

NBER WORKING PAPER SERIES

MARKET SIZE AND TRADE IN MEDICAL SERVICES

Jonathan I. Dingel
Joshua D. Gottlieb
Maya Lozinski
Pauline Mourot

Working Paper 31030
<http://www.nber.org/papers/w31030>

NATIONAL BUREAU OF ECONOMIC RESEARCH

1050 Massachusetts Avenue
Cambridge, MA 02138
March 2023, Revised April 2024

We thank Scott Blatte, Steve Buschbach, Cedric Elkouh, Aaron Haefner, Feng Lin, Luke Motley, Nick Niers, Noah Sobel-Lewin, Vaidehi Parameswaran, and Xiaoyang Zhang for excellent research assistance. We are grateful to Rodrigo Adão, Zarek Brot-Goldberg, Don Davis, Rebecca Diamond, Liran Einav, Ed Glaeser, Tali Han, Tom Holmes, Loukas Karabarbounis, Wojciech Kopczuk, Doug Miller, Sean Nicholson, Bruce Schakman, Martin Schneider, Bradley Setzler, Jon Skinner, Felix Tintelnot, discussants Jan David Bakker, Barthelemy Bonadio, Jessie Handbury, Naomi Hausman, Tim Layton, Elena Patel, and many seminar participants for helpful feedback. We thank the Becker Friedman Institute at the University of Chicago for funding support and enabling us to access the Medicare claims data. We thank Antoine Levy and Jacob Moscona for sharing CBSA-level bedrock-depth data and Amitabh Chandra, Maurice Dalton, and Doug Staiger for hospital quality data and code. Dingel thanks the Cohen and Keenoy Faculty Research Fund at the University of Chicago Booth School of Business for support. Lozinski thanks the NIH (T32GM007281 and T32AG051146). The views expressed herein are those of the authors and do not necessarily reflect the views of the National Bureau of Economic Research.

NBER working papers are circulated for discussion and comment purposes. They have not been peer-reviewed or been subject to the review by the NBER Board of Directors that accompanies official NBER publications.

© 2023 by Jonathan I. Dingel, Joshua D. Gottlieb, Maya Lozinski, and Pauline Mourot. All rights reserved. Short sections of text, not to exceed two paragraphs, may be quoted without explicit permission provided that full credit, including © notice, is given to the source.

Market Size and Trade in Medical Services

Jonathan I. Dingel, Joshua D. Gottlieb, Maya Lozinski, and Pauline Mourot

NBER Working Paper No. 31030

March 2023, Revised April 2024

JEL No. F12,F14,I11,R12

ABSTRACT

We uncover substantial interregional trade in medical services and investigate whether regional increasing returns explain it. In Medicare data, one-fifth of production involves a doctor treating a patient from another region. Larger regions produce greater quantity, quality, and variety of medical services, which they “export” to patients from elsewhere, especially smaller regions. We show that these patterns reflect scale economies: greater demand enables larger regions to improve quality, so they attract patients from elsewhere. Despite concerns about rural access, larger regions have higher marginal returns to spending. We study counterfactual policies that would lower travel costs rather than relocating production.

Jonathan I. Dingel
Booth School of Business
University of Chicago
5807 South Woodlawn Avenue
Chicago, IL 60637
and NBER
jdingel@chicagobooth.edu

Joshua D. Gottlieb
University of Chicago
Harris School of Public Policy
1307 E. 60th St.
Chicago, IL 60637
and NBER
jgottlieb@uchicago.edu

Maya Lozinski
University of Chicago
1307 E 60th St
Chicago, IL 60637
maya.lozinski@gmail.com

Pauline Mourot
University of Chicago
5807 S Woodlawn Ave
Chicago, IL 60637
pmourot@chicagobooth.edu

Rural Americans have worse health outcomes (Finkelstein, Gentzkow, and Williams, 2021), but America’s doctors are disproportionately located in big cities (Rosenblatt and Hart, 2000). This contrast might suggest a spatial mismatch between consumers and producers of medical services, and discussions of whether physicians are geographically “maldistributed” go back decades (Newhouse et al., 1982*a*; Skinner et al., 2019). We consider an alternative explanation: rather than mismatch, do these patterns reflect the benefits of specialization and trade in an industry featuring scale economies?

We find substantial interregional trade in medical services: one-fifth of Medicare production involves a doctor treating a patient from another region. Applying a trade model to these patient flows, we find that scale economies explain larger regions’ production and export of higher-quality medical care. After quantifying regional economies of scale and patients’ willingness to travel, we use the estimated model to explore counterfactual policies aimed at improving access for patients in smaller regions.

If medical services have increasing returns to scale, there are benefits to geographically concentrating production. Indeed, medicine has long been suggested as an industry in which the division of labor is limited by the extent of the market (Arrow, 1963; Baumgardner, 1988). Absent trade, however, the only way to serve patients in smaller regions would be to disperse production across space, foregoing the benefits of scale.¹ This is likely necessary for emergency care. But the vast majority of medical spending is not for emergencies. For example, if patients with cancer can travel across regions in search of the ideal oncologist—one specialized in their particular type of cancer, one with a better reputation, or simply a better personal match—the economic geography of medical care may resemble other tradable industries. In that case, society would face a *proximity-concentration tradeoff*: patients who

¹Many economists assume trade costs for medical services are prohibitively high. Hsieh and Rossi-Hansberg (2021): “Producing many cups of coffee, retail services, or health services in the same location is of no value, since it is impractical to bring them to their final consumers.” Jensen and Kletzer (2005): “Outside of education and healthcare occupations, the typical ‘white-collar’ occupation involves a potentially tradable activity.” Bartik and Erickcek (2007): “An industry can bring in new dollars by selling its goods or services to persons or businesses from outside the local economy (‘export-base production’). . . For health care institutions, demand for services tends to be more local.”

consume medical services produced far from home incur trade costs yet benefit from the higher quality generated by geographically concentrating production.

We investigate these economic forces using Medicare claims data that we introduce in Section 1. Medicare is the federal government’s insurance program for the elderly and disabled and spends four percent of US GDP. Medical service providers submit claims that report the specific medical procedure provided, the treatment location, and where the patient lives. Using hundreds of millions of these claims, we find that “imported” medical procedures—defined as a patient receiving a service provided by a physician in a different region—constitute about one-fifth of US healthcare consumption. Production is geographically concentrated in larger regions, while consumption is much less so. This contrast means exported medical services are disproportionately produced in larger regions and imports are a larger share of consumption in smaller regions. Although some patients travel thousands of kilometers for care, patient flows decline rapidly with distance between regions. While these patterns establish an important role for regional specialization and trade, they do not reveal the strength of regional increasing returns or the consequences of healthcare agglomeration.

To guide our investigation of these mechanisms, Section 2 develops a model of trade in medical services. We adapt standard models to a setting in which the government sets prices, so quality variation clears markets. The model delivers equations that we use to estimate regional service quality, the strength of increasing returns, and market-size effects. If there are regional scale economies, larger markets produce higher-quality care and export it. When economies of scale are sufficiently strong, the model predicts a *strong home-market effect*: greater demand makes larger regions *net* exporters of medical services. Because market size matters more at smaller scales, the model predicts less common medical procedures will respond more to differences in market size. The gravity model of trade flows provides a framework to test this hypothesis, as in Costinot et al. (2019).

In Section 3, we recover revealed-preference estimates of regional service quality by estimating patients’ willingness to travel to each exporting region for medical services. These

estimates based on trade flows align with conventional measures of quality: they are negatively related to external measures of mortality rates and positively related to other hospital rankings and safety scores.

We use these revealed-preference estimates to estimate the scale elasticity: how does serving more patients raise a region's quality of medical care? This parameter plays a key role in examining health policy when there are regional increasing returns. We find substantial regional increasing returns. We estimate the scale elasticity to be about 0.8: a region producing 10% more because of greater demand delivers about 8% higher quality.

The same trade flows show that regional increasing returns are large enough to explain the geographic patterns of production, consumption, and trade documented in Section 1. We find a strong home-market effect in medical services: greater demand induces a greater increase in exports than imports, making larger markets net exporters of medical care. These scale effects cannot be attributed to larger markets having lower input costs or medical production raising population size. This result remains when estimated in the cross-section, a panel, or instrumented with historical characteristics that predict current production levels.

Section 4 investigates potential sources of regional increasing returns in medical services. We examine how market-size effects and returns to scale vary with procedure characteristics. Rare procedures are traded more and over longer distances. For example, half of the patients having left ventricular assist devices (LVADs) implanted to restore their heart function come from outside the surgeon's region, while only 10% of screening colonoscopies are traded. Consistent with the model, the home-market effect is substantially stronger for less common procedures: a larger residential population drives a greater increase in net exports for rare procedures. Moreover, we estimate that rarer procedures exhibit stronger increasing returns.

A variety of mechanisms could generate these increasing returns to scale: finer specialization among physicians, sharing of lumpy capital equipment, knowledge diffusion, or learning by doing (Marshall, 1890). We show that larger markets have more specialized capital equipment and physicians. Critically, trade enables patients across regions to share in these

benefits of scale: imports are more likely to be provided by a specialist—an appropriate and more experienced specialist—and more likely to use rare equipment than locally produced services.

We use our estimates of scale economies and trade costs to explore the proximity-concentration tradeoff quantitatively in Section 5. Counterfactual policies affect regions differently depending on their size and trade patterns. We examine the implications of increasing access to care in one region by either increasing reimbursements or reducing travel costs. Increasing reimbursements has a higher return in more populous regions. When raising reimbursements in the largest regions, the aggregate gain in patient market access per dollar of spending is about 10% higher than when increasing reimbursements in the smallest regions. Increasing reimbursements in one region reduces output quality in neighboring regions, while improving patients’ market access to the extent they import from the treated region. Reducing travel costs for one region increases its import demand, which improves both output quality and market access in neighboring regions. The rich pattern of consequences when subsidizing patients in low-output regions highlights the importance of trade and agglomeration for understanding the incidence of these policies on patients and producers.

The higher-quality care available in larger markets may not benefit all patients equally. We find that patients with lower socioeconomic status (SES) are less likely to travel farther for better care, even when we examine travel patterns within the same billing code. Thus, the proximity-concentration tradeoff varies by SES; the pattern of aggregate returns masks important heterogeneity. We also find that the large geographic scope of the United States explains a meaningful part of the correlation between regional income and health.

This paper builds on research in urban, trade, and health economics. Urban economists have documented skill-biased agglomeration in production as knowledge workers have become more numerous and concentrated in skilled cities (Berry and Glaeser, 2005; Moretti, 2011; Diamond, 2016; Davis and Dingel, 2020; Eckert, Ganapati, and Walsh, 2020). Connecting

this to the production and trade of services has been more difficult. Most studies of the geography of services analyze restaurants and retailers (Davis et al., 2019; Agarwal, Jensen, and Monte, 2020; Allen et al., 2021; Miyauchi, Nakajima, and Redding, 2021; Burstein, Lein, and Vogel, 2022). We show that—even in a service-based economy—the sizes of both local and potential export markets influence production and quality. This suggests that healthcare can serve as an export base for large markets (Bartik and Erickcek, 2007).

The trade literature has examined market-size effects in manufacturing but investigated services much less. Davis and Weinstein (2003), Hanson and Xiang (2004), and Bartelme et al. (2019) link manufactures’ market size to export patterns, in line with the home-market effect of Krugman (1980) and Helpman and Krugman (1985). Dingel (2017) shows that market-size effects drive quality specialization across US cities. Market-size effects for pharmaceuticals have been estimated using demographic variation over time (Acemoglu and Linn, 2004) and across countries (Costinot et al., 2019). Services are much less studied, in part because of the paucity of reliable trade data (Lipsey, 2009; Muñoz, 2022). We advance this literature using the detailed procedure and location information in medical claims data.

The importance of medical care for health and welfare generates substantial public-policy interest. Rural locations have worse health outcomes but fewer doctors per capita. Newhouse et al. (1982*a,b,c*), Newhouse (1990), and Rosenthal, Zaslavsky, and Newhouse (2005) considered this issue and argued against targeting a uniform geographic distribution of physicians. Building on this, we measure interregional trade in medical services, estimate the impact of geography on patient access, and connect this trade to economies of scale. Modern trade theory guides our modeling, estimation strategy, and counterfactual policy analysis.

1 Empirical setting and geographic patterns

This section describes the data, documents how production and consumption of medical services vary with market size, and shows that trade between regions declines with distance.

1.1 Empirical setting

Our primary dataset is claims data from Medicare, the largest source of medical spending in the United States.² Medicare is the federal government’s insurance program for the elderly and disabled. It does not directly employ physicians or run its own hospitals. Instead, it pays bills submitted by independent physicians, physician groups, hospitals, and other medical service providers. These bills—called “claims” in industry terminology—report the specific services provided using 5-digit codes from the Healthcare Common Procedure Coding System (HCPCS). There are over 12,000 distinct HCPCS codes, which identify individual procedures at a granular level.³ In alternative analyses, we use groupings of patient *diagnoses* to account for potential substitution between treatments.

Federal regulation, not pricing decisions by physicians or hospitals, determines the payment for each claim. For physician payments, Medicare sets procedure-specific “reimbursement rates” largely independent of quality, quantity, or region.⁴ We remove regional price differences from our spending measure as described in Appendix A.2. This produces an expenditure measure purged of spatial variation in reimbursement rates. Patients pay a share of these reimbursements through copayments and deductibles, but these cost-sharing rules are constant across regions, and most Medicare patients have a Medigap supplemental or Medicaid insurance that covers most or all of this cost sharing (Cabral and Mahoney, 2019).

Medicare claims report the ZIP codes of both the place of service and the patient’s residence. Using these, we construct a trade matrix for medical services. We study all medical

²Appendix A.1 contains more details about our data and processing.

³Codes distinguish between, *e.g.*, flu vaccines that protect against three or four strains of flu, whether administration is intramuscular or intranasal, and patient ages. There are distinct codes for chest X-rays based on whether the images are of ribs, the breastbone, or the full chest, both sides or one side of the body, and the number of images taken (1, 2, 3, or 4+). By contrast, all of health services are aggregated into one of the nine service sectors available in Canadian data on interprovincial trade (Anderson, Milot, and Yotov, 2014) and the seventeen service sectors in international trade data (Borchert et al., 2021).

⁴While Medicare does have some quality incentive programs, the money at stake is a small share of Medicare’s overall spending (Gupta, 2021). Medicare has some spatial variation in physician reimbursements, but it is not very large, and has diminished over time; details on payment rules are in Clemens and Gottlieb (2014). Hospital payments also reflect local wage indices, the hospital’s occupational composition, and other factors. Gottlieb et al. (2010) show that price variation does not explain spatial variation in expenditure.

care provided by physicians, including both professional and facility charges, outside an emergency room or skilled nursing facility.⁵ Because Medicare rarely reimbursed telehealth before 2020, this trade involves traveling to receive care delivered in-person.⁶ We aggregate the ZIP-code-level information up to 306 hospital referral regions (HRRs), which are geographic units defined by the *Dartmouth Atlas of Health Care* to represent regional markets for tertiary medical care. HRRs are constructed by aggregating residential areas based on where patients were referred for major cardiovascular surgical procedures and for neurosurgery. Each HRR has at least one city where both categories of surgery were performed. Thus, the construction of these geographic units should tend to minimize trade between HRRs.⁷ We construct HRR-to-HRR trade flows by interpreting the patient’s residential HRR as the importing region and the service location’s HRR as the exporting region.⁸

Physicians, hospitals, pharmacies, and other healthcare providers submit different types of claims. We use Traditional Medicare claims data, primarily from 2017, from hospitals, physicians, and outpatient care providers. One year of data includes around 360 million services, representing \$230 billion in spending.⁹ These claims are not perfectly representative of all US healthcare, since Medicare beneficiaries are elderly or disabled.¹⁰ But the geographic distribution of Medicare beneficiaries is quite similar to the overall population, and Medicare alone finances one-fifth of medical spending. These data likely capture the key features

⁵Our results are similar when using only physicians’ professional fees (Appendix Figure D.2 and Appendix Tables D.8 and D.16).

⁶In 2012, Medicare spent only \$5 million—less than 0.001% of its expenditures—on telehealth services (Neufeld and Doarn, 2015), lagging other insurers (Dorsey and Topol, 2016).

⁷We have also used alternative geographies, including core-based statistical areas (CBSAs), metropolitan statistical areas (a subset of CBSAs that excludes the smaller micropolitan areas), and commuting zones (CZs). Because these yield consistent findings, we do not report all such estimates.

⁸The Medicare claims are US patients receiving care at US service facilities. These data do not report any international transactions. Throughout this paper, “imports” and “exports” refer to domestic transactions between regions of the United States.

⁹For one portion of spending (physician services), we observe a random 20% sample of these claims; our spending estimate is scaled up to represent the full Traditional Medicare population.

¹⁰We validate procedure frequencies in the 20% sample of physician services with two alternative data sources: more comprehensive Medicare data (based on all Traditional Medicare patients, but without information on patients’ locations) and Health Care Cost Institute data for privately insured patients (Appendix A.3). We also study the geographic patterns of production by specialty using the physician registry data containing the ZIP code and specialty of all physicians registered to practice in the United States.

of overall healthcare production and consumption. Appendix Table D.1 reports summary statistics of HRR production, consumption, and trade volumes.

1.2 Spatial variation in production and consumption

Figure 1 shows maps of healthcare production and consumption across regions, based on the place of service and patient’s residential address, respectively. The consumption map shows substantial variation that has been well documented by the *Dartmouth Atlas* and related literature on geographic variation in healthcare (Fisher et al., 2003*a,b*; Finkelstein, Gentzkow, and Williams, 2016). The production map shows even more pronounced variation: more production in large urban agglomerations and less in rural areas. There is substantial variation in production even between immediate neighbors, while consumption varies more smoothly.

The subsequent panels show patterns of trade. Nationally, 19% of production is exported to a patient in another HRR, and 21% is traded between CBSAs.¹¹ Panel 1(c) shows the ratio of production to consumption; a value larger than one means an HRR is a net exporter. Net-exporting regions tend to be major urban agglomerations, plus places such as Rochester, Minn. and Durham, N.C. that specialize in healthcare. Panel 1(d) shows gross exports as a share of local production for each HRR. Three-quarters of services produced in the Rochester metropolitan area, home to the top-rated Mayo Clinic, are provided to patients from other regions, who travel an average of 545 km to Rochester.¹² As a major healthcare exporter with a population of merely 220,000, Rochester is an outlier: larger regions are responsible for a disproportionate share of medical services production.

