\mu + 0.5z'_{ij}\omega) = \exp(z'_{ij}\beta^a)$$

with $\beta^a = \beta + \mu + 0.5\omega$, which implies $\mathbb{E}[\eta_{ijs}^a | z_{ij}] = 1$. That is, in this example β^a is the vector of parameters in $\mathbb{E}[x_{ijs} | z_{ij}]$, whereas β is the vector of parameters in $\mathbb{E}[x_{ijs} | z_{ijs}]$.

More generally, and as illustrated by the example above, the difference between the parameters at the two levels of aggregation depends on how the conditional moments of $\exp(\varepsilon'_{ijs}\beta)$ are related to z_{ij} . An important implication of this result is that the elements of β^a can be smaller or larger (in absolute value) than the corresponding elements of β .²¹ Furthermore, assuming that the models include intercepts, we have that in both cases the residuals will have zero mean, with the residuals of the disaggregate model being orthogonal to z_{ijs} , while the residuals of the aggregate model are orthogonal to z_{ij} but only approximately orthogonal to the disaggregate regressors. That is, $\hat{\beta}^a$ is such that the fitted values of the aggregate model approximate some of the characteristics of the fitted values of the regression with disaggregate data, and in that sense $\hat{\beta}^a$ provides an approximation to $\hat{\beta}$.

Combining these results with those for Case 1, we can conclude that the effect of aggregation on the estimated elasticities will be the same as replacing z_{ijs} with z_{ij} in the model for disaggregate data. That is, the aggregate model will estimate β^a rather than

²⁰For example, higher average collected tariff rates tend to be associated with a higher variance across sectors (see, e.g., Pritchett and Sethi, 1994).

²¹This contrasts with the so-called attenuation bias caused by classical measurement error.

β ,²² and it is difficult to predict the magnitude and sign of the differences between the elements of the two vectors unless we have information on how the conditional moments of the omitted variable $\varepsilon'_{ijs}\beta$ vary with z_{ij} .

4.2.4 Models where only the fixed effects vary with s : Case 2 or Case 3?

In the models providing motivating evidence in Section 3, the sectoral fixed effects can be interpreted as a set of dummies that depend on s but with constant coefficients. Therefore, these models can be seen as examples of Case 3. This way of approaching the problem, treating the sector-specific fixed effects as regressors that vary with s , is similar to that of French (2017) in that he also establishes that the consequences of aggregation in this context are equivalent to the omission of a variable. However, because we know from Case 3 that it is difficult to guess how the omission of $\varepsilon'_{ijs}\beta$ will impact the elasticity estimates, this approach is not particularly useful because it does not provide much information on the relation between the estimated parameters and the parameters of interest.

Alternatively, the sectoral fixed effects in the models presented in Section 3 can be interpreted as a set of dummies whose coefficients vary with s . Therefore, in the leading case where these are the only regressors with sectoral variation, the model can be seen as an example of Case 2 where at least the coefficients on the fixed effects vary with s . This alternative interpretation is more useful because it follows from our previous results that in this case the PPML estimates of the aggregate model have a clear interpretation as a weighted average of the sector-specific parameters, something that is illustrated in Figure 1.²³

It is worth noting that, even if the fixed effects are the only coefficients that vary with s , neglecting this heterogeneity will impact the estimates of all coefficients because neglecting the sectoral variation of the fixed effects effectively restricts the set of fixed effects included

²²Except, of course, for the intercept.

²³Naturally, a comparable result is not available for OLS estimation because in that case aggregation leads to inconsistency even if the parameters and the regressors do not vary with s .

in the model and, consequently, alters the remaining coefficient estimates.²⁴ However, these are always weighted averages of the underlying micro parameters,²⁵ and this explains why the PPML estimates of the FTA effect in Table 1 vary with the aggregation level when the fixed effects vary by sector, and also why these variations are relatively minor.

