

The economics of social data

Dirk Bergemann*

Alessandro Bonatti**

and

Tan Gan*

A data intermediary acquires signals from individual consumers regarding their preferences. The intermediary resells the information in a product market wherein firms and consumers tailor their choices to the demand data. The social dimension of the individual data—whereby a consumer's data are predictive of others' behavior—generates a data externality that can reduce the intermediary's cost of acquiring the information. The intermediary optimally preserves the privacy of consumers' identities if and only if doing so increases social surplus. This policy enables the intermediary to capture the total value of the information as the number of consumers becomes large.

1. Introduction

Individual and social data. The rise of large digital platforms—such as Facebook, Google, and Amazon in the United States and JD, Tencent, and Alibaba in China—has led to the unprecedented collection and commercial use of individual data. The steadily increasing user bases of these platforms generate massive amounts of data about individual consumers, including their preferences, locations, friends, and political views. In turn, many of the services provided by large Internet platforms rely critically on these data. The availability of individual-level data permits refined search results, personalized product recommendations, informative ratings, timely traffic data, and targeted advertisements.

The recent disclosures on the use and misuse of social data by digital platforms have prompted regulators to limit the largely unsupervised use of individual data by these companies. As a result, nearly all proposed and enacted regulation to date aims to ensure that consumers

* Yale University; dirk.bergemann@yale.edu, tan.gan@yale.edu.

** MIT Sloan School of Management; bonatti@mit.edu.

Bergemann and Bonatti acknowledge financial support through NSF Grant SES-1948692. Bergemann and Gan acknowledge financial support from the Omidyar and Knight Foundation. We thank Joseph Abadi, Daron Acemoglu, Susan Athey, Steve Berry, Nima Haghpahan, Nicole Immorlica, Al Klevorick, Scott Kominers, Annie Liang, Roger McNamee, Jeanine Miklós-Thal, Enrico Moretti, Stephen Morris, Denis Nekipelov, Asu Özdağlar, Fiona Scott-Morton, Shoshana Vasserman, Glen Weyl, and Kai-Hao Yang for helpful discussions. We also thank Michelle Fang and Miho Hong for valuable research assistance and the audiences at numerous seminars and conferences for their productive comments.

retain control over their data. However, the *digital privacy paradox* (e.g., Athey, Catalini, and Tucker (2017)) indicates that even small monetary incentives may lead individuals to relinquish their private data. The low cost of acquiring private data—seemingly at odds with consumers’ stated preferences over their privacy—drives the appetite of platforms to gather information and may undermine the efficacy of regulation.¹

This article suggests a unified explanation for the digital privacy paradox and the selective use of data for price and product choices. A key observation is that *individual data* are actually *social data*: Data captured from an individual user are informative not only about that user but also about other users with similar characteristics or behaviors. The social dimension of the data generates a *data externality*, the sign and magnitude of which depend on the structure of the data and on the use of the gained information.

Many digital platforms use social data to increase the value provided by their services. Google search results, for example, are informed by the choices of previous users. Indeed, the search engine is only successful because it mediates information among many consumers. Likewise, Amazon uses data collected from other consumers to curate a “recommended for you” list of products and explicitly suggests goods that are “frequently bought together.” Finally, YouTube and Waze tailor their suggestions of videos and traffic directions to each user’s preferences, as estimated by combining network data and individual data. However, social data also facilitate surplus extraction, for example, individual shopping data convey information about the willingness to pay of consumers with similar purchase histories.

We analyze three critical aspects of the economics of social data. First, we consider how the collection and transmission of individual data change the terms of trade among consumers, firms (e.g., advertisers), and data intermediaries (e.g., large Internet platforms that sell targeted advertising space). Second, we examine how the social dimension of the data magnifies the value of individual data for platforms and facilitates the acquisition of large datasets. Third, we analyze how data intermediaries with market power manipulate the trade-offs induced by social data through the aggregation and the precision of the information that they provide about consumers.

A model of data intermediation. We develop a framework to evaluate the flow and allocation of individual data in the presence of data externalities. Our model focuses on three types of economic agents: consumers, firms, and data intermediaries. These agents interact in two distinct but linked markets: a *data market* and a *product market*.

In the product market, each consumer (she) determines the quantity that she wishes to purchase, and a single producer (he) sets the unit price at which he offers a product to the consumers. Initially, each consumer has private information about her willingness to pay for the firm’s product. This information consists of a signal with two additive components: a *fundamental* component and a *noise* component. The fundamental component represents her willingness to pay, and the noise component reflects that her initial information might be imperfect. Both components can be correlated across consumers: In practice, different consumers’ preferences can exhibit common traits, and consumers might undergo similar experiences that induce correlation in their errors.

In the data market, a monopolist intermediary acquires demand information from the individual consumers in exchange for a monetary payment. The intermediary then chooses how much information to share with the other consumers and how much information to sell to the producer. Sharing data with consumers allows them to tailor their demand to their true preferences. Sharing data with the producer enables more tailored and possibly personalized prices.