Figure 2 plots the average production and consumption per beneficiary across HRRs of different sizes. Production per beneficiary increases with population, with an elasticity of 0.05. Consumption per beneficiary is virtually uncorrelated with population; in fact, the

¹¹42% of physician production is exported outside the provider’s Hospital Service Area (HSA), a smaller *Dartmouth Atlas* geographic unit. For comparison, 6% of employees work outside their commuting zone of residence (Monte, Redding, and Rossi-Hansberg, 2018, Appendix Table B.1). For manufactured goods, we obtain a 68% export share for metropolitan areas covered by the Commodity Flow Survey.

¹²These top hospitals tend to have top physicians; Mourot (2024) documents positive assortative matching.

estimated elasticity is slightly negative. The difference between production and consumption is net trade: larger markets are net exporters and smaller markets are net importers. Gross trade flows exceed net trade flows, with imports comprising about one-third of consumption in the smallest regions. Imports per beneficiary decline with an elasticity of -0.35 with respect to population. Exports per beneficiary are approximately flat, which means total exports are increasing with local population.

1.3 Bilateral trade and bilateral distance

Despite the clear patterns in Figure 2, geographic variation in trade is far from entirely explained by market size. The regions with the lowest export shares are Anchorage, Honolulu, and Yakima, Wash., likely reflecting their remoteness. To account for these geographic patterns, we examine bilateral trade flows.

Figure 3 depicts how trade varies with the distance between the patient and place of service. Panel (a) shows the distribution of distances patients travel for care, distinguishing between care provided in the patient’s home region and in other regions. Within HRRs, there is a narrow distribution of distances that peaks around 10 km. When visiting providers in a different HRR, patients travel a great variety of distances. There is a local plateau between approximately 30–100 km, suggesting a fair amount of travel to nearby HRRs, perhaps indicating regional medical centers. There is substantial long-distance travel for care: many patients travel hundreds or thousands of kilometers.¹³

Panel 3(b) shows that trade declines with distance. The blue curve depicts trade volume against distance (for pairs of HRRs with positive trade flows) after removing fixed effects for each exporter and each importer.¹⁴ This intensive-margin relationship is roughly log-linear. The red curve shows the extensive margin: the share of pairs with positive trade as a function

¹³Patients’ choices to travel these distance underpin our revealed-preference estimates of regional service quality in Section 3.1. The average patient travels 640 km to Chicago and 740 km to New York City, compared with less than 146 km to Urbana-Champaign, Ill. or Charlottesville, Va. An older literature cited in Dranove and Satterthwaite (2000) finds that patients who travel farther incur higher hospital costs.

¹⁴This application of the Frisch-Waugh-Lovell theorem is only feasible for positive trade volumes.

of distance. This is 100% for nearby pairs and under 80% for the most distant pairs. These patterns motivate the inclusion of distance covariates in our gravity-based analysis.

Patients may vary in their ability or willingness to travel, especially by socioeconomic status. We quantify it here, to the extent feasible in our data, to enrich counterfactual analysis and interpret welfare implications. Panel 3(c) depicts distance elasticities—how fast bilateral trade declines with the distance between regions—estimated separately by neighborhood income decile.¹⁵ We find a strong, nearly monotonic relationship between socioeconomic status and the distance elasticity: patients from the highest neighborhood-income decile exhibit a distance elasticity 25% smaller than those in the lowest decile.¹⁶ This means patients from higher-income neighborhoods are more amenable to traveling farther for medical care. If society faces a proximity-concentration tradeoff, this tradeoff varies by socioeconomic status. This is especially notable given the empirical setting: Medicare insures the near-universe of elderly and disabled Americans.

2 Theoretical framework

This section develops a model of regional trade in medical services tailored to our analysis of US healthcare. We develop a competitive model of a market for one medical procedure that has a fixed price.¹⁷ With fixed prices, quality variation plays a key role in clearing markets. Patients select quality-differentiated services and face trade costs, so patient flows between regions follow a gravity equation. Regional increasing returns cause the quality-adjusted cost of producing a service to decline with scale. The model delivers equations that allow us to estimate the strength of regional increasing returns and the home-market effect. We use

¹⁵We estimate these distance elasticities using equation (8) as described in Section 2.6.

¹⁶These estimates are consistent with the interaction that Silver and Zhang (2022) estimate between income and distance to care. These differences in distance elasticities are not driven by differences in the composition of procedures. When we estimate elasticities separately for rare and common services—or even for individual procedures (see Appendix Table D.5)—the income gradient of distance elasticities persists.

¹⁷For brevity, we present a competitive model, but the consequences of regional increasing returns for trade flows in a fixed-price environment do not hinge on this assumption. Appendix B.1 shows that a monopolistic-competition model with one medical provider in each region delivers the same predictions. As in flexible-price models, many market structures can give rise to a home-market effect (Costinot et al., 2019).

the estimated model to quantify outcomes in counterfactual policy scenarios.

2.1 Demand

We use a logit model of individuals choosing providers for a given service. Providers and patients are in regions indexed by i or j , with \mathcal{I} denoting the set of regions. All providers in a region are identical. Let N_j denote the number of patients residing in region j who make a choice.¹⁸ Patient k in region j choosing a provider in region i obtains utility

$$U_{ik} = \ln \delta_i + \ln \rho_{ij(k)} + \epsilon_{ik}.$$

The provider-region component δ_i would usually include a product's characteristics and price. With fixed prices, δ_i is simply the quality of care available in region i . The region-pair component ρ_{ij} represents bilateral inverse trade costs (proximity). The idiosyncratic component ϵ_{ik} is independently and identically drawn from a standard Gumbel distribution, so the probability that patient k selects a provider in region i is

$$\Pr(U_{ik} > U_{i'k} \ \forall i' \neq i) = \frac{\exp(\ln \delta_i + \ln \rho_{ij(k)})}{\sum_{i' \in 0 \cup \mathcal{I}} \exp(\ln \delta_{i'} + \ln \rho_{i'j(k)})}.$$

There is an outside option denoted by $i = 0$, which represents individuals choosing to forgo care, and we normalize its common component to zero, $\ln \delta_0 = \ln \rho_{0j(k)} = 0 \ \forall k$.¹⁹

This choice probability implies a gravity equation for the quantity of trade between any two regions when we aggregate patients' decisions. Let Q_{ij} denote the quantity of procedures supplied by providers in i to patients residing in j , and let Q_{0j} denote the number of patients in j selecting the outside option. Because each patient selects at most one

¹⁸Appendix B.2 extends the model to have multiple patient types.

¹⁹This formulation of demand is familiar from the hospital competition literature, which has been surveyed by, e.g., Gaynor and Town (2011) and Gaynor, Ho, and Town (2015). Our competitive model does not distinguish between hospitals or physicians within a region.

provider, $N_j = \sum_{i \in \mathcal{I} \cup \{0\}} Q_{ij}$. The demand by patients in j for procedures performed in i is

$$Q_{ij} = \delta_i \frac{N_j}{\Phi_j} \rho_{ij}, \quad (1)$$

where $\Phi_j \equiv \sum_{i' \in 0 \cup \mathcal{I}} \delta_{i'} \rho_{i'j}$ is the expected value of the choice set for patients in region j . We call this Φ_j “patient market access,” as it captures the quality of care available to patients and their costs of accessing it. Equation (1) is a gravity equation with an origin i component, a destination j component, and an ij pair component. Total demand for care produced in i is

$$Q_i = \delta_i \sum_{j \in \mathcal{I}} \frac{N_j}{\Phi_j} \rho_{ij}. \quad (2)$$

2.2 Production

We assume competitive production of services with free entry and regional increasing returns that are external to individual providers. That is, each price-taking provider chooses its output quality and quantity given total regional production, an exogenous factor price, and an exogenous productivity shifter. A provider in region i that employs L units of the composite input to produce service of quality δ produces the following output quantity:

$$A_i \frac{H(Q_i)}{K(\delta)} L.$$

Improving quality is costly so $K(\delta)$ is increasing. Regional increasing returns to scale are a weakly increasing, concave function $H(Q_i)$ of total regional production, Q_i , which competitive producers take as given (Chipman, 1970). The regional productivity shifter A_i captures any other influences, such as past investments. Provider size L is indeterminate (and unimportant) given the linear production function, external economies of scale, and price-taking behavior. The composite input is supplied to region i at factor price w_i .²⁰ Thus, the unit

²⁰If the regional factor supply were less than perfectly elastic, we would estimate increasing returns net of the cost of hiring additional inputs. That is, if the factor supply elasticity were β , our estimate of the scale elasticity α in equation (4) would instead be an estimate of the effective scale elasticity $\tilde{\alpha} \equiv \alpha - \frac{\beta}{1+\beta}$.

cost of producing quality δ in region i is

$$C(Q_i, \delta_i; w_i, A_i) \equiv \frac{w_i K(\delta_i)}{A_i H(Q_i)}.$$

In our setting, output prices are fixed: the Medicare-determined “reimbursement rate” \bar{R} is independent of quality, quantity, and region. Each provider that produces output of the highest quality produced in region i earns revenue \bar{R} per unit.

Provider optimization and free entry require the unit cost to equal the reimbursement rate in each region. Given factor price w_i and productivity shifter A_i , the free-entry condition

$$C(Q_i, \delta_i; w_i, A_i) = \bar{R} \tag{3}$$

defines a regional isocost curve: the set of quantity-quality combinations for which the average cost of production equals the reimbursement rate. Regional increasing returns make the isocost curve upward-sloping in (Q, δ) space. With free entry and fixed prices, the benefits of scale are realized as higher-quality services in higher-output regions.

While our assumptions thus far suffice for qualitative results, we later specify functional forms for additional predictions and empirical quantification; specifically, $K(\delta_i) = \delta_i$ and $H(Q_i) = Q_i^\alpha$, with a scale elasticity of $\alpha \geq 0$.²¹ In this case, the free-entry condition (3) is

$$\bar{R} = \frac{w_i \delta_i}{A_i Q_i^\alpha}. \tag{4}$$

2.3 Equilibrium

Equilibrium equates supply and demand in each region, $Q_i = \sum_j Q_{ij}$. Given exogenous parameters \bar{R} , $\{w_i, A_i, N_i\}_{i \in \mathcal{I}}$, and $\{\rho_{ij}\}_{(i,j) \in (\mathcal{I}, \mathcal{I})}$, an equilibrium is a set of quantities and qualities $\{Q_i, \delta_i\}_{i \in \mathcal{I}}$ that simultaneously satisfy equations (2) and (3).

²¹There are increasing returns to scale if $\alpha > 0$. Note that this formulation is compatible with an equilibrium in which region i does not produce the procedure at all: $Q_i = 0 \implies \delta_i = 0$.

2.4 Scale effects in autarky

In autarky, patients choose whether to receive care, but they cannot travel between regions ($\rho_{ij} = 0$ for $i \notin \{0, j\}$). In this case, all demand is local and equation (2) simplifies to

$$Q_{jj} = \frac{\delta_j \rho_{jj}}{1 + \delta_j \rho_{jj}} N_j. \quad (5)$$

The autarkic equilibrium is at the intersection of the demand curve given by equation (5) and the free-entry isocost curve given by equation (3).²² An increase in population size, $\Delta N_j > 0$, affects equilibrium outcomes by shifting the demand curve.

Figure 4 illustrates how greater demand affects quality in autarky. Panel (a) shows the role of returns to scale. The vertical axis shows quality δ_i and the horizontal axis shows quantity Q_i (on logarithmic scales). Higher quality attracts more patients, so demand is upward-sloping.²³ We show two cases of the free-entry isocost curve defined by equation (4): the horizontal line depicts constant returns ($\alpha = 0$) and the upward-sloping line depicts increasing returns ($\alpha > 0$). With constant returns, a rightward shift in demand ($\Delta N_j > 0$) causes a proportional increase in quantity produced and no change in output quality. With increasing returns, higher demand moves producers up the isocost curve. This quality improvement causes a more-than-proportional increase in quantity produced.

Panel 4(b) shows that an increase in demand raises quality more as the demand curve is increasingly elastic. The panel depicts two demand curves: the one on the left is more elastic, as we would expect for a less common procedure.²⁴ Shifting each demand curve to the right raises the equilibrium quality of each procedure because of increasing returns to scale. This market-size effect is larger for the less common procedure with more elastic

²²For the equilibrium to be Marshallian stable, the demand curve must be steeper than the isocost curve at the intersection. There is a stable equilibrium because equation (5) means $Q_{jj} \rightarrow N_j$ as $\delta_j \rightarrow \infty$.

²³For visual clarity, we draw a log-linear demand curve. The logit demand function (5) is in fact log-convex, which is consistent with all the comparative statics illustrated in Figure 4.

²⁴The demand function (5) is log-convex, so demand is indeed more elastic at lower quality. This is a fixed-price counterpart of Marshall's second law that demand is more elastic at higher prices.

demand because the demand shift is amplified by a larger increase in quantity demanded.²⁵

2.5 Market-size effects on trade flows

We now consider trade. With multiple regions and finite trade costs ($\rho_{ij} > 0$), some patients will import—*i.e.*, select a provider located in another region. This trade stems from two sources. First, in the logit demand system with finite trade costs, patients have idiosyncratic preferences that yield a strictly positive probability of choosing every region. Second, when quality varies, regions producing higher-quality services attract more patients.

Fixing the qualities produced in other regions, an increase in one region’s demand affects its trade flows through three mechanisms. First, greater total demand for services proportionally increases a region’s demand for imports through the N_j term in equation (1). Second, with increasing returns, an increase in N_i elicits an increase in quality δ_i , which raises region i ’s *gross* exports to each region. Costinot et al. (2019) call this the “weak home-market effect.” Third, if increasing returns are sufficiently strong, quality δ_i improves so much that region i ’s patient market access Φ_i rises such that $\ln \delta_i$ rises more than $\ln(N_i/\Phi_i)$ does. That is, the increase in region i ’s gross exports exceeds any increase in its gross imports. This is the “strong home-market effect”: an increase in local demand raises a region’s *net* exports.

Figures 4(c) and 4(d) introduce trade and illustrate the distinction between weak and strong home-market effects.²⁶ Panel (c) depicts the quality and quantity produced in one region under two scale elasticities. Comparing points B and C , a given increase in demand elicits a larger quality improvement when increasing returns are stronger. Panel (d) depicts equilibrium exports and imports as a function of the region’s demand shifter N_j . The import

²⁵Alternatively, one could obtain this prediction by assuming that demand is log-linear and the isocost curve is log-concave, so that rare procedures have stronger scale economies. For example, introducing a fixed cost F would yield the isocost curve $\frac{w_i \delta_i}{A_i Q_i^\alpha} + F Q_i^{-1} = \bar{R}$. In this case, the scale elasticity $\frac{\partial \delta_i}{\partial Q_i} \frac{Q_i}{\delta_i} = \frac{\alpha \bar{R} + (1-\alpha) F Q_i^{-1}}{\bar{R} - F Q_i^{-1}}$ is decreasing in quantity produced Q_i for $\alpha \in (0, 1)$. A rightward shift in demand would cause a larger (log) difference in quality for the low-volume procedure on the steeper part of the isocost curve.

²⁶These diagrams are fixed-price analogues of Figures II and III in Costinot et al. (2019). See their discussion of the assumption that one region is large enough to affect its own quality but too small to affect the quality produced in other regions. This assumption is only made for this figure.

curves slope up because an increase in local demand raises demand for imports. The export curves slope up because of increasing returns: higher local demand increases quality, which increases gross exports. This is the weak home-market effect. When the scale elasticity α is larger—the free-entry isocost curve in Panel (c) is steeper—greater demand elicits a larger increase in output quality. This steepens the export curve and flattens the import curve in Panel (d). When the export curve is steeper than the import curve, there is a strong home-market effect: the increase in demand raises exports more than imports.

As in the autarkic case, we predict larger effects of market size for less common procedures. An increase in demand raises quality more when demand is more elastic, leading to a stronger home-market effect for rare procedures. If rare procedures also have greater economies of scale (higher α) that would amplify this contrast. This result motivates a research design comparing home-market effects for common and rare procedures.

These results also hold when an increase in demand in one region affects equilibrium outcomes in other regions. To demonstrate this, we examine the home-market effect in the neighborhood of a symmetric equilibrium. If all regions are the same size, $N_i = \bar{N} \forall i$, and trade costs are symmetric, $\rho_{ii} = 1$ and $\rho_{ij} = \rho \in (0, 1) \forall i \notin \{0, j\}$, there is a symmetric equilibrium with quality $\bar{\delta}$ and patient market access $\bar{\Phi}$ in each region. As detailed in Appendix B.3, we totally differentiate the system of equations in terms of $\{d\delta_i, dN_i\}_{i=1}^I$ and evaluate a change in one region's demand N_1 at the symmetric equilibrium.

With increasing returns of any magnitude, there is a weak home-market effect; with sufficiently strong increasing returns, there is a strong home-market effect. When $\alpha \in (0, 1)$, an increase in region 1's demand elicits an increase its relative service quality:

$$d \ln \delta_1 - d \ln \delta_{j \neq 1} = \left[\frac{1 - \alpha (\bar{\Phi} - 1)}{\alpha (1 - \rho) \bar{\delta}} + \frac{(1 - \rho) \bar{\delta}}{\bar{\Phi}} \right]^{-1} d \ln N_1 > 0.$$

This higher quality causes region 1 to export more to every other region: $\frac{d \ln Q_{1j}}{d \ln N_1} > 0$. The

effect on the region's net exports is

$$d \ln Q_{1,j \neq 1} - d \ln Q_{j \neq 1,1} = \left[\frac{1 - \frac{1-\alpha}{\alpha} \frac{1+(\mathcal{I}-1)\rho}{1-\rho}}{\frac{1-\alpha}{\alpha} \frac{1+(\mathcal{I}-1)\rho}{(1-\rho)} + \frac{(1-\rho)\bar{\delta}}{1+(1+(\mathcal{I}-1)\rho)\bar{\delta}}} \right] d \ln N_1. \quad (6)$$

Net exports increase with population size if and only if $\frac{\alpha}{1-\alpha} > \frac{1+(\mathcal{I}-1)\rho}{1-\rho}$. This occurs if increasing returns are sufficiently strong (α is large enough) and trade costs are sufficiently large (ρ is small enough). Otherwise, there is a weak home-market effect but not a strong one. Given a strong home-market effect, the effect in equation (6) is diminishing in the population of patients \bar{N} , so we predict a stronger home-market effect for rare procedures.