In summary, the leading case where the fixed effects are the only variables to vary by sector can be reinterpreted as a situation in which the parameters vary by sector but the regressors do not, and this approach is more informative about the effects of aggregation. Naturally, this interpretation does not extend to the case where other variables vary by sector such as tariffs (see, e.g., Amiti, Redding, and Weinstein, 2019). Our results for Case 3 suggest that in such situations it is essential to use disaggregated data on trade flows and tariffs, as done by Amiti, Redding, and Weinstein (2019), since even with PPML the magnitude and the direction of the bias resulting from aggregation are unknown and difficult to interpret.

4.2.5 Case 4: Regressors and parameters vary with s

Finally we consider the case where both the regressors and the parameters vary with s . This problem can be addressed by combining earlier results. As we know from Case 3, the effect of replacing z_{ijs} with z_{ij} in the regressions for each s is that in each case we estimate a vector β^a that is an approximation to β_s . From Case 2 we know that imposing the same coefficients for all s will lead to a weighted average of these individual estimates. Finally, the invariance result for Case 1 shows that aggregation will only change the intercept. Therefore, in Case 4 we estimate a weighted average of the approximations to β_s .

Since estimating a weighted average of approximations to the true coefficient is unlikely to be useful in most practical applications, the implications of Case 4 are clear. If the regressors vary at the micro level (such as tariffs) and the coefficients on those regressors

²⁴For example, when we eliminate sector-level fixed effects from Table 1, we change the set of included fixed effects from importer-sector-year, exporter-sector-year and exporter-importer-sector fixed effects to importer-year, exporter-year and exporter-importer fixed effects (see the notes to Table 1).

²⁵As (16) makes clear, the weights in these averages are matrices, which is another way to see that the entire vector of estimates can be affected even if only a single coefficient varies with s .

are also likely to vary across products (e.g., because price elasticities vary across products), there is no alternative to using appropriately disaggregated data.²⁶

5 Implications for gravity-based forecasts: An application to free trade agreements

Given that coefficient estimates often depend on the level of aggregation at which they are estimated, an important question that arises is what the consequences are for the use of gravity equations for trade policy-related questions. One of the policy questions for which the gravity equation has been used extensively is the impact of free trade agreements on trade flows, which is of course the application that we use throughout this paper to illustrate our findings. Given the results from the previous sections, a natural next step is to ask whether aggregation also matters for predictions of the trade flow increases expected after the implementation of FTAs. That is, if we estimate (9) and (10) at the aggregate level and predict the total trade impact of a free trade agreement using the estimated coefficients, do we obtain different effects compared to the alternative of estimating the same equations at the sector level, predicting sector-level trade flow changes and then adding up to the country level? Put differently, would a researcher who has access to trade data at the sector level reach the same conclusion as another researcher who only has country-level trade data available?

In trying to answer this question, and as in Section 3, we consider again two estimation methods (OLS and PPML), three levels of aggregation, and models that impose coefficient homogeneity or allow the estimates to vary at the sector level.²⁷ Note that there are three

²⁶An example of work which fits the setting of Case 4 is Bas, Mayer and Thoenig (2017) who use firm-product-level trade flows combined with product-level tariff data to obtain price elasticity estimates that potentially vary by product. They regress firm-product-level trade flows on (product-level) tariffs separately for each of the products in their data. Our results for Case 4 suggest that there is no alternative to using such disaggregated data since relying on more aggregate data (e.g., by regressing bilateral trade flows on average bilateral tariffs) would render the resulting elasticity estimates uninformative, irrespective of whether PPML or OLS estimation is used. Note, however, that our results do not imply the necessity of firm-level data but only of data that have variation at the micro level of interest (in this case, the product level).