We introduce a rich model for the structure of individual data, which allows for correlation in the fundamentals as well as in the noise terms across the individuals. We view this richness in the data structure as the defining element of data in digital platforms. In contrast, we adopt a more

¹ The recent report by Furman et al. (2019) identifies “the central importance of data as a driver of concentration and a barrier to competition in digital markets”—a theme echoed in the reports by Cr  m  r, de Montjoye, and Schweitzer (2019) and by the Stigler Committee on Digital Platforms (2019).

specific model for product market interaction whereby a monopolist seller charges linear prices for variable quantity. However, our main insights apply to any product market where: (a) data sharing teaches consumers about their preferences, and (b) a firm seeks to extract the consumers' surplus.²

The value of social data. Collecting data from multiple consumers helps any market participant to filter the individual signals. This process can occur through two channels. First, in a market where the noise terms are largely idiosyncratic, a large sample size filters out errors and identifies any common fundamentals. Second, in a market with largely idiosyncratic fundamentals, many observations filter out demand shocks and identify common noise terms, thereby estimating individual fundamentals by differencing (Proposition 1).

However, the choice by each consumer to share her data with the intermediary is guided only by her private benefits and costs, not by the information gains she generates with her actions. Thus, the intermediary must compensate each individual consumer only to the extent that the disclosed information affects her own welfare *on the margin*. Critically, the platform does not have to compensate the individual consumer for any changes in her welfare caused by the information deduced from other consumers' signals.

Therefore, social data drive a wedge between the socially efficient and profitable uses of information. First, the cost of acquiring individual data can be substantially less than the value of the information to the platform. Second, although many uses of consumer information exhibit positive externalities, very little prevents the platform from trading data for profitable uses that are in fact harmful to consumers (Proposition 2).³

More generally, the data externality can induce too much or too little trade in data. Indeed, the condition for data intermediation to yield positive profits is qualitatively different from the condition for intermediation to yield social welfare gains. In particular, the intermediary obtains positive profits when a large number of consumers exhibit a strong degree of correlation in their preferences: This allows the intermediary to acquire the consumers' data in exchange for minimal compensation. On the other hand, welfare improvements depend on how much additional information each consumer can obtain about her own preferences from the other consumers' signals. Therefore, when many consumers have strongly correlated preferences and very precise signals, data sharing is detrimental to welfare but profitable to the intermediary. Conversely, there are data structures (e.g., ones with independent fundamentals and strongly correlated error terms) for which data sharing is beneficial to consumers but unprofitable for the data intermediary.

Equilibrium data-sharing policies. A natural question is then whether the data market imposes any limitations at all on equilibrium information sharing. To shed light on this issue, we consider the choice of whether to reveal the consumers' identities to the producer or to collect anonymous data. When consumers are homogeneous *ex ante*, we show that the intermediary collects anonymous data if and only if the transmission of identity data reduces total surplus. Therefore, even if the data *transmission* may be socially detrimental, the optimal choice of data *anonymization* is socially efficient (Proposition 3).

In our model of linear price discrimination, collecting anonymous data amounts to selling aggregate, market-level information to the producer. With this choice, the intermediary does not

² In fact, the value of social data in Section 3 can be computed for alternative specification of the product market and a general result for anonymization is obtained in Section 4. Finally in Section 5, we consider a richer environment where a merchant offers different varieties to consumers with heterogeneous tastes for his products. In that case, the firm's actions both generate value and attempt to capture it.

³ Recent empirical work on the effects of privacy regulation such as the European Union's General Data Protection Regulation (e.g., Aridor, Che, and Salz (2020) and Johnson, Shriver, and Goldberg (2020)), indicates that data externalities are relevant for consumers' and businesses' decisions to share their data. In the United States, legislators are also increasingly aware of the consequences of data externalities. In particular, the US House Committee on the Judiciary (2020) reports that "[...] the social data gathered through [a platform's] services may exceed their economic value to consumers."

enable the producer to set personalized prices: the data are transmitted but disconnected from the users' personal profiles. In other words, the role of social data provides a more nuanced ability to determine the modality of information acquisition and use.

Under anonymized data intermediation, the gap between the social value of the data and the price of the data widens when the number of consumers increases. In particular, as the sources of data multiply, the contribution of each individual consumer to the aggregate information shrinks, which drives down the individual payments to consumers, and possibly the total payment as well (Proposition 5).

We develop a general anonymization result (Proposition 8) and extend the model in several directions. We first introduce consumer heterogeneity by considering multiple groups of consumers. We find that the intermediary aggregates the data at least to the level of the coarsest partition of homogeneous consumers, although further aggregation is profitable when the number of consumers is small. The resulting group pricing—discriminatory pricing based on observable characteristics, such as location—has welfare consequences between those of complete privacy and those of price personalization (Proposition 9).

We then consider a model in which the producer can choose prices and product characteristics to match an additional horizontal (taste) dimension of the consumers' preferences. The resulting data policy then aggregates the vertical dimension but not the horizontal dimension, thereby enabling the producer to offer personalized product recommendations but not personalized prices (Proposition 10). Finally, in the online Appendix, we explore the data intermediary's ability to offer privacy guarantees by collecting less than perfect information about the consumers' signals.