2.6 Estimating quality, scale economies, and home-market effects

The model yields equations for estimating quality δ_i , the scale elasticity α , and detecting home-market effects. Equation (1) provides a revealed-preference measure of each region's quality, δ_i . We assume the region-pair component is $\ln \rho_{ij} = \gamma X_{ij} + v_{ij}$, where X_{ij} is a vector of observed trade-cost shifters and v_{ij} is an orthogonal unobserved component. Taking expectations and then logs of equation (1) yields gross bilateral trade flows:

$$\ln \mathbb{E}[S_{ij}] = \ln \delta_i + \ln \left(\frac{N_j}{\Phi_j} \right) + \gamma X_{ij}. \quad (7)$$

S_{ij} is the value of procedures exported from region i to patients residing in j .²⁷ The right side contains three terms: exporter quality $\ln \delta_i$, importer demand shifter $\ln \left(\frac{N_j}{\Phi_j} \right)$, and observed trade costs γX_{ij} . These trade costs may reflect the opportunity cost of patients' time, fiscal costs of travel, and any adverse health consequences of travel. Since our empirical work excludes emergency care, direct health costs are not likely to be critical.

To estimate each region's service quality, we replace the first two terms with exporter

²⁷For most of the analysis, we define this as total value across all care, computed as described in Section 1.1. In some cases, we consider subsets of care or even individual procedures. In the procedure-specific case, the dependent variable is the procedure count Q_{ij} rather than spending S_{ij} because no aggregation is required.

fixed effects and importer fixed effects, respectively:

$$\ln \mathbb{E}[S_{ij}] = \underbrace{\ln \delta_i}_{\text{exporter FE}} + \underbrace{\ln \theta_j}_{\text{importer FE}} + \gamma X_{ij}. \quad (8)$$

The exporter fixed effects are a revealed-preference measure of regional quality, which we validate using external quality metrics. The importer fixed effects, combined with an assumption about the number of potential patients (see Appendix C.2), enable us to compute $\Phi_j = N_j/\theta_j$, a measure of patient market access for those who reside in location j .

The relationship between quality and quantity reveals the scale elasticity of the regional production function, α . Per the free-entry condition (4), quality is an isoelastic function of the quantity produced, conditional on price, cost, and productivity shifters. We take the log of (4), replace $\ln \delta_i$ with its estimate $\widehat{\ln \delta_i}$ from (8), and rearrange terms to obtain the estimating equation:²⁸

$$\ln \delta_i = \alpha \ln Q_i + \ln \bar{R} - \ln w_i + \ln A_i. \quad (9)$$

To learn whether these regional increasing returns explain the pattern of trade described in Section 1, we estimate a more parsimonious gravity regression that uses each region's population as a demand shifter. Following Costinot et al. (2019), we differentiate the system of equations (2) and (3) around the symmetric equilibrium and replace $\ln \delta_i$ and $\ln \left(\frac{N_j}{\Phi_j} \right)$ in (7) by log population in the producing and consuming regions, respectively, yielding:

$$\ln \mathbb{E}[S_{ij}] = \lambda_X \ln \text{population}_i + \lambda_M \ln \text{population}_j + \gamma X_{ij}. \quad (10)$$

A positive coefficient $\lambda_X > 0$ implies a weak home-market effect as defined in Costinot et al. (2019): *gross* exports increase with market size. If $\lambda_X > \lambda_M$, the home-market effect is strong: *net* exports increase with market size.

²⁸Appendix A.5 quantifies the potential bias resulting from our observing only the quantity produced for Traditional Medicare beneficiaries, rather than the total quantity produced for all patients. It shows that the bias is small: the estimates in Table 1 should be deflated by about 5%.

To estimate these specifications, we must parameterize observed trade costs γX_{ij} . We first use log distance and a same-region dummy, *i.e.* $\gamma X_{ij} = \gamma_1 \ln \text{distance}_{ij} + \gamma_0 \mathbf{1}(i = j)$. Alternative specifications add $(\ln \text{distance}_{ij})^2$ or replace these continuous distance covariates with indicators for distance deciles. Since observed bilateral trade is zero for many pairs of regions, especially when looking at trade in individual procedures, we use the Poisson pseudo-maximum-likelihood (PPML) estimator (Santos Silva and Tenreyro, 2006).

In other analyses, we estimate slight variants of these specifications. First, we leverage population growth over time using a panel regression with a region-pair fixed effect:

$$\ln \mathbb{E}[S_{ijt}] = \lambda_X \ln \text{population}_{it} + \lambda_M \ln \text{population}_{jt} + \phi_{ij} + \gamma_t X_{ij}, \quad (11)$$

where S_{ijt} denotes the gross bilateral trade flow in year t . In this specification, changes in population across time periods t identify λ_X and λ_M . Second, we contrast common and rare services by dividing procedures into two groups based on whether their national volume is above or below that of the median procedure.²⁹ We estimate equation (8) separately for these two groups to obtain the qualities of rare and common procedures. Let S_{ijc} denote exports from i to j in category of care $c \in \{\text{common}, \text{rare}\}$. The model predicts a stronger home-market effect for rare procedures, so we estimate the following specification:

$$\begin{aligned} \ln \mathbb{E}[S_{ijc}] &= \lambda_X \ln \text{population}_i + \lambda_M \ln \text{population}_j + \gamma X_{ij} \\ &+ (\mu_X \ln \text{population}_i + \mu_M \ln \text{population}_j + \psi X_{ij}) \cdot \mathbf{1}(c = \text{rare}). \end{aligned} \quad (12)$$

In the presence of an overall strong home-market effect, theory predicts $\mu_X > \mu_M$. An alternate specification introduces ij -pair fixed effects, which absorb all the covariates not interacted with $\mathbf{1}(c = \text{rare})$. We also estimate similar models with heterogeneity in other dimensions, such as procedure intensity and frequency of patient engagement, and for specific procedures.

²⁹Appendix Figure D.1(b) plots the distribution of import shares across regions for the two groups.

3 Regional increasing returns in medical services

This section tests for increasing returns in healthcare and the implications for trade flows. Section 3.1 estimates region-level quality of medical care and shows that external quality indicators line up with revealed-preference measures based on trade flows. Using these measures, Section 3.2 finds substantial regional increasing returns in healthcare production. Section 3.3 demonstrates a strong home-market effect, implying that scale economies indeed generate the observed patterns of production and trade.

3.1 Quality estimates

Our revealed-preference measures of regional service quality are the exporter fixed effects $\widehat{\ln \delta_i}$ from estimating equation (8). Figure 5 relates these fixed effects to external measures of regional hospital quality. We aggregate estimated hospital mortality from Chandra, Dalton, and Staiger (2023) and the Centers for Medicare and Medicaid Services (CMS) by HRR. We also count the number of times each region’s hospitals appear on *U.S. News Best Hospitals*.³⁰ Panels 5(a)–(c) establish that patients travel farther to obtain care from regions with lower hospital mortality rates or better *U.S. News* rankings. Appendix Table D.2 reports the corresponding regressions. It also shows a regression relating the exporter fixed effects simultaneously to a variety of other hospital safety, mortality, accreditation, and other quality measures. This regression has an adjusted R^2 of 0.60. While our revealed-preference estimates capture all attributes that influence patient choices, their relationship to these other quality measures suggest that our estimates capture many characteristics of clinical interest.

³⁰Compared with the CMS measures, Chandra, Dalton, and Staiger (2023) use empirical Bayes estimation to account for differences in hospital volume and quality drift over time. Appendix A.1 explains how we use the *U.S. News* rankings. Appendix Figure D.3 also presents results using Hospital Safety Grades from the Leapfrog Group.

3.2 Estimating the scale elasticity

3.2.1 Empirical implementation

We use equation (9) to estimate the strength of regional increasing returns. Regional output quality depends on the volume of production and exogenous shifters of the isocost curve. One potential concern with estimating equation (9) by ordinary least squares is reverse causality. Shifts of the isocost curve would cause movements along the upward-sloping demand curve, biasing the estimated scale elasticity upwards. We address this by using three instruments for production and differences over time.

Our first instrument is current population. Population is relevant for healthcare output, since larger populations clearly require more healthcare.³¹ The instrument is valid if it shifts only demand and not supply. We show that urban amenities do not make healthcare labor inputs cheaper in larger markets.

Nevertheless, one potential concern with this instrument is reverse causality. Suppose that success in exporting medical services serves as an employment base that raises current population size, as epitomized by “anchor institutions.”

We use two further instruments to address this concern. Our second instrument is historical population. In 1940, medicine was a far smaller industry and did not drive local population in the way it might today. Because local populations persist, population in 1940 predicts current population, so we use the former as an instrument for the latter.

The next instrument goes farther back than 1940 and uses local geology to predict healthcare production. Rosenthal and Strange (2008) and Levy and Moscona (2020) show that shallower subterranean bedrock makes construction easier, increasing population density. Bedrock depth also predicts population size, so it is our third instrument for local demand.³²

³¹The variation in demand predicted by population is overwhelmingly driven by headcount, not income per capita: the population elasticity of income per capita is below 0.1, and the income elasticity of health spending is likely below 1.0 (Acemoglu, Finkelstein, and Notowidigdo, 2013).

³²This instrument is currently only available for CBSAs and CZs, not HRRs. Our main results hold when defining regions as CBSAs and CZs. Levy and Moscona (2020) show that the bedrock instrument has ample first-stage power for predicting CBSA population density; the same is true for our endogenous variables (population levels).

Our final research design uses variation over time. We estimate equation (9) in first differences using changes from 2013 to 2017. This uses the relationship between quality growth and production growth across regions to estimate the scale elasticity. To similarly confound both the cross-sectional and first-difference regressions, omitted supply shifters would have to exhibit similar cross-sectional and temporal variation.

3.2.2 Scale improves quality

Estimated service quality $\widehat{\ln \delta}_i$ rises substantially with the regional production volume $\ln Q_i$. Figure 5(d) depicts this relationship, which appears consistent with our isoelastic specification. Table 1 reports regression estimates. The estimated scale elasticity is around 0.8 and stable under various estimation approaches. The first row shows our baseline estimate for the cross-section, and the second uses variation over time. The third and fourth rows instrument for output using contemporaneous or historical population, which exhibit strong first-stage regressions. The second column omits the diagonal S_{ii} observations when estimating the gravity equation (8), to avoid any bias from having a region’s own local consumption influence both the quality measures and output. The third column controls for (small) spatial variation in reimbursements, mean two-bedroom property value, and mean annual earnings for non-healthcare workers. Instrumenting for output tends to reduce the estimated scale elasticity. Excluding the diagonal of the trade matrix when estimating quality tends to raise it. Across all 12 estimates, the lowest elasticity is 0.63 and the highest is 1.04.³³ The results for CBSAs and CZs, reported in Appendix Tables D.3 and D.4, are also stable across specifications and when using the bedrock instrument.

The rest of this paper demonstrates the economic implications of these scale economies in three ways. First, Section 3.3 shows that they are large enough to generate the observed patterns of production and trade in medical services. Second, Section 4 demonstrates mech-

³³These estimates lie in the middle of other estimated agglomeration elasticities. Kline and Moretti (2013) estimate an elasticity of 0.4–0.47 from the Tennessee Valley Authority’s investments. In manufacturing, Greenstone, Hornbeck, and Moretti (2010) report an analogous elasticity above 1 (a 12% increase in total factor productivity caused by adding a plant representing 8.6% of the county’s prior output).

anisms that contribute to regional increasing returns. Third, Section 5 uses the parameter estimates to analyze counterfactual scenarios that illustrate how these forces work and provide guidance for healthcare policy in a world with increasing returns.

3.3 A strong home-market effect in medical services

While we have established there are regional increasing returns, we have yet to see if they are sufficiently strong to explain the fact that larger markets are net exporters. This question determines whether market size enables larger regions to specialize in healthcare production and to share the benefits of scale with patients in other regions. We now test for a strong home-market effect in medical services by estimating equation (10).

Table 2 presents the results. The first column shows significant, positive coefficients on both provider- and patient-market population.³⁴ The coefficient on provider-market population is substantially greater than that on patient-market population. This demonstrates a *strong* home-market effect (Costinot et al., 2019). Not only does a larger population increase gross exports, but it does so more than it increases gross imports by local patients.

The distance elasticity of medical services trade between hospital referral regions is -1.6. This is substantially larger than the distance elasticity of -0.95 estimated for trade in manufactures between CBSAs (Dingel, 2017).³⁵ This suggests that trade in personal services incurs greater distance-related costs, relative to the degree of product differentiation across regions, than trade in manufactured goods. The most obvious difference is that patients themselves must travel to the provider.

The next two columns of Table 2 demonstrate that more flexible distance-covariate specifications do not meaningfully alter the estimated home-market effect. Column 2 introduces the square of log distance as an additional covariate. Column 3 replaces the parametric distance controls with dummies for deciles of distance. The result is stable across the

³⁴In all estimates using trade matrices, we two-way cluster standard errors by patient and provider market.

³⁵We find a distance elasticity of medical services trade between CBSAs of -2.4. The analogous elasticity of health care and social assistance services trade between Canadian provinces is -1.42 (Anderson, Milot, and Yotov, 2014). The distance elasticity of international trade is typically near -0.9 (Disdier and Head, 2008).

columns: both gross and net exports increase with market size.

The fourth column of Table 2 uses the historical population instrument to address concerns about reverse causality. We obtain similar home-market-effect estimates to our baseline results. Appendix Tables D.6 and D.7 report similar results when using CBSAs or CZs rather than HRRs as our geographic unit. Appendix Tables D.6 and D.7 also show the results are robust to instrumenting with historical population and with bedrock depth.

The final column of Table 2 uses changes in population over time to proxy for changes in demand, per equation (11). The market-pair fixed effects in this specification absorb all cross-sectional variation. We find a strong home-market effect that is similar to the cross-sectional estimate.

Competing explanations. Our instruments address reverse causality, but all of them operate through population size. So other channels could generate relationships similar to the market-size effect we estimate. Most significantly, if doctors prefer to live in big cities (Lee, 2010), as college graduates generally do (Diamond, 2016), they could accept lower nominal wages and thus reduce healthcare production costs w_i in such cities. This would raise quality in large markets, but through a different mechanism than increasing returns.

Before we address this problem, first note what is *not* a problem: physicians preferring to work in larger regions for job-related reasons. If a region’s scale enables it to support large academic medical centers, which attract workers, this is an agglomeration benefit that the scale elasticity α ought to capture.³⁶ More patients enable physicians to specialize, conduct research, and train medical students. These forces operate through the scale of healthcare production in the region, and academic medical centers are part of this production function.

The challenge to our interpretation arises if physicians prefer larger markets for non-

³⁶The most salient example is Cornell University: after an abortive attempt to have medical training in both Ithaca and New York City, the Cornell Trustees quickly closed down the Ithaca location and centered the medical school in New York—where the patients and doctors were more abundant—in the early 20th century (Flexner, 1910; Gotto and Moon, 2016). As this history illustrates, the potential local demand for care can drive the location of medical training. In general education, in contrast, university placement induces economic growth (Moretti, 2004).

medical reasons, and this labor supply shift increases quality or reduces costs in larger markets. We investigate whether this mechanism is strong enough to reduce net costs in larger markets. Doctors are cheaper in larger markets (Gottlieb et al., 2023), but other costs rise with population size. Appendix Figure D.4 shows that the population elasticity of doctors' earnings is -0.01, but that for non-physicians is 0.043.³⁷ To compute the population elasticity of overall labor costs, note that non-physician labor's share of healthcare production is twice the physician labor's share (Appendix Table D.9). The population elasticity of labor costs is thus positive. The higher cost of real estate in larger markets reinforces these higher labor costs. This spatial variation undermines the idea that amenities make production cheaper in larger markets. Section 4 further shows that specialty-specific income elasticities are close to zero on average and unrelated to the number of specialists.

One final concern is measurement error in Medicare's records of patients' residences. To address this, Appendix A.4 first demonstrates our results are robust to excluding states with large seasonal populations. Second, we examine how far dialysis patients appear to travel. We find that residential measurement error is limited and does not drive our results.

4 Mechanisms

This section explores potential sources of scale economies and asks if trade expands the geographic scope of these mechanisms. We focus on professional fees for which payments are made at the service level.³⁸ Section 4.1 examines how market-size effects vary with procedure characteristics, especially the rare-vs-common comparisons suggested by the model. Consistent with our theoretical predictions, rare services exhibit stronger market-size effects. They also exhibit stronger increasing returns to scale. Section 4.2 shows that division of labor and lumpy capital may be substantial sources of regional increasing returns for medical services. Section 4.3 directly connects trade to these mechanisms: imported care is more likely to be

³⁷Appendix A.1 discusses subtleties of the income data.

³⁸We observe 210 million claim lines representing \$89 billion in spending on professional services.

performed by specialists, by appropriate specialists, by more experienced physicians, and using scarce capital equipment, especially for smaller regions’ imports. Trade thus expands the population benefiting from specialization.

4.1 How market-size effects vary with procedure characteristics

4.1.1 Spatial variation in production and consumption by frequency

Larger regions produce a wider variety of procedures. In Figure 6(a), we count the number of procedures produced in each region and graph these counts within seven procedure categories against regional population. Larger regions produce a greater variety of care in all seven categories. Panel (b) shows that this is true for consumption, but to a lesser extent: Aggregating across all seven categories, HRRs in the top decile of population produce 3 times, but consume only 2 times, as many unique procedures as those in the bottom decile.

To better understand the characteristics of the procedures performed in larger markets, we estimate each procedure’s population elasticity of production per Medicare beneficiary.³⁹ Let Q_{pi} denote the count of procedure p produced in region i . Let M_i denote the number of Medicare beneficiaries residing in i . For each procedure p , we estimate the following relationship across regions:

$$\ln \mathbb{E} \left[\frac{Q_{pi}}{M_i} \right] = \zeta_p + \beta_p \ln \text{population}_i. \quad (13)$$

The estimated elasticity, $\hat{\beta}_p$, describes how production varies with market size, and we estimate it using Poisson pseudo-maximum-likelihood.⁴⁰ If the quantity produced were simply proportional to population, β_p would be zero. Our model suggests that scale effects play a larger role for rare procedures. It predicts less common services will have higher population

³⁹Davis and Dingel (2020) relate population elasticities to other measures of geographic concentration, such as location quotients, and estimate population elasticities of employment for various skills and sectors.

⁴⁰In a robustness check, we have also estimated a zero-inflated Poisson model, to account for the possibility that fixed costs are especially important for the decision of whether to provide the first instance of a service in a region. These results (not reported here) are quite similar.

elasticities of production.