²⁷See the description of Figure 1 and Table 1 for details.

types of counterfactuals we can perform. First, we could ask by how much trade flows between existing FTA partners are higher because of the FTAs in place. Second, we might be interested in finding out by how much trade would be larger if countries without FTAs put such agreements in place. Third, we could consider the change in trade moving from a situation without FTAs to a situation with FTAs in place between all countries. Conceptually, the first counterfactual corresponds to the average treatment effect on the treated (ATT), the second captures the average treatment effect on the untreated (ATU), and the third captures the average treatment effect (ATE).²⁸

Denoting by $x_{ijst,1}$ the value of trade for country pair ij in sector s at time t in the presence of an FTA, and by $x_{ijst,0}$ the same flow in the absence of an FTA, the relevant counterfactuals are easily computed from the estimates of β_1 , β_2 and β_3 obtained either from (9) with OLS or (10) with PPML. For example, trade among FTA partners is simply the observed trade flow, $x_{ijst,1,FTA_{ijt}=1} = \exp\left(\hat{\alpha}_{ist} + \hat{\alpha}_{jst} + \hat{\alpha}_{ijs} + \hat{\beta}_{1s} + \hat{\beta}_{2s} + \hat{\beta}_{3s}\right) \hat{\eta}_{ijst}$, where we let the estimated parameters vary with s and we have assumed that the agreement is fully phased in (for the purpose of this illustration, we drop all pairs for which there is an FTA that is not fully phased in).²⁹ The (counterfactual) trade flow between the partners in the absence of an agreement is then given by $x_{ijst,0,FTA_{ijt}=1} = \exp\left(\hat{\alpha}_{ist} + \hat{\alpha}_{jst} + \hat{\alpha}_{ijs}\right) \hat{\eta}_{ijst} = x_{ijst,1,FTA_{ijt}=1} \times \exp\left(-\hat{\beta}_{1s} - \hat{\beta}_{2s} - \hat{\beta}_{3s}\right)$, which can be computed using data on actual trade flows and the coefficient estimates from (9) or (10) obtained at the disaggregated level. Likewise, current trade among non-FTA partners can be expressed as $x_{ijst,0,FTA_{ijt}=0} = \exp\left(\hat{\alpha}_{ist} + \hat{\alpha}_{jst} + \hat{\alpha}_{ijs}\right) \hat{\eta}_{ijst}$, and the (counterfactual) trade in the presence of an FTA would be $x_{ijst,1,FTA_{ijt}=0} = x_{ijst,0,FTA_{ijt}=0} \times \exp\left(\hat{\beta}_{1s} + \hat{\beta}_{2s} + \hat{\beta}_{3s}\right)$. Once we have computed these counterfactuals, we can calculate the implied percentage

²⁸Note, however, that we are interested in changes in total trade flows rather than the average change in bilateral flows. That is, if we allow for sectoral coefficient heterogeneity, estimates for sectors with more trade get more weight. We also note that we are only concerned with the direct trade cost effects (i.e., what Head and Mayer, 2014, call “the partial trade impact”), not the indirect general equilibrium effects that operate through price indices, income and expenditure.

²⁹Using standard Neyman–Rubin notation, the subscript FTA_{ijt} indicates whether or not countries i and j have an FTA in place at time t . Thus, $x_{ijst,1,FTA_{ijt}=1}$ is the trade flow with an FTA for country pair ij in sector s at time t , given that country pair ij has an FTA in place. Note that this is of course simply the observed trade flow. By contrast, $x_{ijst,0,FTA_{ijt}=1}$ is the trade flow for country pair ij in sector s at time t without an FTA, which is a counterfactual trade flow given there is currently an FTA in place.