Production per beneficiary indeed rises with market size, especially for less common procedures. The blue dots in Figure 6(c) relates the population elasticity of production per beneficiary $\hat{\beta}_p$ to the procedure’s national volume, $Q_p = \sum_i Q_{pi}$. Across values of Q_p , procedure output per beneficiary increases with market size. Less common procedures have higher elasticities, consistent with stronger economies of scale in rare procedures.

This finding raises questions about patients’ access to care. What happens to patients who live in smaller markets but need rare services? To investigate this, we separately estimate the population elasticity of *consumption* per beneficiary for each procedure. Let G_{pi} denote the count of procedure p consumed by patients *residing* in region i . We then estimate an analogue of equation (13) where we replace the dependent variable with $\ln \mathbb{E}[G_{pi}/M_i]$, and denote by β_p^C the resulting coefficient. If $\beta_p^C \neq \beta_p$, there is size-predicted net trade in procedure p .

The population elasticity of consumption per beneficiary is smaller for the vast majority of procedures and less steeply related to a procedure’s national frequency. The red squares in Figure 6(c) plot the population elasticity of consumption per beneficiary $\hat{\beta}_p^C$ of a procedure against its national volume $\ln Q_p$. While the relationship is negative, the slope for consumption is only one-third that for production. Appendix Table D.10 reports the production, consumption, and trade patterns for two exemplar procedures: screening colonoscopy and LVAD implantation. Colonoscopies are common and geographically dispersed, while LVAD procedures are rare, geographically concentrated, and traded over longer distances.

We have thus far modeled patients as demanding (and providers as producing) specific service codes. An alternative view is that patients have a particular medical condition that requires treatment, but the patients may not know what particular care they need; they simply know they require care. As physicians might use different treatments across regions for the same condition, our procedure-level results could reflect substitution across procedures. We address this by conducting a similar analysis across diagnoses.

Figure 6(d) shows production and consumption elasticities by diagnosis, rather than by procedure. The key patterns remain similar: production elasticities are higher, and decline more rapidly with national patient volume, than consumption elasticities. Both sets of elasticities have less steep relationships with national volume than for procedures. This could reflect measurement error within each category: the 482 diagnosis categories we use are far coarser than the 8,210 procedures in Figure 6(c). Alternatively, it could indicate true substitution among procedures within a condition that varies with location.

The contrasting population elasticities of production and consumption summarized in Figure 6 imply trade in medical services between markets of different sizes. Just as theories of trade with scale effects would predict, larger markets export rare procedures and smaller markets import them. For almost all procedures, production increases more than proportionately with market size. Consumption also increases more than proportionately with market size, but much less so than production. The differences between these elasticities mean net exports vary with market size. The implied net trade between markets of different sizes is particularly large for procedures that have small national volumes.

4.1.2 Market-size effects are stronger for rare and resource-intensive procedures

We now use trade flows to explicitly understand where rare medical services are provided and why. Since Section 3.3 found a strong home-market effect overall, the model predicts that this should be especially true for rare services. Trade is indeed more prevalent in rare services (Appendix Figure D.1). We now use the specific pattern of these trade flows to test the model’s prediction that scale economies drive the production and trade of rare services.

Table 3 reports gravity regressions in which each pair of locations has two observations: one for rare services and one for common. Column 1 repeats our baseline regression from Table 2 for professional fees and reports similar results. Column 2 limits the sample to pairs of location that have positive trade in at least one of the two procedure groups, which is the estimation sample used in the remainder of the table. We then estimate equation (12),

interacting both provider-market and patient-market population with an indicator for rare services.

The home-market effect is stronger for rare services. In column 3 of Table 3, the coefficient on provider-market population increases by about 60% relative to common services. The coefficient on patient-market population shrinks by more than half. Column 4 introduces location-pair fixed effects. Columns 5 and 6 are analogues of the previous two, but add a quadratic log-distance term. These results are statistically indistinguishable from the previous columns. Columns 7 and 8 demonstrate that our result holds when we look across diagnoses rather than procedures. As with the production and consumption elasticities in Figure 6(d), the magnitude of the difference between rare and common care shrinks. This could reflect substitution across care within a diagnosis or a less precise classification of diagnoses than of procedures. But the qualitative pattern holds and remains significant, consistent with the model’s prediction.

Appendix Table D.11 shows that these results are robust to using our instruments for market size. Columns 1 and 2 show estimates for common and rare services, respectively, when instrumenting for population in each region by its 1940 population. Columns 3 and 4 repeat the exercise using CBSAs rather than HRRs, and columns 5 and 6 switch to the bedrock-depth instrument. The results are consistent regardless of geographic unit or instrument.⁴¹ This stability suggests that neither the aggregate result nor the variation with procedure frequency is driven by anchor institutions or similar omitted variables.

The finding that less common procedures exhibit stronger home-market effects is robust to different ways of defining rare and common care. Appendix Figure D.5 shows that the pattern holds when splitting by deciles of national frequency. Appendix Table D.13 shows the same pattern among illustrative procedures.

The home-market effect is also stronger for rare procedures when controlling for how often an individual patient receives the same procedure, which we call a procedure’s “engagement”.

⁴¹Appendix Table D.12 reports analogous results based on commuting zones.

If patients are less willing to travel for high-engagement services and these services are more common, higher engagement could drive the stronger home-market effect we observe for rare procedures.⁴² Appendix Table D.14 shows that patients are indeed more sensitive to distance for high-engagement procedures, but controlling for this heterogeneity does not meaningfully alter the estimated differential impacts of population size for rare procedures.

The model predicts larger home-market effects for rare services when they have more elastic demand or larger scale economies. Fixed costs are one reason the scale elasticity may be larger for rarer services (see footnote 25). When estimating equation (9) separately by procedure frequency, we find that scale elasticity is substantially larger for rare services. It is near 1.0, as shown in the second panel of Appendix Table D.16.

Some procedures require more capital and specialized knowledge, and the consequences of market size for these resource-intensive procedures may illuminate sources of increasing returns. Relative value units (RVUs) capture the extent of physician work, expertise, and additional resources required to perform the procedure and determine the procedure’s Medicare payment. We divide procedures into terciles by intensity as measured by the total RVUs CMS assigns to a procedure. We find that the home-market effect is increasing in procedure intensity (Appendix Table D.15). This motivates our investigation of physician specialization and equipment usage in Section 4.2.

The potential concern about omitted cost shifters from Section 3.3 has an analogue here: Do the doctors who provide rare services benefit more from urban amenities than those providing common ones, lowering the cost of producing rare services in larger markets? This has facial plausibility if rare services are produced by elite specialists who earn more and might be more willing to pay for urban amenities through lower compensation.

Examining the population elasticities of physician earnings for each specialty alleviates this concern. If urban amenities drive specialists’ locations, earnings elasticities should be

⁴²In fact, the national frequency of a service has a very low correlation with various measures of engagement for that service, so it does not confound this result. The correlation between the share of patients who had more than one claim for the procedure in a given year and the procedure’s frequency is 0.14.

negative, especially for rare specialties. But Figure 7(a) shows that the income elasticities are close to zero on average and uncorrelated with the specialty’s national abundance. However urban amenities affect physicians’ choices, these choices do not exhibit the compensating differentials necessary to explain the relationship between market size and specialization.

4.2 Mechanisms for strong regional increasing returns

This section shows that division of labor and lumpy capital are mechanisms behind regional increasing returns to scale in medical services. This evidence does not preclude other agglomeration mechanisms from also playing a role. Knowledge diffusion (Baicker and Chandra, 2010) and thicker input markets could also be important productivity benefits of scale. We focus on specialization, physician experience, and equipment use because they can be measured using claims data and are closely related to procedure-level returns to scale.

4.2.1 Scale facilitates the division of labor

One source of increasing returns could be division of labor among physicians. In particular, the specialized labor required to produce rare services could drive the patterns we found across treatments and diagnoses. Specialized services may require physicians with specific training, whose presence may require high demand (Dranove, Shanley, and Simon, 1992).

To study this mechanism, we estimate the population elasticity of physicians per capita for each specialty and relate it to the number of physicians in the specialty. Let Y_{si} be the number of doctors of specialty s in location i .⁴³ We estimate a Poisson model,

$$\ln \mathbb{E} \left[\frac{Y_{si}}{\text{population}_i} \right] = \zeta_s^S + \beta_s^S \ln \text{population}_i, \quad (14)$$

for each specialty s by maximum likelihood.

Figure 7(b) shows a clear negative relationship between a specialty’s per capita population

⁴³Data are from the National Plan and Provider Enumeration System (NPPES); see Appendix A.1.

elasticity $\hat{\beta}_s^S$ and the national number of physicians in that specialty.⁴⁴ A natural explanation for rare procedures and rare specialties both being geographically concentrated in larger regions is that the size of the market limits the division of labor. To the extent that producing rare procedures requires specialized physicians, a larger volume of patients makes production economically viable.

4.2.2 Lumpy capital as a source of regional increasing returns

Another potential source of increasing returns is lumpy capital, namely expensive medical equipment. For instance, if a piece of medical equipment is rarely used in procedures, small markets may not use it enough to justify the investment.

We provide evidence consistent with this hypothesis using a CMS dataset that enumerates the types of equipment used by each medical procedure (HCPCS code).⁴⁵ We measure equipment use by linking this dataset to procedure-level production data. The frequency of equipment use H_{di} for a piece of equipment d in region i is defined as the sum of the procedure volume in region i across all procedures that use equipment d . We use a Poisson model to estimate the population elasticity of per-capita equipment use:

$$\ln \mathbb{E} \left[\frac{H_{di}}{\text{population}_i} \right] = \zeta_d^K + \beta_d^K \ln \text{population}_i. \quad (15)$$

Figure 7(c) plots the population elasticity of per-capita equipment use β_d^K against log national use, $\ln H_d$. The strong negative correlation indicates that equipment used less often tends to concentrate in larger markets. Similar to the patterns for physician specializations, rare capital equipment is disproportionately employed in larger regions.

⁴⁴This pattern is not attributable to spatial sorting driven by rare specialties commanding higher earnings. In fact, a specialty’s number of physicians and mean earnings are uncorrelated. Appendix Table D.17 shows that controlling for a specialty’s earnings has no effect on the negative relationship between population elasticity and number of physicians across specialties.

⁴⁵See Appendix A.1 for details.

4.3 Travel to access specialized services

We next ask whether the distribution of specialist physicians helps explain trade. Figure 8(a) shows the share of imports and of locally produced consumption that are provided by specialists as a function of regional population.⁴⁶ Imports are significantly more specialist-intensive than local production. This difference is especially pronounced in the smallest regions, and it remains true throughout the population distribution.

Does trade match patients with the appropriate specialist? Among all specialty care, we determine the two most common specialties to provide each unique service and label these the “standard” specialties for that care. We then determine whether each instance of the treatment was provided by a standard or non-standard specialist.

Figure 8(b) shows the share of imports and of locally produced care provided by non-standard specialties. Imports are less likely to come from non-standard specialists than local care, and the distinction is especially pronounced in the smallest regions. The difference is substantial: Local care in the smallest regions is 40% more likely to be provided by a non-standard specialist than in the largest regions (7% vs. 5%). When importing medical services, this share falls to 5%—indistinguishable from the largest regions’ locally produced care.

We conduct a similar analysis based on provider experience. Using the public Medicare provider data (based on all Traditional Medicare patients), we count the number of times the physician billed for the specific service in the previous year. We divide this experience measure by the procedure’s national mean and average it across all procedures provided to patients in an HRR. We then rescale HRR-level experience by the mean across HRRs. Figure 8(c) shows that, at all population sizes, care imported from other regions is produced by more experienced providers than locally produced care.⁴⁷ Patients in larger regions see

⁴⁶We define “specialist” to mean all physicians except those whose primary specialty is internal medicine, general practice, or family practice.

⁴⁷This comparison restricts attention to the 143 procedures performed in all HRRs. Thus, regional variation does not reflect the fact that larger markets produce a greater number of distinct codes (Figure 6(a)).

more experienced providers for both imported and locally produced care.

Specialists are disproportionately located in larger markets, as are physicians with more experience in their procedures. Since imported care is predominantly specialized and provides patients access to higher experience, we conclude that visiting the appropriate specialist based on training or experience is part of the value proposition for trade in medical care.

Figure 8(d) shows analogous patterns for rare capital equipment. We measure the share of care that requires a piece of rare capital equipment and plot these shares separately for local and imported care against regional population. While locally produced care in larger markets is more likely to use rare equipment, imported care has a higher rare-equipment share throughout the population distribution.

This analysis also provides a second validation of our interpretation that trade reflects quality variation. Patients travel to regions with highly-ranked hospitals, which larger markets tend to have—along with rare equipment and the ability to provide rare services. This market-size effect strongly predicts gross and net exports. Together, this suggests that economies of scale play an important role in increasing the quality of care, and trade between regions enables patients from many regions to share the benefits of this agglomeration.

5 Health policy with trade and increasing returns

Given the estimated strength of regional increasing returns, geographically concentrating healthcare production has substantial benefits. Larger regions support more specialists, physician experience, specialized equipment, and procedures. But geographic concentration implies that patients in smaller regions may have limited access to care. We use our estimates of the scale elasticity α , region-specific qualities δ_i , and observed trade flows to quantify how various counterfactual policy scenarios would change each region’s patient market access for non-emergency care.⁴⁸ Our results underline the importance of distinguishing between the

⁴⁸Because our estimation sample excludes emergency services, our estimates and counterfactual scenarios omit any complementarities between emergency and non-emergency care. Our partial-equilibrium model

quality of locally produced services and the quality of services to which local residents have access. They also show how geography may contribute to striking patterns of inequality in Americans' health.

5.1 Method to compute counterfactual outcomes

We compute counterfactual equilibrium outcomes relative to the baseline equilibrium. For the baseline equilibrium, define export shares $x_{ij} \equiv \frac{Q_{ij}}{\sum_{j'} Q_{ij'}}$ and import shares $m_{ij} \equiv \frac{Q_{ij}}{N_j}$. For every variable or parameter y , denote the ratio of its counterfactual value y' to its baseline value y by $\hat{y} \equiv \frac{y'}{y}$. Appendix C.1 shows how we solve for the relative counterfactual endogenous qualities ($\hat{\delta}$) using baseline equilibrium shares (x_{ij}, m_{ij}), the scale elasticity (α), and relative counterfactual exogenous parameters ($\hat{A}, \hat{R}, \hat{w}, \hat{\rho}, \hat{N}$). In particular, counterfactual qualities are given by a system of \mathcal{I} equations with unknowns $\{\hat{\delta}_i\}_{i=1}^{\mathcal{I}}$:

$$\hat{\delta}_i = \left(\hat{R}_i \hat{A}_i / \hat{w}_i \right)^{\frac{1}{1-\alpha}} \left(\sum_{j \in \mathcal{I}} \frac{x_{ij} \hat{\rho}_{ij} \hat{N}_j}{m_{0j} + \sum_{i' \in \mathcal{I}} m_{i'j} \hat{\delta}_{i'} \hat{\rho}_{i'j}} \right)^{\frac{\alpha}{1-\alpha}}.$$

The first term of this expression, $\left(\hat{R}_i \hat{A}_i / \hat{w}_i \right)^{\frac{1}{1-\alpha}}$, shows that the scale elasticity α governs the effect of exogenous supplier shifters, including reimbursements \hat{R}_i , on quality produced in a region. Reimbursement rates shift the scale of production, and stronger scale economies (higher α) amplify these shifts. The second term shows how changes in other regions influence local outcomes through trade, combined with scale. Thus, our counterfactual scenarios rely on both our estimates of the scale elasticity α and observed trade patterns.^{49,50}

assumes elastic factor supplies. Our model and estimates represent long-run elasticities, omitting any adjustment costs or short-run diseconomies of scale due to crowding or queuing.

⁴⁹The qualitative results are not sensitive to our exact estimate of α . For instance, cutting α by half changes the magnitudes but not the geographic patterns.

⁵⁰To compute import shares, we assume that the number of potential patients is proportional to the number of Traditional Medicare beneficiaries and infer the share choosing the outside option in each region. See Appendix C.2. The qualitative and spatial patterns of counterfactual outcomes do not depend on what share we assume choose the outside option. Appendices C.3 and C.4 generalize this method of computing counterfactual outcomes to the model with multiple types of patients introduced in Appendix B.2.

5.2 Counterfactual reimbursement rates and travel costs

Local reimbursement increase. Policymakers concerned about a region’s healthcare access might consider local production subsidies. Figure 9 contrasts the consequences of raising reimbursements by 30% in Boston and in Paducah, Ky. For the Boston scenario, Panel (a) depicts the impact on quality of care in each region relative to its baseline value. Free entry means that higher reimbursements translate to higher-quality care produced in Boston. Quality declines in the rest of New England as patients substitute away and scale economies translate lower volumes into lower quality (an “agglomeration shadow”, as in Fujita and Krugman, 1995). These effects diminish with distance to Boston.

Yet regions with larger declines in output quality due to Boston’s expansion experience larger improvements in patient market access, $\hat{\Phi}_i$ (Panel 9(b)). This measure accounts for changes in the quality of care patients receive and their costs of traveling when they choose to do so. Patients in Boston benefit the most from the higher reimbursement of their local production. Outside Boston, regional changes in patient market access are nearly opposite the changes in local output quality. Nearby regions import enough that the benefits of improved quality in Boston exceed the declines in the quality of local production, so their patient market access improves. Regions closer to Boston experience larger declines in the quality of local production precisely because their residents’ choice sets improve more, spurring more substitution. In more distant regions, the welfare impacts are virtually zero.

Patients who live in higher-income neighborhoods benefit more. The value of market access increases by 120% more for third- than for first-tercile patients nationwide (Appendix Table D.18). This reflects greater consumption of Boston production (column 3).

The consequences of higher reimbursement rates in Paducah, Ky. exhibit very different spatial patterns than in Boston. Figures 9(c) and 9(d) depict the regional changes in output quality and patient market access, respectively, caused by a 30% reimbursement increase in Paducah. Unlike Boston, Paducah is a net importer: its consumption of medical services exceeds local production by more than one-third. Higher reimbursements that improve out-

put quality in Paducah cause Paducahans to reduce their imports from neighboring regions. This reduces the quantity produced in neighboring regions, lowering their output quality. But Panel (d) shows that—unlike the pattern of outcomes in the Boston scenario—those regions where output quality declines more are the regions where patient market access declines more.

The contrasting outcomes reflect trade flows in the baseline equilibrium: Boston is a net exporter of medical services and Paducah is a net importer. Higher reimbursements in Boston cause output quality declines in nearby regions—largely because residents of those regions import more when Boston’s quality improves. In contrast, higher reimbursements in Paducah reduce neighboring regions’ output quality largely because Paducah residents demand fewer exports from these regions when Paducah’s quality improves. Nearby regions import little from Paducah, so they benefit little from its improved quality. Appendix Figure D.6 shows that the lessons from Boston and Paducah generalize: the pattern of spillovers from increasing reimbursements in one region is driven by that region’s net trade in medical care. To summarize, the spillover consequences of subsidizing production in one region depend on trade patterns; changes in regional output quality need not align with changes in regional patient market access.

The distributional consequences of region-specific subsidies depend on which region is subsidized. We compute the nationwide gains in market access from subsidizing production in each region, one at a time. Figure 10 shows this gain, scaled by the increase in total spending, as a function of region size. The aggregate gain in market access per dollar spent is higher in larger markets: further concentration of production has larger benefits. The graph also shows the gains per dollar separately by income tercile. Subsidizing production in less populous regions benefits lower-income ZIP codes more. These contrasts reflect geographic divides in incomes: lower-income patients are more likely to live in and near smaller regions.

Reducing travel costs. Rather than subsidizing local production, policies might improve patient market access in a particular region by facilitating trade. We examine the consequences of a policy that reduces travel costs for Paducahans obtaining care elsewhere.⁵¹ Figure 11 shows that, unlike an increase in Paducah reimbursements, this policy has positive spillovers on neighboring regions. These regions increase their exports to Paducah, and thus their own scale and quality. This improves their residents’ market access.

Reducing travel costs for Paducahans benefits both Paducah and its neighbors—though we do not estimate the fiscal costs of subsidizing travel. Paducahans benefit even though facilitating travel reduces the quantity—and thus the quality—produced in Paducah. Analysts looking at the impact of travel subsidies on the quantity or quality of care provided in Paducah itself would reach very different conclusions than those looking at the impact on patient market access.

These counterfactual scenarios are subject to significant caveats, and we have not attempted to identify the optimal policy. Even so, this simple model rationalizes important aspects of the economic geography of US healthcare policy. The counterfactual scenarios highlight our main findings: Healthcare production has substantial local increasing returns, and patient travel plays a meaningful role in enabling access to higher-quality care. Given these economic mechanisms, regional spillovers are larger when economies of scale are stronger, depend on the pattern of trade flows, and differ depending on whether policies subsidize production or travel. This shows the importance of distinguishing between regional output quality and regional patient access when evaluating healthcare policies.

5.3 Geography and healthcare inequality across countries

Looking at healthcare geography more broadly, a striking fact is that the United States has a steeper relationship between regional income and mortality than other rich economies.

⁵¹Specifically, $\hat{\rho}_{i,\text{Paducah}} = 1.3$ when $i \neq \text{Paducah}$. The impact of this change on Paducah residents’ market access $\hat{\Phi}_{\text{Paducah}}$ is similar to a 3% increase in reimbursements in Paducah.

Appendix Figure D.8(a) compares the US health-income gradient with Germany’s. We use our model to explore how much geography can explain this difference in gradients.

To do so, we compute the changes in trade costs that would make the distribution of market potential across US regions comparable to that across German regions.⁵² This requires broad reductions in travel costs, effectively condensing the United States. It requires particularly large declines in travel costs for remote US regions.

This transformation reduces the Φ_i -income elasticity by 21% (Appendix Figure D.8(b)). This suggests that a substantial part of the inequality in market access can be explained by the US’s greater geographic area. If the US were Germany’s size, patients from lower-income areas would see particular benefit: these areas are disproportionately rural, and thus experience the largest gains in this counterfactual. Viewed from the opposite perspective, 27% of the larger Φ_i -income elasticity in the United States is because of its geographic scope.

6 Conclusion

Smaller markets have fewer specialized physicians, produce less medical care per capita, and have worse health outcomes than larger markets. Thanks to trade in medical services, less production does not translate one for one into less consumption of medical services. Instead, trade affords patients who live in smaller markets access to higher-quality care. This quality comes in part from consuming services that would otherwise be unavailable, visiting appropriate specialists, and accessing experienced physicians.

This trade amplifies the scale advantages of large markets and hence the quality of care they produce. This means the healthcare industry can serve as an export base for large cities. Substantial scale economies also imply that policies to reallocate care across regions may impact the quality of care available. The rich and varied patterns of consequences when subsidizing production or travel in “under-served” markets highlight the importance of trade and agglomeration for the incidence of these policies on patients and producers.

⁵²See Appendix C.5 for details. The analysis is similar when we compare the United States to France.

References

- AAPC. 2021. “What Is CPT?” Available online at <https://www.aapc.com/resources/medical-coding/cpt.aspx>.
- Acemoglu, Daron, Amy Finkelstein, and Matthew J. Notowidigdo. 2013. “Income and Health Spending: Evidence from Oil Price Shocks.” *The Review of Economics and Statistics*, 95(4): 1079–1095.
- Acemoglu, Daron, and Joshua Linn. 2004. “Market Size in Innovation: Theory and Evidence from the Pharmaceutical Industry.” *The Quarterly Journal of Economics*, 119(3): 1049–1090.
- Agarwal, Sumit, J. Bradford Jensen, and Ferdinando Monte. 2020. “Consumer Mobility and the Local Structure of Consumption Industries.” NBER Working Paper 23616.
- Allen, Treb, Simon Fuchs, Sharat Ganapati, Alberto Graziano, Rocio Madera, and Judit Montoriol-Garriga. 2021. “Urban Welfare: Tourism in Barcelona.”
- Anderson, James E., Catherine A. Milot, and Yoto V. Yotov. 2014. “How much does geography deflect services trade? Canadian answers.” *Int’l Econ. Rev.*, 55(3): 791–818.
- Arrow, Kenneth J. 1963. “Uncertainty and the Welfare Economics of Medical Care.” *American Economic Review*, 53(5): 941–973.
- Baicker, Katherine, and Amitabh Chandra. 2010. “Understanding Agglomerations in Health Care.” *Agglomeration Economics*, 211–236. University of Chicago Press.
- Bartelme, Dominick G., Arnaud Costinot, Dave Donaldson, and Andrés Rodríguez-Clare. 2019. “The Textbook Case for Industrial Policy: Theory Meets Data.” NBER Working Paper 26193.
- Bartik, Timothy, and George Ericckek. 2007. “Higher Education, the Health Care Industry, and Metropolitan Regional Economic Development: What Can ‘Eds & Meds’ Do for the Economic Fortunes of a Metro Area’s Residents?” Upjohn Inst. Working Paper 08-140.
- Baumgardner, James R. 1988. “Physicians’ Services and the Division of Labor across Local Markets.” *Journal of Political Economy*, 96(5): 948–982.
- Berenson, Robert A., Jonathan H. Sunshine, David Helms, and Emily Lawton. 2015. “Why Medicare Advantage plans pay hospitals traditional Medicare prices.” *Health Affairs*, 34(8): 1289–1295.
- Berry, Christopher R., and Edward L. Glaeser. 2005. “The divergence of human capital levels across cities.” *Papers in Regional Science*, 84(3): 407–444.
- Borchert, Ingo, Mario Larch, Serge Shikher, and Yoto V. Yotov. 2021. “The International Trade and Production Database for Estimation (ITPD-E).” *International Economics*, 166(C): 140–166.

- Burstein, Ariel, Sarah Lein, and Jonathan Vogel. 2022. “Cross-border shopping: evidence and welfare implications for Switzerland.”
- Cabral, Marika, and Neale Mahoney. 2019. “Externalities and taxation of supplemental insurance: A study of Medicare and Medigap.” *American Economic Journal: Applied Economics*, 11(2): 37–73.
- Centers for Medicare and Medicaid Services. 2022. “National Health Expenditure Data.”
- Chandra, Amitabh, Maurice Dalton, and Douglas O Staiger. 2023. “Are Hospital Quality Indicators Causal?” NBER Working Paper 31789.
- Chipman, John S. 1970. “External Economies of Scale and Competitive Equilibrium.” *The Quarterly Journal of Economics*, 84(3): 347–385.
- Clemens, Jeffrey, and Joshua D. Gottlieb. 2014. “Do Physicians’ Financial Incentives Affect Treatment Patterns and Patient Health?” *American Economic Review*, 104(4): 1320–1349.
- Costinot, Arnaud, Dave Donaldson, Margaret Kyle, and Heidi Williams. 2019. “The More We Die, The More We Sell? A Simple Test of the Home-Market Effect.” *The Quarterly Journal of Economics*, 134(2): 843–894.
- Davis, Donald R., and David E. Weinstein. 2003. “Market access, economic geography and comparative advantage: an empirical test.” *J. of International Economics*, 59(1): 1–23.
- Davis, Donald R., and Jonathan I. Dingel. 2020. “The comparative advantage of cities.” *Journal of International Economics*, 123(C).
- Davis, Donald R., Jonathan I. Dingel, Joan Monras, and Eduardo Morales. 2019. “How Segregated Is Urban Consumption?” *Journal of Political Economy*, 127(4): 1684–1738.
- Diamond, Rebecca. 2016. “The determinants and welfare implications of US workers’ diverging location choices by skill: 1980-2000.” *American Economic Review*, 106(3): 479–524.
- Dingel, Jonathan I. 2017. “The Determinants of Quality Specialization.” *Review of Economic Studies*, 84(4): 1551–1582.
- Dingel, Jonathan I., and Felix Tintelnot. 2021. “Spatial Economics for Granular Settings.” NBER Working Paper 27287.
- Disdier, Anne-Célia, and Keith Head. 2008. “The Puzzling Persistence of the Distance Effect on Bilateral Trade.” *The Review of Economics and Statistics*, 90(1): 37–48.
- Dorsey, E. Ray, and Eric J. Topol. 2016. “State of Telehealth.” *New England Journal of Medicine*, 375(2): 154–161. PMID: 27410924.
- Dranove, David, and Mark A. Satterthwaite. 2000. “The Industrial Organization of Health Care Markets.” *Handbook of Health Economics*, 1: 1093–1139.

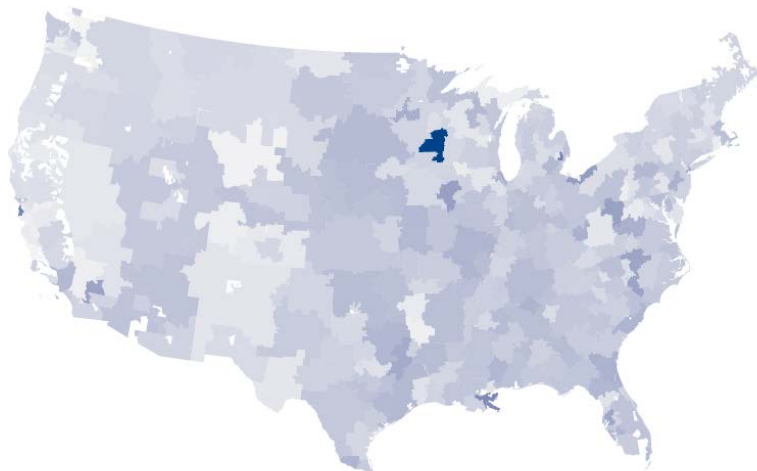
- Dranove, David, Mark Shanley, and Carol Simon. 1992. “Is Hospital Competition Wasteful?” *The RAND Journal of Economics*, 23(2): 247–262.
- Eckert, Fabian, Sharat Ganapati, and Conor Walsh. 2020. “Skilled Scalable Services: The New Urban Bias in Economic Growth.” CESifo Working Paper Series 8705.
- Finkelstein, Amy, Matthew Gentzkow, and Heidi Williams. 2016. “Sources of geographic variation in health care: Evidence from patient migration.” *QJE*, 131(4): 1681–1726.
- Finkelstein, Amy, Matthew Gentzkow, and Heidi Williams. 2021. “Place-Based Drivers of Mortality: Evidence from Migration.” *American Economic Review*, 111(8): 2697–2735.
- Fisher, Elliott S., David E. Wennberg, Thérèse A. Stukel, Daniel J. Gottlieb, F. Lee Lucas, and Etoile L. Pinder. 2003a. “The implications of regional variations in Medicare spending. Part 1: the content, quality, and accessibility of care.” *Ann. Intern. Med.*, 138(4): 273–287.
- Fisher, Elliott S., David E. Wennberg, Thérèse A. Stukel, Daniel J. Gottlieb, F. Lee Lucas, and Etoile L. Pinder. 2003b. “The implications of regional variations in Medicare spending. Part 2: health outcomes and satisfaction with care.” *Ann. Intern. Med.*, 138(4): 288–298.
- Flexner, Abraham. 1910. “Medical Education in the United States and Canada: A Report to the Carnegie Foundation for the Advancement of Teaching.” Carnegie Fnd. Bulletin 4.
- Fowler, Christopher S, and Leif Jensen. 2020. “Bridging the gap between geographic concept and the data we have: The case of labor markets in the USA.” *Environment and Planning A: Economy and Space*, 52(7): 1395–1414.
- Fujita, Masahisa, and Paul Krugman. 1995. “When is the economy monocentric?: von Thünen and Chamberlin unified.” *Regional Science and Urban Econ.*, 25(4): 505–528.
- Gaynor, Martin, and Robert J Town. 2011. “Competition in health care markets.” *Handbook of Health Economics* Vol. 2, 499–637. Elsevier.
- Gaynor, Martin, Kate Ho, and Robert J Town. 2015. “The industrial organization of health-care markets.” *Journal of Economic Literature*, 53(2): 235–284.
- Gottlieb, Daniel J, Weiping Zhou, Yunjie Song, Kathryn Gilman Andrews, Jonathan S Skinner, and Jason M Sutherland. 2010. “Prices don’t drive regional Medicare spending variations.” *Health Affairs*, 29(3): 537–543.
- Gottlieb, Joshua D., Maria Polyakova, Kevin Rinz, Hugh Shiple, and Victoria Udalova. 2023. “Who Values Human Capitalists’ Human Capital? The Earnings and Labor Supply of U.S. Physicians.” National Bureau of Economic Research Working Paper 31469.
- Gotto, Antonio M., and Jennifer Moon. 2016. *Weill Cornell Medicine: A History of Cornell’s Medical School*. Ithaca: Cornell University Press.
- Greenstone, Michael, Richard Hornbeck, and Enrico Moretti. 2010. “Identifying agglomeration spillovers: Evidence from winners and losers of large plant openings.” *Journal of Political Economy*, 118(3): 536–598.

- Gupta, Atul. 2021. “Impacts of Performance Pay for Hospitals: The Readmissions Reduction Program.” *American Economic Review*, 111(4): 1241–83.
- Hanson, Gordon H., and Chong Xiang. 2004. “The Home-Market Effect and Bilateral Trade Patterns.” *American Economic Review*, 94(4): 1108–1129.
- Helpman, Elhanan, and Paul R. Krugman. 1985. *Market Structure and Foreign Trade*. MIT Press.
- Hsieh, Chang-Tai, and Esteban Rossi-Hansberg. 2021. “The Industrial Revolution in Services.” Center for Economic Studies, U.S. Census Bureau Working Papers 21-34.
- Jensen, J. Bradford, and Lori G. Kletzer. 2005. “Tradable Services: Understanding the Scope and Impact of Services Offshoring.” *Brookings Trade Forum*, 75–116.
- Kline, Patrick, and Enrico Moretti. 2013. “Local Economic Development, Agglomeration Economies, and the Big Push: 100 Years of Evidence from the Tennessee Valley Authority.” *The Quarterly Journal of Economics*, 129(1): 275–331.
- Krugman, Paul. 1980. “Scale Economies, Product Differentiation, and the Pattern of Trade.” *American Economic Review*, 70(5): 950–59.
- Lee, Sanghoon. 2010. “Ability sorting and consumer city.” *Journal of Urban Economics*, 68(1): 20–33.
- Levy, Antoine, and Jacob Moscona. 2020. “Specializing in Density: Spatial Sorting and the Pattern of Trade.” Mimeo, MIT.
- Lipsey, Robert E. 2009. “Measuring International Trade in Services.” In *International Trade in Services and Intangibles in the Era of Globalization*, ed. Marshall B. Reinsdorf and Matthew J. Slaughter, 27–74. University of Chicago Press.
- Lopez, Eric, and Gretchen Jacobson. 2020. “How Much More Than Medicare Do Private Insurers Pay? A Review of the Literature.”
- Marshall, Alfred. 1890. *Principles of Economics*. London: Macmillan and Co.
- Miyauchi, Yuhei, Kentaro Nakajima, and Stephen J. Redding. 2021. “The Economics of Spatial Mobility: Theory and Evidence Using Smartphone Data.” NBER Working Paper 28497.
- Monte, Ferdinando, Stephen J Redding, and Esteban Rossi-Hansberg. 2018. “Commuting, migration, and local employment elasticities.” *American Economic Review*, 108(12): 3855–3890.
- Montiel Olea, José Luis, and Carolin Pflueger. 2013. “A robust test for weak instruments.” *Journal of Business & Economic Statistics*, 31(3): 358–369.
- Moretti, Enrico. 2004. “Estimating the social return to higher education: evidence from longitudinal and repeated cross-sectional data.” *J. of Econometrics*, 121(1-2): 175–212.