changes in trade flows as

$$\begin{aligned} ATT &= \frac{\sum_{ijst} x_{ijst,1,FTA_{ijt}=1}}{\sum_{ijst} x_{ijst,0,FTA_{ijt}=1}} - 1 \\ &= \frac{\sum_{ijst} x_{ijst,1,FTA_{ijt}=1}}{\sum_{ijst} x_{ijst,1,FTA_{ijt}=1} \times \exp\left(-\hat{\beta}_{1s} - \hat{\beta}_{2s} - \hat{\beta}_{3s}\right)} - 1, \end{aligned}$$

$$\begin{aligned} ATU &= \frac{\sum_{ijst} x_{ijst,1,FTA_{ijt}=0}}{\sum_{ijst} x_{ijst,0,FTA_{ijt}=0}} - 1 \\ &= \frac{\sum_{ijst} x_{ijst,0,FTA_{ijt}=0} \times \exp\left(\hat{\beta}_{1s} + \hat{\beta}_{2s} + \hat{\beta}_{3s}\right)}{\sum_{ijst} x_{ijst,0,FTA_{ijt}=0}} - 1, \end{aligned}$$

and

$$\begin{aligned} ATE &= \frac{\sum_{ijst} x_{ijst,1}}{\sum_{ijst} x_{ijst,0}} - 1 = \frac{\sum_{ijst} [x_{ijst,1,FTA_{ijt}=1} + x_{ijst,1,FTA_{ijt}=0}]}{\sum_{ijst} [x_{ijst,0,FTA_{ijt}=1} + x_{ijst,0,FTA_{ijt}=0}]} - 1 \\ &= \frac{\sum_{ijst} \left[x_{ijst,1,FTA_{ijt}=1} + x_{ijst,0,FTA_{ijt}=0} \times \exp\left(\hat{\beta}_{1s} + \hat{\beta}_{2s} + \hat{\beta}_{3s}\right) \right]}{\sum_{ijst} \left[x_{ijst,1,FTA_{ijt}=1} \times \exp\left(-\hat{\beta}_{1s} - \hat{\beta}_{2s} - \hat{\beta}_{3s}\right) + x_{ijst,0,FTA_{ijt}=0} \right]} - 1, \end{aligned}$$

where the summations are over all country pairs ij , sectors s and time periods t in our data.³⁰

Table 3 presents the results of this exercise. The first row of the table shows the predicted increases in trade flows when we estimate our FTA coefficients using country-level data (i.e., 0-digit). As there is no sector-level dimension, we have that $ATT=ATU=ATE$. We also have that $ATT=ATU=ATE$ whenever we impose coefficient homogeneity. The reason is obvious on inspection of the relevant expressions above. Indeed, if the coefficient

³⁰Since we compute treatment effects as percentage changes, the above definition of the ATE yields the same results as the more traditional ATE definition in terms of the average effect of a treatment (here: the presence of an FTA) across the units in a population (here: all country pairs, sectors and time periods) when the effect is expressed relative to the average baseline trade flows without FTAs. To see this write

$$ATE = \frac{\sum_{ijst} (x_{ijst,1} - x_{ijst,0})}{\sum_{ijst} x_{ijst,0}} = \frac{\sum_{ijst} x_{ijst,1} - \sum_{ijst} x_{ijst,0}}{\sum_{ijst} x_{ijst,0}} = \frac{\sum_{ijst} x_{ijst,1}}{\sum_{ijst} x_{ijst,0}} - 1.$$

estimates do not vary by sector s , the trade flow terms in the numerator and denominator of the ATT and ATU cancel so that the estimated treatment effect is simply the exponential of the sum of the coefficients for both the ATT and ATU. Since the ATE is simply a weighted mean of the ATT and the ATU, it will also be equal to whatever value the ATT and ATU take.

Table 3: Estimated treatment effects at different aggregation levels

Aggregation level	Heterogeneous coefficients	Treatment effect type	Estimator	
			OLS	PPML
SITC 0-digit	No	ATT=ATU=ATE	104.2%	80.7%
SITC 2-digit	Yes	ATT	56.3%	66.2%
SITC 2-digit	Yes	ATU	61.7%	83.4%
SITC 2-digit	Yes	ATE	60.5%	79.8%
SITC 2-digit	No	ATT=ATU=ATE	71.9%	80.7%
SITC 4-digit	Yes	ATT	61.0%	57.6%
SITC 4-digit	Yes	ATU	92.0%	105.4%
SITC 4-digit	Yes	ATE	85.3%	95.0%
SITC 4-digit	No	ATT=ATU=ATE	61.8%	80.7%

Notes: The table shows the predicted effect of FTAs at the 0-digit, 2-digit and 4-digit levels of aggregation. ATT is average treatment effect on the treated, ATU is average treatment effect on the untreated, ATE is average treatment effect. See text for details.