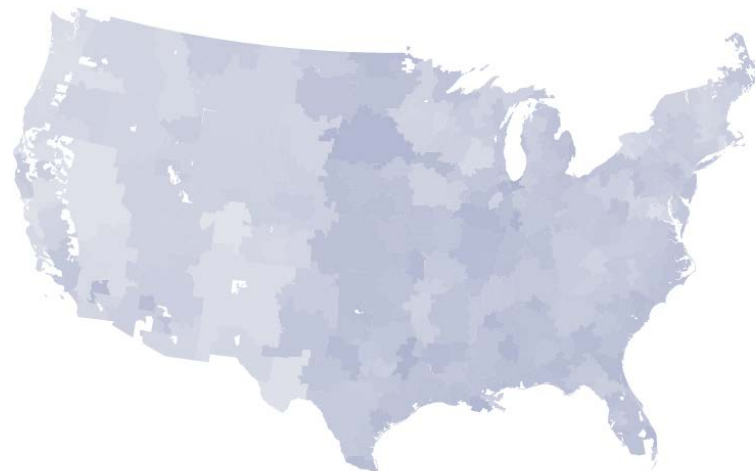
- Moretti, Enrico. 2011. "Local labor markets." In *Handbook of Labor Economics*. Vol. 4, 1237–1313. Elsevier.
- Mourot, Pauline. 2024. "Should Top Surgeons Practice at Top Hospitals? Sorting and Complementarities in Healthcare." University of Chicago, mimeo.
- Muñoz, Mathilde. 2022. "Trading Non-Tradables: The Implications of Europe's Job Posting Policy."
- Neufeld, Jonathan D., and Charles R. Doarn. 2015. "Telemedicine Spending by Medicare: A Snapshot from 2012." *Telemed J E Health*, 21(8): 686–693.
- Newhouse, Joseph P. 1990. "Geographic access to physician services." *Annual Review of Public Health*, 11(1): 207–230.
- Newhouse, Joseph P., Albert P. Williams, Bruce W. Bennett, and William B. Schwartz. 1982a. "Does the Geographical Distribution of Physicians Reflect Market Failure?" *The Bell Journal of Economics*, 13(2): 493–505.
- Newhouse, Joseph P., Albert P. Williams, Bruce W. Bennett, and William B. Schwartz. 1982b. "The geographic distribution of physicians: Is the conventional wisdom correct?" Santa Monica: RAND Corp. Publ. No. R-2734.
- Newhouse, Joseph P., Albert P. Williams, Bruce W. Bennett, and William B. Schwartz. 1982c. "Where have all the doctors gone?" *JAMA*, 247(17): 2392–2396.
- Rosenblatt, R. A., and L. G. Hart. 2000. "Physicians and rural America." *The Western Journal of Medicine*, 173(5): 348–351.
- Rosenthal, Meredith B., Alan Zaslavsky, and Joseph P. Newhouse. 2005. "The Geographic Distribution of Physicians Revisited." *Health Services Research*, 40(6p1): 1931–1952.
- Rosenthal, Stuart S., and William C. Strange. 2008. "The attenuation of human capital spillovers." *Journal of Urban Economics*, 64(2): 373–389.
- Ruggles, Steven, Sarah Flood, Ronald Goeken, Megan Schouweiler, and Matthew Sobek. 2022. "IPUMS USA: Version 12.0." Minneapolis, MN: IPUMS.
- Santos Silva, J. M. C., and Silvana Tenreyro. 2006. "The Log of Gravity." *The Review of Economics and Statistics*, 88(4): 641–658.
- Silver, David, and Jonathan Zhang. 2022. "Impacts of Basic Income on Health and Economic Well-Being: Evidence from the VA's Disability Compensation Program." NBER Working Paper 29877.
- Skinner, Lucy, Douglas O. Staiger, David I. Auerbach, and Peter I. Buerhaus. 2019. "Implications of an aging rural physician workforce." *NEJM*, 381(4): 299–301.
- Zuckerman, Stephen, Laura Skopec, and Joshua Aarons. 2021. "Medicaid Physician Fees Remained Substantially Below Fees Paid By Medicare In 2019." *Health Affairs*, 40(2): 343–348.

Figure 1: Production, consumption, and trade across regions

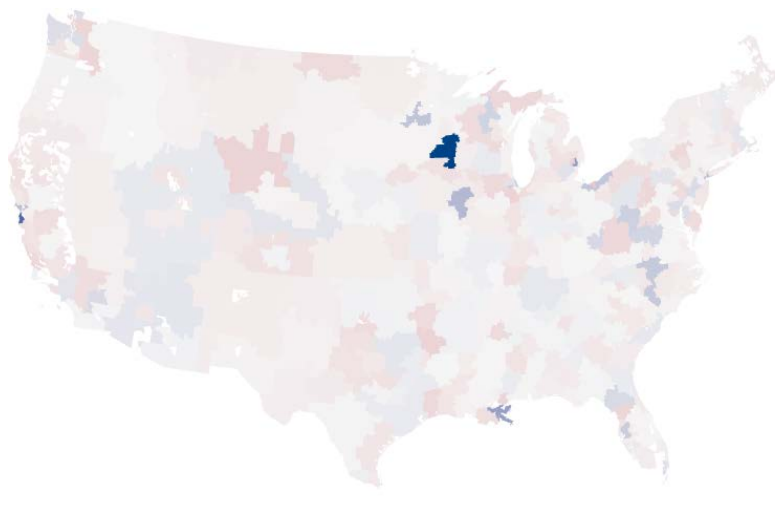
(a) Production per capita



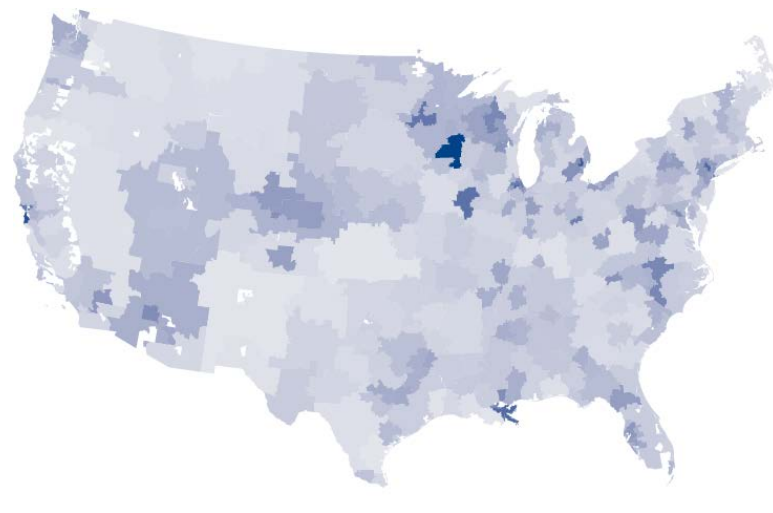
(b) Consumption per capita



(c) Production divided by consumption

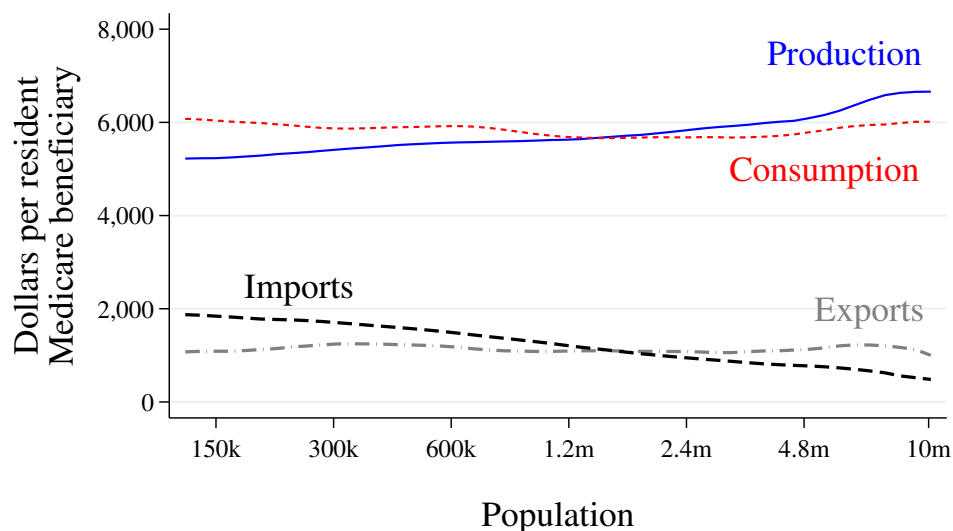


(d) Gross exports relative to production



Notes: These maps depict production, consumption, and trade by hospital referral region (HRR). Panel (a) shows dollars of production per capita. The HRR of production is the location where the service is provided. Panel (b) shows dollars of consumption per capita. The HRR of consumption is based on the patient's residential address. Panel (c) shows the ratio of production per capita to consumption per capita. Panel (d) shows gross exports as a share of production. Section 1.1 and Appendices A.1 and A.2 detail how we compute trade flows in physician services (excluding emergency-room care and skilled nursing facilities) at standardized prices from the Medicare 20% carrier, 100% MedPAR, and 100% outpatient claims Research Identifiable Files. HRR definitions are from the Dartmouth Atlas Project. The Anchorage and Honolulu HRRs are not depicted.

Figure 2: Production and consumption of medical care across regions

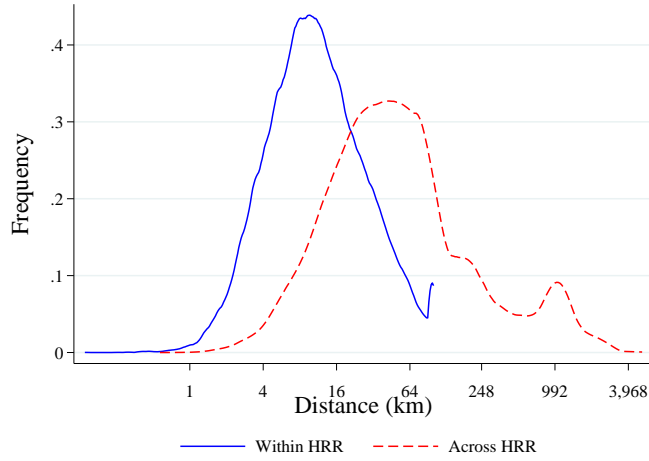


Population elasticity (log–log regression slope) of transactions per resident Medicare beneficiary:
 Production: 0.05 (0.01), Consumption: -0.02 (0.01)
 Exports: -0.02 (0.04), Imports: -0.35 (0.03)

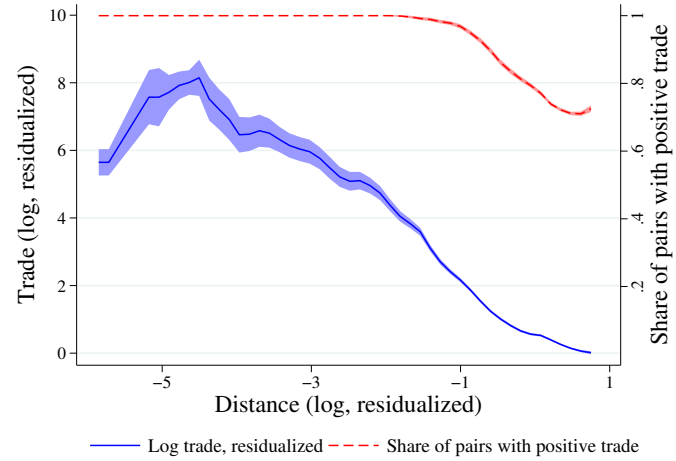
Notes: This figure shows production, consumption, and trade per capita of Medicare services across hospital referral regions (HRRs) of different sizes, all smoothed via local averages. The blue series shows production of medical care per Medicare beneficiary residing in the HRR of production. The red series shows consumption of medical care per Medicare beneficiary residing in the HRR of consumption. The dashed gray series shows interregional “exports” of medical care and the dashed black series shows interregional “imports” of medical care, again per Medicare beneficiary. Section 1.1 and Appendices A.1 and A.2 detail how we compute trade flows in physician services (excluding emergency-room care and skilled nursing facilities) at standardized prices from the Medicare 20% carrier, 100% MedPAR, and 100% outpatient claims Research Identifiable Files. HRR definitions are from the Dartmouth Atlas Project.

Figure 3: Patients travel between regions and trade declines with distance, more so for lower-income patients

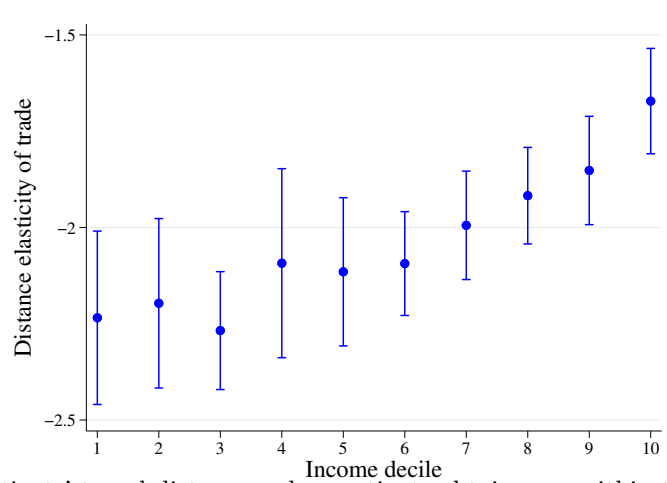
(a) Distribution of travel distances within and across HRRs



(b) Trade volume and extensive margin by distance

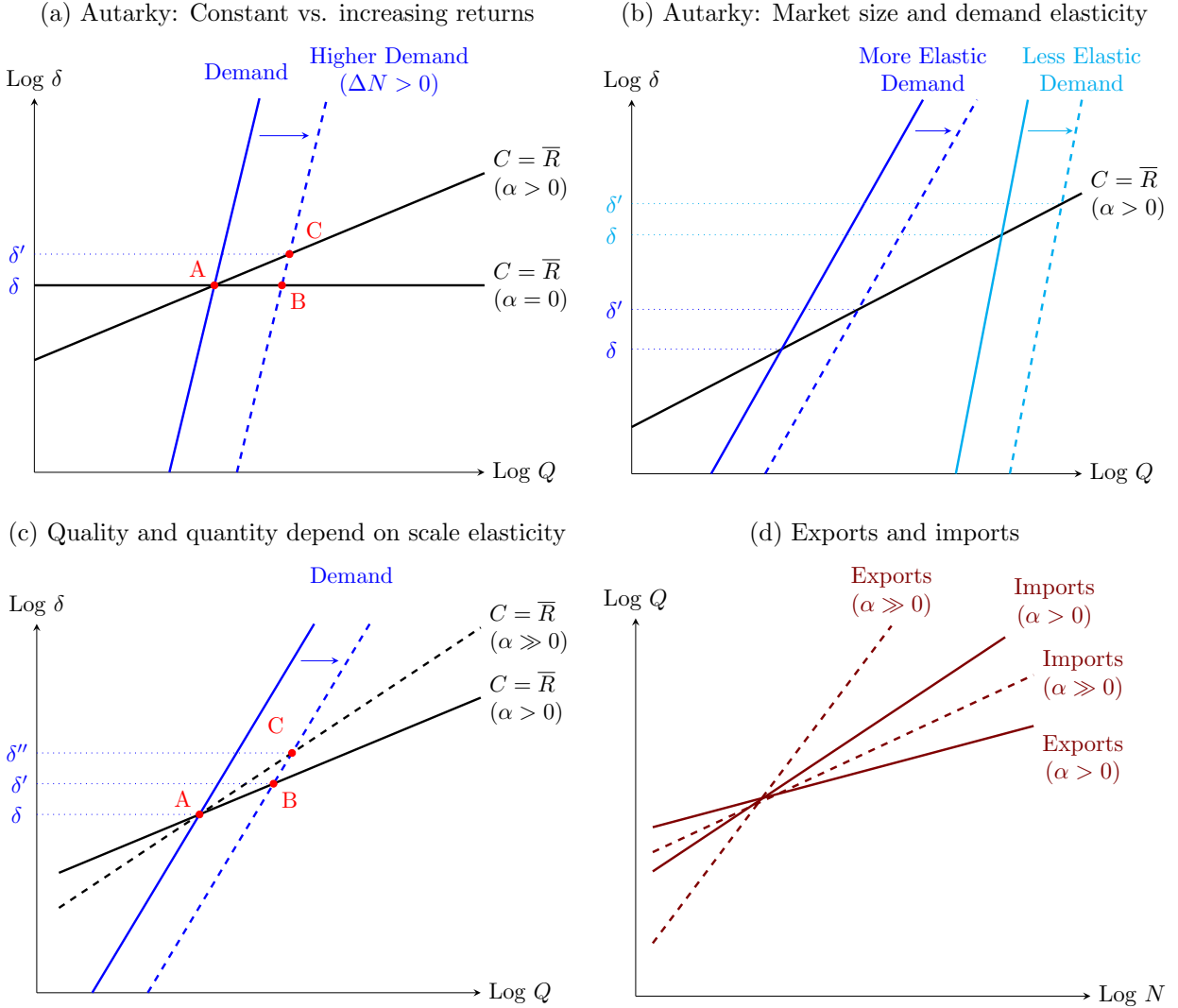


(c) Higher-income patients are less sensitive to distance



Notes: Panel (a) shows the distribution of patients' travel distances when patients obtain care within their home HRR (blue distribution) and when they travel across HRRs (red distribution). Travel distances measure the distance between home and treatment locations. For travel within a hospital referral region, the distance measure reflects the distance between the centroid of the patient's residential ZIP code and the ZIP code of the service location. We use ZCTA-to-ZCTA distances downloaded from the National Bureau of Economic Research; those exceeding 160 kilometers are winsorized at 160 kilometers. For travel across HRRs, we use ZCTA-to-ZCTA distances when they are within 160 kilometers and (for computational ease) use HRR-to-HRR distances beyond 160 kilometers. In Panel (b), the blue series depicts the volume of trade against distance, after conditioning out the fixed effects in equation (8), for positive-trade pairs of locations. The red series shows the share of HRR pairs with positive trade as a function of the distance between them, after conditioning out the importer fixed effects and exporter fixed effects, as in equation (8). Panel (c) depicts the coefficient on log distance obtained by estimating equation (8) separately for each decile of the national ZIP-level median-household-income distribution. The 95% confidence intervals are computed using standard errors two-way clustered by both patient HRR and provider HRR. Patients from higher-income ZIP codes are less sensitive to distance.