After these preliminary observations, we now move on to the more interesting comparison of how predicted trade flow increases vary with the level of aggregation and the underlying estimation method. Consistent with our results from Case 1, which demonstrated the invariance of PPML estimates when coefficient estimates do not vary at the sector level, Table 3 shows that the predicted trade flow increase under PPML with homogeneous estimates is the same regardless of whether we use aggregate, 2-digit or 4-digit data (it is always 80.7%). However, the same is not true for predictions based on OLS

estimates, even if we impose coefficient homogeneity. Specifically, with OLS the estimated ATE is 61.8% when using 4-digit data but 71.9% when using 2-digit data and 104.2% when using country-level data. Thus, using country-level data instead of 4-digit sector-level data can lead to substantially different predictions regarding the trade effects of FTAs when based on traditional OLS estimation.

As expected from our results for Case 2 above, however, aggregation matters even with PPML when the underlying sector-level elasticities are heterogeneous. Looking at the results in Table 3, when we use 4-digit data we estimate an ATE of 95.0%. When we instead use 2-digit data (allowing coefficient estimates to vary at that level), the estimated ATE is 79.8%. When we aggregate up further to bilateral trade at the country level, we obtain an ATE of 80.7% as mentioned previously.³¹ The corresponding results for OLS estimation are considerably more heterogeneous, with estimated ATEs of 85.3% when using 4-digit data, 60.5% when using 2-digit data and 104.2% when using aggregate data. This variability reflects the fact that, as noted before, OLS combines the aggregation bias and the bias resulting from log-linearization, and that these biases can partially offset or compound each other.

6 Aggregation in gravity equations: A practitioner's guide

Having systematically analyzed the issue of aggregation in gravity equations, we now present a number of recommendations for applied work resulting from our findings.

A first lesson is that clearly there will be situations where there is no good substitute for using disaggregated data. As we have shown, recovering micro-level elasticities from macro-level data will not be possible if these elasticities vary at the micro level, and forecasts based on elasticities estimated on the basis of macro data may prove inaccurate. However, even in this case our results suggest that using PPML rather than OLS esti-

³¹Note that with aggregate bilateral trade, there is of course no sector dimension and so we cannot allow for sector heterogeneity in our FTA estimates. Accordingly, Table 3 only reports results without coefficient heterogeneity at the aggregate (0-digit) level.

mation is preferable because PPML will recover a (trade-weighted) average of the true micro-level elasticities, whereas OLS estimates will be altogether uninformative.

The result that PPML estimation is to be preferred to OLS holds *a fortiori* when we expect no micro-level variation in the elasticities of interest. In this case, we have shown that PPML is able to recover the micro-level elasticities even with aggregate data, whereas OLS is not. But if the regressors vary at the sector level (as will be the case, for example, for bilateral tariffs), the interpretation of results is hard even with PPML because we can provide guidance on neither the sign nor the magnitude of the resulting bias, irrespective of whether the underlying elasticities also vary at the micro level. In such cases, as well as when the objective is to predict the effects of policy changes, there is no good alternative to estimating the corresponding models on the basis of micro-level data.

7 Conclusion

In this paper, we have investigated the consequences of aggregation for the estimation of gravity equations, using both PPML and OLS. We started by asking two related sets of questions. First, is it possible to infer micro-level elasticities and other parameters from aggregate-level gravity regressions? Second, what are the implications of aggregation for the use of gravity equations in evaluating policy changes? We provided motivating evidence on the consequences of aggregation using the classic question of the impact of free trade agreements on trade flows.