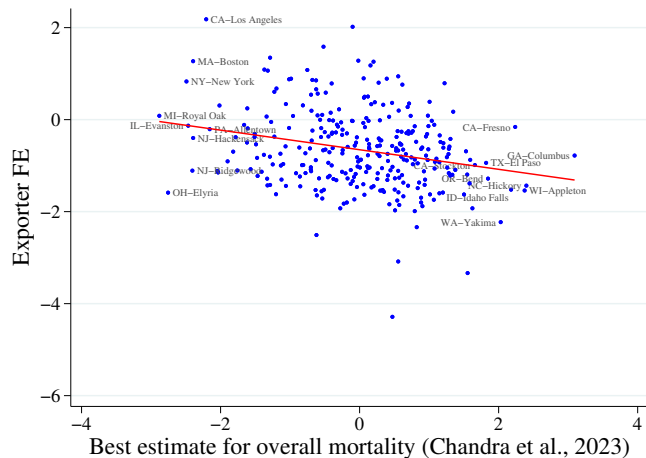
Figure 4: Illustrative model diagrams



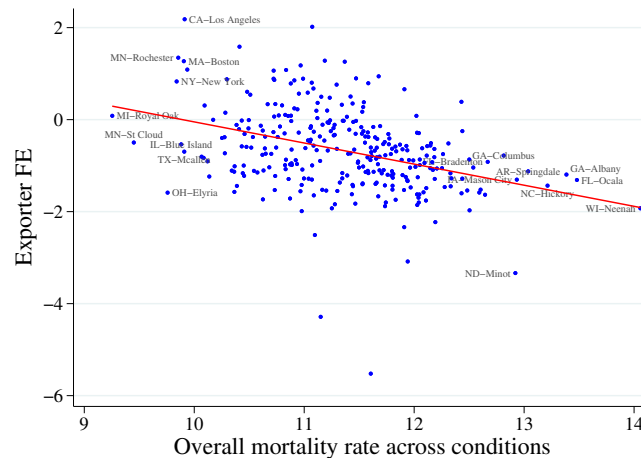
Notes: This figure depicts how increasing demand in one region affects its equilibrium outcomes. In Panels (a)–(c), quantity produced Q is on the horizontal axis and service quality δ is on the vertical axis. The black lines depict the free-entry isocost curve, $C = \bar{R}$, given by equation (3). The blue and cyan lines depict demand for the region’s service, which we depict as log-linear for visual clarity. (The logit demand function is actually log-convex, which is consistent with all the depicted comparative statistics.) Equilibrium is the intersection of the demand and isocost curves. An increase in demand is the rightward shift from the solid to the dashed demand curve. This shift increases equilibrium quality from δ to δ' . Panel (a) shows that higher demand elicits higher quality if there are increasing returns to scale. Panel (b) shows that this quality improvement is larger when demand is more elastic. Panels (c) and (d) introduce trade and compare the extent of quality improvement under two different magnitudes of increasing returns ($\alpha > 0$ and $\alpha \gg 0$). These magnitudes govern the patterns of interregional trade, shown in Panel (d) as a function of the number of potential patients N . Imports from other regions rise with N . With increasing returns to scale ($\alpha > 0$), exports to other regions also rise with N (a weak home-market effect). When the scale elasticity α is larger ($\alpha \gg 0$), the import curve is flatter and the export curve is steeper. With sufficiently strong increasing returns, an increase in local demand causes a greater increase in exports than imports (a strong home-market effect).

Figure 5: Estimated quality is positively correlated with total output and external quality metrics

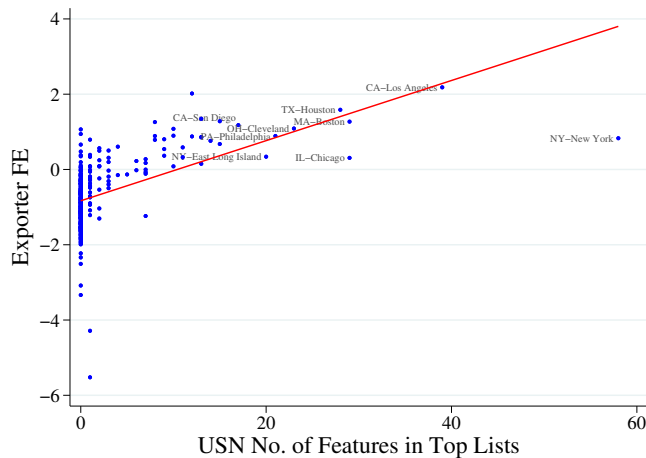
(a) Quality estimates vs. mortality rate (Chandra et al.)



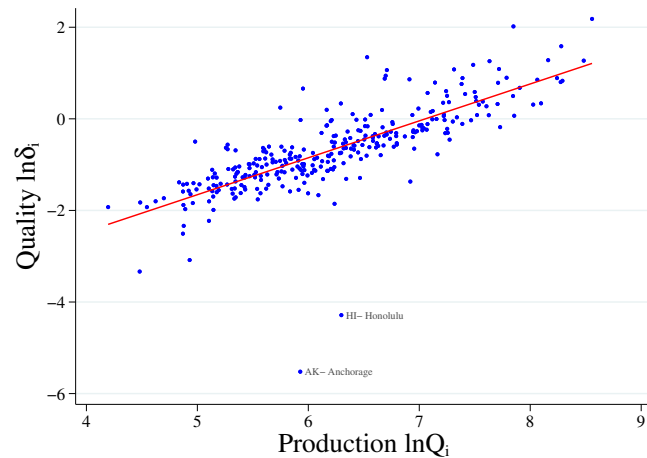
(b) Quality estimates vs. mortality rate (CMS)



(c) *U.S. News* vs. estimated quality

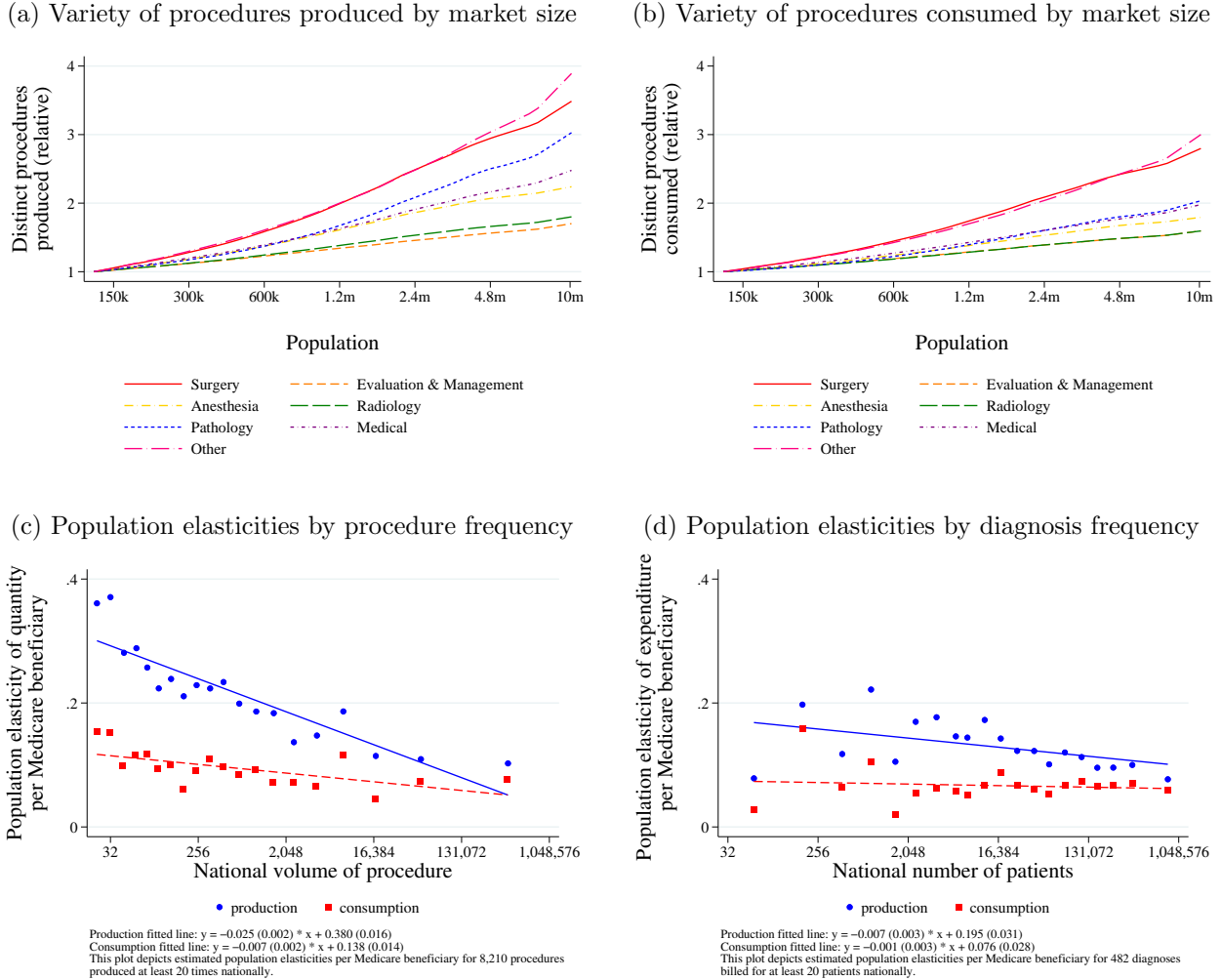


(d) Quality is higher in regions producing more output



Notes: The first three panels show the relationship between the exporter fixed effects (our revealed-preference measure of quality) and external quality measures. The vertical axis shows the exporter fixed effects for each HRR estimated from equation (8). The horizontal axis in Panel (a) is an average of the hospital mortality rates estimated by Chandra, Dalton, and Staiger (2023) in a region, and in Panel (b) an average of those estimated by CMS. The negative correlations indicate patients travel farther to obtain care from regions with lower hospital mortality rates. The horizontal axis in Panel (c) is a count of the number of times each region’s hospitals appear on the *U.S. News* list of best hospitals. *U.S. News* produces an overall ranking as well as rankings for 12 particular specialties. We count the number of times each HRR’s hospitals appear on any of these 13 lists. The relationship is positive, indicating that patients travel farther to obtain care from regions highly ranked by *U.S. News*. The horizontal axis in Panel (d) is the volume of production. Section 1.1 and Appendices A.1 and A.2 detail how we compute production and trade of physician services (excluding emergency-room care and skilled nursing facilities) at standardized prices from the Medicare 20% carrier, 100% MedPAR, and 100% outpatient claims Research Identifiable Files.

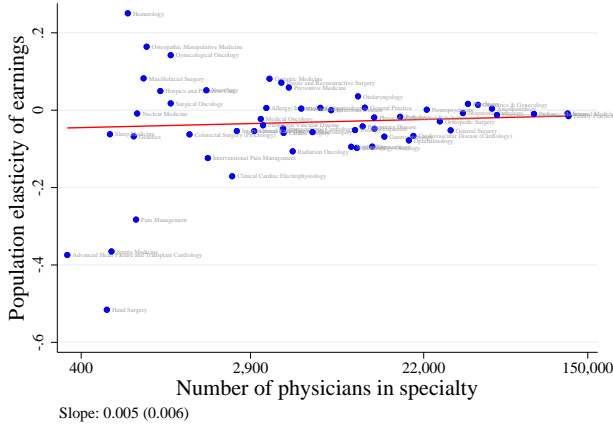
Figure 6: Production and consumption of medical care by market size and service



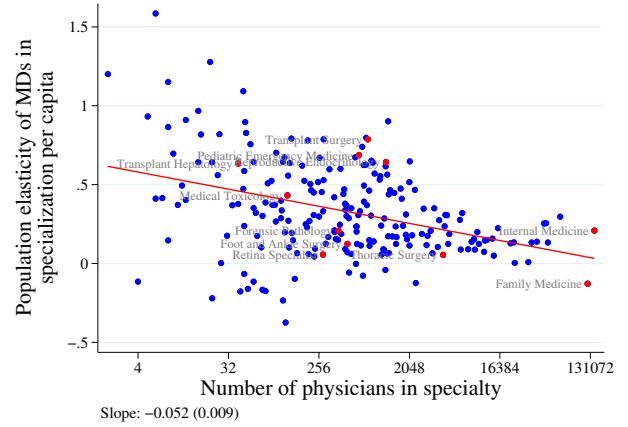
Notes: Panel (a) depicts the number of distinct services produced in an HRR as a function of population. We use procedure classifications from the American Academy of Professional Coders, which groups codes into surgeries, anesthesia, radiology, pathology, medical, and evaluation & management services (AAPC, 2021). We combine Category II codes, Category III codes and Multianalyte Assays into “other.” Within each group, the procedure count is relative to that of the smallest regions. More populous HRRs produce a greater number of services. Panel (b) depicts the number of distinct services consumed in an HRR as a function of population. In Panels (c) and (d), the vertical axes are the population elasticities of quantity of medical care produced (blue dots) and consumed (red squares) per local Medicare beneficiary. The elasticities are computed using the Poisson model in equation (13) based on place of service and patients’ residential location, respectively. Panel (c) estimates these elasticities for each of the procedures provided at least 20 times nationally in the Medicare data. Panel (d) estimates the elasticities for care provided to treat each of the Clinical Classifications Software Refined (CCSR) diagnoses billed for at least 20 patients nationally in the Medicare data. Each panel depicts these population elasticities as a function of the national volume (of procedures and diagnoses, respectively). Section 1.1 and Appendices A.1 and A.2 detail how we compute production and consumption of physician services (excluding emergency-room care and skilled nursing facilities) at standardized prices from the Medicare 20% carrier, 100% MedPAR, and 100% outpatient claims Research Identifiable Files. The contrasting population elasticities of production and consumption imply trade in medical services between markets of different sizes, with more net trade for rare procedures and rare diseases.

Figure 7: Larger regions produce more specialized care, with more specialized equipment

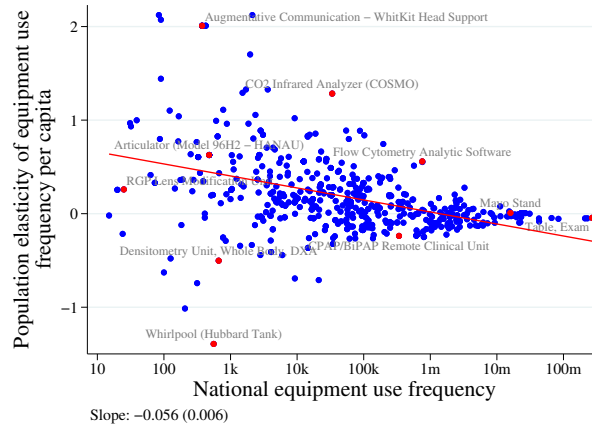
(a) Population elasticities of specialties' earnings



(b) Population elasticities of physician specializations

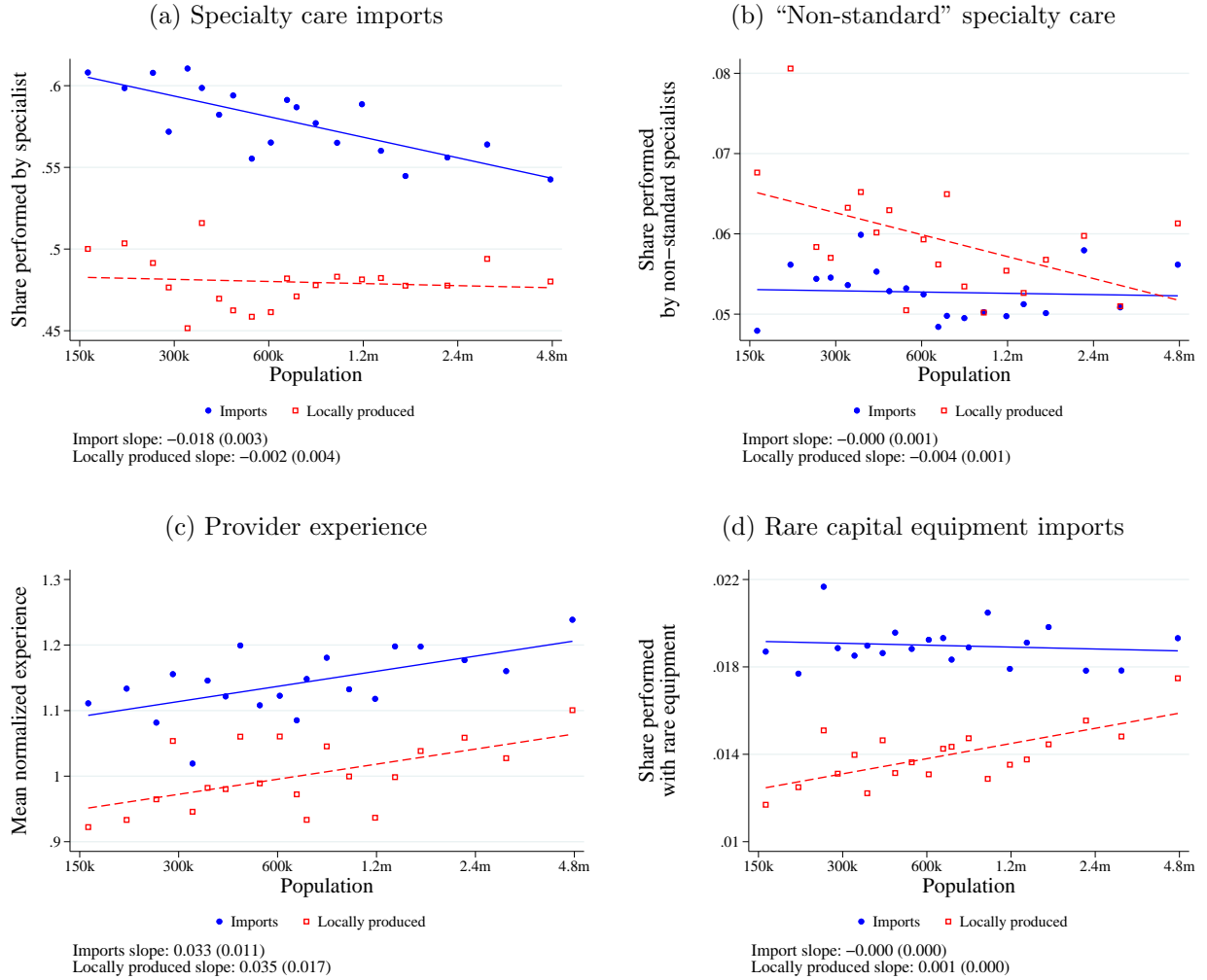


(c) Population elasticities of equipment use



Notes: Panel (a) shows the population elasticity of income for different medical specialties against the total number of physicians in those specialties. For each specialty, we estimate the elasticity of income with respect to population across commuting zones, using data from Gottlieb et al. (2023). The regression line weights each specialty with the number of physicians in that specialty. The graph shows that these elasticities are unrelated to the total national count of physicians in those specialties. The vertical axis of Panel (b) depicts the population elasticities of quantity of physicians in an HRR. The population elasticities are computed for each specialty using the Poisson model in equation (14). The horizontal axis shows the nationwide number of physicians in each specialty. The negative relationship indicates that rare specialties are disproportionately concentrated in high-population regions. Panel (c) plots the population elasticity of equipment use frequency per capita against the national frequency of equipment use, for each piece of equipment. Pieces of equipment with elasticities above the 98th percentile are excluded. The graph shows that rarely-used equipment is used more intensively in large markets.