We then examined the aggregation properties of gravity equations from an econometric point of view, distinguishing four different cases. In the simplest case (Case 1), neither the regressors nor the parameters vary at the micro (i.e., product or sector) level. In Case 2, the parameters vary across products but the regressors do not. In Case 3, the regressors vary across products but the parameters do not. For Case 4, we assume that both the regressors and the parameters vary.

In Case 1, when gravity equations are estimated in the original multiplicative form with PPML, we obtain an invariance result. This means that aggregation is innocuous in the sense that the micro-level elasticities can be recovered using aggregate data. However, a comparable result is not generally available when gravity equations are estimated in their log-linear form using OLS, even when neither regressors nor parameters vary across sectors. These findings demonstrate that, in this particular case, the negative consequences of using aggregate data for the estimation of constant-elasticity models (such as the gravity model) are eliminated when the model is estimated in its multiplicative form by PPML.

When we allow the parameters to vary across sectors (Case 2), it is obviously impossible to recover the micro-level parameters. However, we showed that also in this case gravity estimation by PPML is more informative than OLS estimation. Specifically, PPML estimates are trade-weighted averages of the underlying micro-level parameters so that they still provide economically meaningful information. By contrast, no such result exists for OLS since OLS estimation combines aggregation bias with bias resulting from log-linearization, rendering the corresponding parameter estimates uninformative about the parameters of interest.

However, when the regressors vary across sectors (Case 3), even the PPML estimates no longer provide much useful information about the underlying parameters. We showed that the estimates obtained from aggregate data are generally different from the true parameters, and it is impossible to establish the sign or magnitude of the corresponding bias. Finally, when both regressors and parameters vary across sectors (Case 4), it is again impossible to recover useful information about the micro-level parameters of interest, as PPML only estimates a weighted average of approximations to the true micro-level parameters.

Having established these theoretical results, we then argued that they have straightforward implications for the use of gravity equations in applied policy analysis. We again used the effect of trade agreements as an example. Consistent with our theoretical results, we showed that when the parameters and regressors do not vary at the micro level, pre-

dictions based on PPML estimates are robust to aggregation in the sense that we predict the same effect on trade flows irrespective of whether we use aggregate (country-level) or disaggregated (product-level) data. This aggregation property does not carry over to OLS estimation, nor to the case where there is heterogeneity in the micro-level parameters or regressors.

We concluded our analysis by drawing lessons for applied researchers who are interested in the estimation of gravity equations but who might not have disaggregated micro-level data at their disposal. Santos Silva and Tenreyro (2006) have shown that PPML estimation is superior even when disaggregated data are available. In our paper we showed why in situations where only aggregate data are available, researchers are likely to obtain more informative estimates when using PPML as opposed to OLS. We therefore see our results as further strengthening the case for estimating gravity equations in their multiplicative form using PPML, irrespective of whether researchers have access to aggregate or disaggregated data.

References

- Amiti, M., Redding, S., and Weinstein, D. (2019). “The Impact of the 2018 Trade War on U.S. Prices and Welfare,” *Journal of Economic Perspectives* 33, 187-210.
- Amrhein, C.G, and Flowerdew, R. (1992). “The Effect of Data Aggregation on a Poisson Regression Model of Canadian Migration,” *Environment and Planning A: Economy and Space* 24, 1381-1391.
- Anderson, J.E. (2011). “The Gravity Model,” *Annual Review of Economics* 3, 133-160.
- Baier, S.L., and Bergstrand, J.H. (2007). “Do Free Trade Agreements Actually Increase Members’ International Trade?,” *Journal of International Economics* 71, 72-95.
- Bas, M., Mayer, T., and Thoenig, M. (2017). “From Micro to Macro: Demand, Supply, and Heterogeneity in the Trade Elasticity,” *Journal of International Economics* 108, 1-19.

- Cameron, A.C., and Trivedi, P.K. (2013). *Regression Analysis of Count Data*, 2nd edition, New York (NY): Cambridge University Press.
- Carroll, R.J., Ruppert, D., Stefanski, L.A., and Crainiceanu, C. (2006). *Measurement Error in Nonlinear Models: A Modern Perspective*, 2nd Edition, Boca Raton (FL): CRC.
- Drake, C., and McQuarrie, A. (1995). "A Note on the Bias due to Omitted Confounders," *Biometrika* 82, 633-638.
- Feenstra, R., Luck, P., Obstfeld, M., and Russ, K. (2018). "In Search of the Armington Elasticity," *Review of Economics and Statistics* 100, 135-150.
- French, S. (2017). "Comparative Advantage and Biased Gravity," UNSW Business School Research Paper No. 2017-03.
- Gourieroux, C., Monfort, A., and Trognon, A. (1984). "Pseudo Maximum Likelihood Methods: Applications to Poisson Models," *Econometrica* 52, 701-720.
- Head, K., and Mayer, T. (2014). "Gravity Equations: Workhorse, Toolkit, and Cookbook," in G. Gopinath, E. Helpman, and K. Rogoff, eds., *Handbook of International Economics*, Volume 4, Chapter 3, Amsterdam: Elsevier, 131-195.
- Helpman, E., Melitz, M., and Rubinstein, Y. (2008). "Estimating Trade Flows: Trading Partners and Trading Volumes," *Quarterly Journal of Economics* 123, 441-487.
- Imbs, J., and Mejean, I. (2015). "Elasticity Optimism," *American Economic Journal: Macroeconomics* 7, 43-83.
- Kiefer, N.M., and Skoog, G.R. (1984). "Local Asymptotic Specification Error Analysis," *Econometrica* 52, 873-885.
- Kukush, A., Schneeweis, H., and Wolf, R. (2004). "Three Estimators for the Poisson Regression Model with Measurement Errors," *Statistical Papers* 45, 351-368.
- Lewbel, A. (1992). "Aggregation with Log-Linear Models," *Review Economic Studies* 59, 635-642.

- Nakamura, T. (1990). “Corrected Score Function for Errors-in-Variables Models: Methodology and Application to Generalized Linear Models,” *Biometrika* 77, 127-137.
- Neuhaus, J.M., and Jewell, N.P. (1993). “A Geometric Approach to Assess Bias Due to Omitted Covariates in Generalized Linear Models,” *Biometrika* 80, 807-815.
- Pesaran, M.H., and Smith, R. (1995). “Estimating Long-Run Relationships from Dynamic Heterogeneous Panels,” *Journal of Econometrics* 68, 79-113.
- Pritchett, L., and Sethi, G. (1994). “Tariff Rates, Tariff Revenue, and Tariff Reform: Some New Facts,” *World Bank Economic Review* 8, 1-16.
- Redding, S., and Weinstein, D. (2019a). “Aggregation and the Gravity Equation,” *AEA Papers and Proceedings* 109, 450-455.
- Redding, S., and Weinstein, D. (2019b). “Aggregation and the Gravity Equation,” NBER Working Paper 25464.
- Santos Silva, J.M.C., and Tenreyro, S. (2006). “The Log of Gravity,” *Review of Economics and Statistics* 88, 641-658.
- Santos Silva, J.M.C., and Tenreyro, S. (2015). “Trading Partners and Trading Volumes: Implementing the Helpman-Melitz-Rubinstein Model Empirically,” *Oxford Bulletin of Economics and Statistics* 77, 93-105.
- van Garderen, K.J., Lee, K., and Pesaran, M.H. (2000). “Cross-Sectional Aggregation of Non-Linear Models,” *Journal of Econometrics* 95, 285-331.
- Yotov, Y.V., Piermartini, R., Monteiro, J.-A., and Larch, M. (2016). *An Advanced Guide to Trade Policy Analysis: The Structural Gravity Model*, Geneva: World Trade Organization.