Figure 8: Imports are specialist-intensive, especially in smaller regions

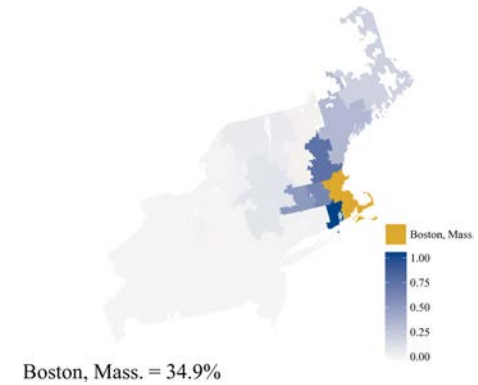
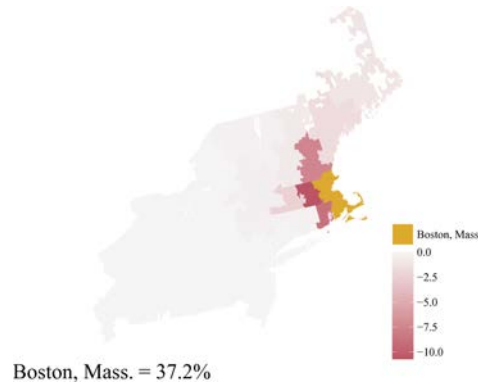


Notes: Panel (a) shows the share of procedures that are performed by a specialist, for imports and locally produced procedures, by market size. We define generalists as internal-medicine, general-practice, and family-practice physicians and define specialists as all other physicians. Imports are more likely to be performed by a specialist, and smaller markets' imports especially so. Panel (b) examines procedures that are typically performed by specialists, and classifies the "standard" specialists as the top two specialties performing the procedure nationally. It shows the shares of procedures performed by the "non-standard" specialties in imported specialty care and locally produced specialty care as a function of local population size. Imports are less likely to be performed by "non-standard" specialties, especially for smaller regions. Panel (c) shows the mean relative experience of providers for care produced locally and imported by population size of the patient's region. This panel describes only procedures that are performed in all hospital referral regions (143 procedures). In public-use Medicare data, we define a provider's experience for a given procedure as the number of times they performed the procedure for Traditional Medicare patients in the prior calendar year. Before aggregating to the regional level, we rescale experience in each procedure so that its mean is one. We further normalize at the regional level so that the experience of the average HRR is one. On average, patients in larger markets obtain treatment from more experienced providers. At all population levels, imported care is produced by more experienced providers than local care. Panel (d) shows the share of procedures that use rare capital equipment in imported and locally produced care by market size. Rare equipment is defined as those pieces of equipment with below-median use in the Medicare data, defined as in Appendix A.1. We see that locally produced care in larger markets is more likely to use rare equipment, and imported care has a higher rare-equipment share throughout the population distribution

Figure 9: Counterfactual outcomes for higher reimbursements in one region

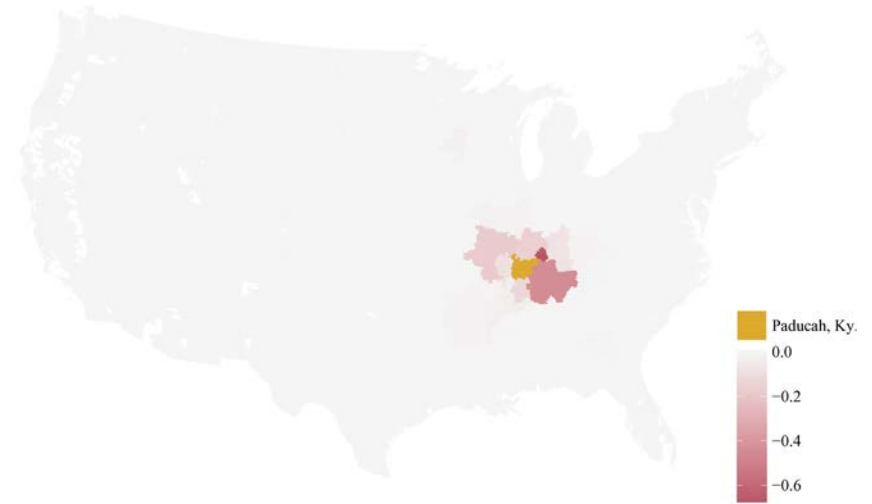
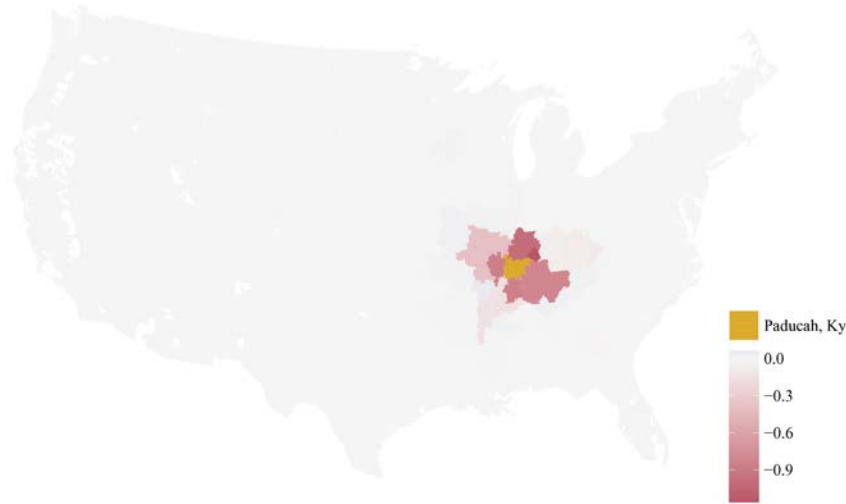
(a) Change (%) in output quality δ_i : higher reimbursement in Boston, Mass.

(b) Change (%) in market access Φ_i : higher reimbursement in Boston, Mass.



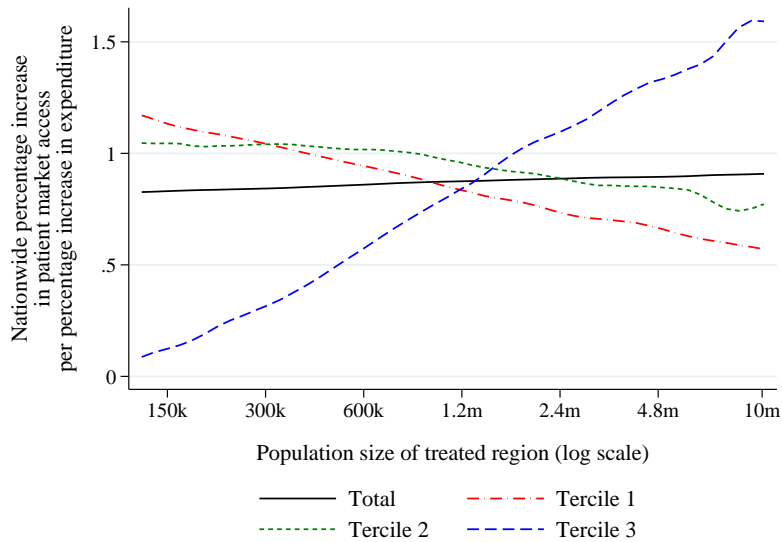
(c) Change (%) in output quality δ_i : higher reimbursement in Paducah, Ky.

(d) Change (%) in market access Φ_i : higher reimbursement in Paducah, Ky.



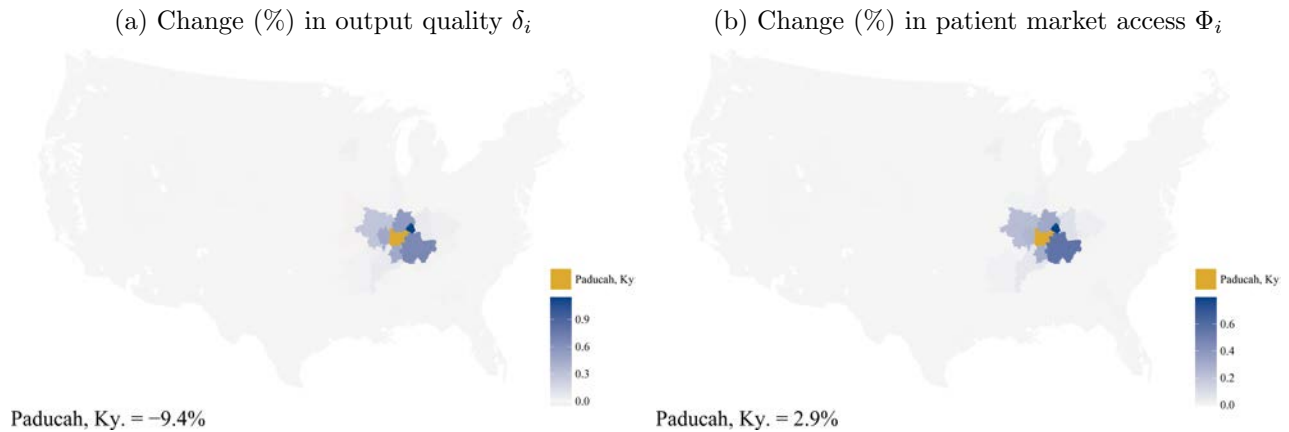
Notes: Panels (a) and (b) show the impacts of increasing reimbursements by 30% in the Boston, Mass. HRR ($\hat{R}_i = 1.3$) based on our estimated model. Panel (a) illustrates the percentage change in quality of care δ_i provided in each region. Panel (b) illustrates the percentage change in the value of market access Φ_i for patients who live in an region. Panels (c) and (d) are analogous, but for a 30% increase in reimbursements in Paducah, Ky., a net importer. In all panels, the predicted change for the region whose reimbursement changes (“treated region”) is listed on the map itself. In both cases, the quality produced in neighboring regions declines (Panels (a) and (c)). Patients in regions near Boston benefit from increased access to the treated region (Panel (b)), so there is a negative relationship between the percentage changes in δ and Φ across regions. In contrast, patients in regions near Paducah suffer a decrease in access (Panel (d)). The contrasting outcomes stem from Boston being a net exporter and Paducah being a net importer in the baseline equilibrium. The exercise is described in detail in Section 5.

Figure 10: Returns to higher reimbursements in one region by region size



Notes: This figure summarizes the counterfactual outcomes of 30% higher reimbursements in one HRR as a function that HRR’s population size. The nationwide return is the percentage increase in patient market access $\sum_{\kappa} \sum_j N_{j\kappa} \Phi_{j\kappa}$ per percentage increase in nationwide expenditures $\sum_i Q_i R_i$. The tercile-specific return is the increase in tercile-specific patient market access $\sum_j N_{j\kappa} \Phi_{j\kappa}$. Increasing reimbursements in more populous HRRs has the highest return when measured as impact on aggregate market access. Subsidies in less populous regions favor lower-income patients, primarily because there are more low-income patients living in and close to smaller regions.

Figure 11: Counterfactual outcomes when changing travel costs for Paducah, Ky. residents



Notes: This figure shows the impacts of a counterfactual 30% decrease in travel costs for Paducah residents ($\hat{\rho}_{ij} = 1.3 \forall i \neq \text{Paducah}$). Panel (a) depicts the percentage change in quality of care δ_i produced in each region. Panel (b) depicts the percentage change in the value of market access Φ_i for patients who live in a region. The notes report the changes for Paducah.

Table 1: Scale elasticity estimates

All services	Baseline	No Diagonal	Controls
OLS: 2017	0.806 (0.031)	0.961 (0.047)	0.786 (0.041)
OLS: 2013–2017 difference	0.999 (0.079)	1.045 (0.083)	1.018 (0.082)
2SLS: population (log)	0.800 (0.037) [2141]	0.905 (0.057) [2141]	0.777 (0.050) [1621]
2SLS: population (1940, log)	0.697 (0.063) [163]	0.924 (0.093) [163]	0.633 (0.070) [206]

Notes: This table reports estimates of α from ordinary least squares (OLS) or two-stage least squares (2SLS) regressions of the form $\widehat{\ln \delta}_i = \alpha \ln Q_i + \ln R_i + \ln w_i + u_i$ using Hospital Referral Regions (HRRs) as the geographic units. The dependent variable $\widehat{\ln \delta}_i$ is estimated in equation (8) using a same-region dummy and a quadratic function of log distance. Q_i is region i 's total production for Medicare beneficiaries, R_i is Medicare's Geographic Adjustment Factor, the w_i covariate includes mean two-bedroom property value and mean annual earnings for non-healthcare workers, and u_i is an error term. In the "no diagonal" column, S_{ii} observations were omitted when estimating $\widehat{\ln \delta}_i$ in equation (8). In the third column, the $\ln R_i$ and $\ln w_i$ controls are included in the regressions (coefficients not reported). The first-difference specification estimates $\Delta \widehat{\ln \delta}_i = \alpha \Delta \ln Q_i + \Delta u_i$. In the rows labeled "2SLS," we instrument for $\ln Q_i$ using the specified instruments. The standard errors in parentheses are robust to heteroskedasticity. For 2SLS estimates, first-stage effective F -statistics (Montiel Olea and Pflueger, 2013) are reported in square brackets. All estimates reveal substantial scale economies.

Table 2: Aggregate medical services exhibit a strong home-market effect

	(1)	(2)	(3)	(4)	(5)
	Cross-sectional PPML			IV: 1940 population	2013–2017 panel
λ_X Provider-market population (log)	0.671 (0.0543)	0.681 (0.0505)	0.671 (0.0366)	0.757 (0.0547)	0.939 (0.151)
λ_M Patient-market population (log)	0.260 (0.0547)	0.252 (0.0501)	0.286 (0.0346)	0.284 (0.0467)	-0.205 (0.148)
Distance (log)	-1.627 (0.0489)	0.344 (0.304)		0.377 (0.250)	
Distance (log, squared)		-0.199 (0.0305)		-0.201 (0.0247)	
Distance (log) \times 2017					-0.00117 (0.00667)
p-value for $H_0: \lambda_X \leq \lambda_M$	<0.001	<0.001	<0.001	<0.001	<0.001
Observations	93,636	93,636	93,636	93,636	162,678
Fixed effects					ij
Distance elasticity at mean		-1.59		-1.57	
Distance deciles			Yes		

Notes: This table reports estimates of equations (10) and (11), which evaluate the presence of weak or strong home-market effects. The sample is all HRR pairs ($N = 306^2$), multiplied by two for the last column. The dependent variable in all regressions is the value of trade when including professional and facility (inpatient and outpatient) fees at national average prices. Columns 1 through 4 use 2017 data. In column 1, the independent variables are patient- and provider-market log population, log distance between HRRs, and an indicator for same-HRR observations ($i = j$). The positive coefficient on provider-market log population implies a weak home-market effect, and the fact that this coefficient exceeds that on patient-market population implies a strong home-market effect. Column 2 makes the distance coefficient more flexible by adding a control for the square of log distance. Column 3 replaces parametric distance specifications with fixed effects for each decile of the distance distribution. Column 4 uses the provider-market and patient-market log populations in 1940 as instruments for the contemporaneous log populations when estimating by generalized method of moments. Column 5 presents estimates of equation (11) including ij fixed effects using data from 2013 and 2017. Section 1.1 and Appendices A.1 and A.2 detail how we compute trade flows in physician services (excluding emergency-room care and skilled nursing facilities) at standardized prices from the Medicare 20% carrier, 100% MedPAR, and 100% outpatient claims Research Identifiable Files. HRR definitions are from the Dartmouth Atlas Project. Standard errors (in parentheses) are two-way clustered by patient market and provider market.

Table 3: The home-market effect is stronger for rare procedures and diagnoses

	Procedure		Procedure		Procedure		Diagnosis	
	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)
λ_X Provider-market population (log)	0.618 (0.0516)	0.605 (0.0493)	0.603 (0.0493)		0.606 (0.0488)		0.601 (0.0489)	
λ_M Patient-market population (log)	0.360 (0.0519)	0.364 (0.0492)	0.366 (0.0492)		0.364 (0.0486)		0.371 (0.0486)	
μ_X Provider-market population (log) \times rare			0.344 (0.0447)	0.329 (0.0405)	0.362 (0.0452)	0.317 (0.0392)	0.120 (0.0232)	0.110 (0.0206)
μ_M Patient-market population (log) \times rare			-0.241 (0.0606)	-0.239 (0.0587)	-0.250 (0.0612)	-0.220 (0.0564)	-0.0986 (0.0186)	-0.0915 (0.0172)
p-value for $H_0: \lambda_X \leq \lambda_M$	0.005	0.006	0.007		0.005		0.008	
p-value for $H_0: \mu_X \leq \mu_M$			<0.001	<0.001	<0.001	<0.001	<0.001	<0.001
Observations	187,272	110,402	110,402	110,402	110,402	110,402	109,658	109,658
Distance [linear] controls	Yes	Yes	Yes	Yes				
Distance [quadratic] controls					Yes	Yes	Yes	Yes
Patient-provider-market-pair FEs				Yes		Yes		Yes

Notes: This table reports estimates of equation (12), which introduces interactions with an indicator for whether a procedure or diagnosis is “rare” (provided less often than the median, when adding up all procedures provided nationally). The interactions with patient- and provider-market population reveal whether the home-market effect is larger for rare procedures/diagnoses. The unit of observation is {rare indicator, exporting HRR, importing HRR} so the number of observations is 2×306^2 in column 1, and the dependent variable in all regressions is the value of trade. Columns 1–6 reflect estimates by procedure and Columns 7 and 8 reflect estimates by diagnosis. Columns 2 onwards drop HRR pairs with zero trade in both procedure groups, and column 2 shows that this restriction has a negligible impact on the estimated log population coefficients. Columns 3 onwards include the rare indicator interacted with patient- and provider-market populations and distance covariates. Columns 1–4 control for distance using the log of distance between HRRs. Columns 5–8 add a control for the square of log distance. Columns 4, 6, and 8 introduce a fixed effect for each ij pair of patient market and provider market, so these omit all covariates that are not interacted with the rare indicator. The positive coefficient on provider-market population \times rare across all columns indicates that the home-market effect is stronger for rare than for common services. The negative coefficient on patient-market population \times rare across all columns indicates that the *strong* home-market effect has a larger magnitude for rare services. Valid primary diagnoses observed in 1,000 distinct claims or more nationally in the professional fees 20% sample are included. Section 1.1 and Appendices A.1 and A.2 detail how we compute these trade flows for physician services (excluding emergency-room care and skilled nursing facilities) at standardized prices from the Medicare 20% carrier, 100% MedPAR, and 100% outpatient claims Research Identifiable Files. HRR definitions are from the Dartmouth Atlas Project. Standard errors (in parentheses) are two-way clustered by patient market and provider market.