Implications and applications. There exist a wide variety of data intermediaries and associated business models. The point of our article is not to try to match any specific business model but rather to speak to the general principles that seem to be in effect in many markets. In particular, our model assumes that each consumer is compensated directly with a monetary transfer for her individual data. There exist a few concrete examples of such transactions (e.g., Nielsen offers monetary rewards to consumers for access to their browsing and purchasing data). However, most data intermediaries compensate their users via the quality of the services they offer, for example, social networks, search, mail, video. These services are nominally free, but in practice (through more or less transparent terms and conditions), they are fueled by the data generated by their users.

Likewise, digital platforms occasionally transfer consumers' data to merchants for a fee. Much more often, however, they monetize their data by selling access to targeted advertising campaigns—see the report by the Competition & Markets Authority (2020). This enables the merchants to reap the value of information, by conditioning their messages and their prices on the consumers' preferences, without directly observing their data. These transactions amount to indirect sales of information, as discussed in Bergemann and Bonatti (2019).

Our model's predictions for the nature of the externality, for the equilibrium allocation of data, and for its welfare properties then depend on how targeted advertising affects total surplus in any given market. Specifically, if advertising generates value by matching of buyers and sellers, as in Bergemann and Bonatti (2011), our model would predict the intermediation of individual-level information through very detailed targeting categories. If, instead, targeted advertisements facilitate inefficient price discrimination, the model would predict coarser targeting categories that induce market-level or group-level pricing.

Related literature. This article contributes to the growing literature on data markets recently surveyed in Bergemann and Bonatti (2019). In particular, the role of data externalities in the socially excessive diffusion of personal data has been a central concern in Choi, Jeon, and Kim (2019) and Acemoglu et al. (2021), both considering a model with many buyers and one firm that acts as an integrated data intermediary and seller.

Choi, Jeon, and Kim (2019) introduce information externalities into a model of monopoly pricing with unit demand. Each consumer is described by two independent random variables: her willingness to pay for the monopolist's service and her sensitivity to a loss of privacy. The purchase of the service by the consumer requires the transmission of personal data. From the collected data, the seller gains additional revenue, depending on the proportion of units sold and the volume of data collected. The total nuisance cost paid by each consumer depends on the total number of consumers sharing their personal data. Thus, the optimal pricing policy of the monopolist yields excessive loss of privacy, relative to the social welfare maximizing policy.

Acemoglu et al. (2021) also analyze data acquisition in the presence of information externalities. As in Choi, Jeon, and Kim (2019), they consider a model with many consumers and a single data-acquiring firm. Like the current analysis, Acemoglu et al. (2021) propose an explicit statistical model for their data; the model allows the authors to assess the loss of privacy for the consumer and the gains in prediction accuracy for the firm. Their analysis then pursues a different, and largely complementary, direction from ours. In particular, they analyze how consumers with heterogeneous privacy concerns trade information with a data platform and derive conditions under which the equilibrium allocation of information is (in)efficient.

In contrast to these two important contributions, we explicitly introduce a data intermediary with objectives distinct from either the consumers or the seller. We consider rich data structures for fundamentals and noise terms that capture the wide range of social data on digital platforms. This allows us to endogenize privacy concerns and to quantify the downstream welfare impact of data intermediation. In addition, we can investigate when and how privacy can be partially or fully preserved through aggregation, anonymization, and noise. Thus, we augment the analysis in the above contributions with additional insights regarding data flows and data intermediation. In particular, we show that the data externality can have a positive or a negative impact on consumer and social surplus depending on the data structure. In consequence, a monopolist intermediary can induce either too little or too much information sharing in equilibrium.

Acquisti, Taylor, and Wagman (2016) survey the literature on the economics of privacy in great detail. An early and influential article is Taylor (2004), who analyzes the sales of consumer purchase histories without data externalities. More recently, Cummings et al. (2016) investigate how privacy policies affect user and advertiser behavior in a simple model of targeted advertising. The low level of compensation that users command for their personal data is discussed in Arrieta-Ibarra et al. (2018), who propose sources of countervailing market power, and in a growing body of empirical work. In particular, Tang (2019) uses large-scale field experiments to estimate the value that online borrowers assign to several pieces of personal data. Lin (2019) separates intrinsic and instrumental preferences for privacy through a lab experiment that also uncovers data externalities.

Fainmesser, Galeotti, and Momot (2020) provide a digital privacy model in which data collection improves the service provided to consumers. However, as the collected data can also leak to third parties and thus impose privacy costs, an optimal digital privacy policy must be established. Similarly, Jullien, Lefouili, and Riordan (2020) analyze the equilibrium privacy policy of websites that monetize information collected from users by charging third parties for targeted access. Gradwohl (2017) considers a network game in which the level of beneficial information sharing among the players is limited by the possibility of leakage and a decrease in informational interdependence. Ali, Lewis, and Vasserman (2019) study a model of personalized pricing with disclosure by an informed consumer, and they analyze how different disclosure policies affect consumer surplus. Ichihashi (2020b) studies both personalized pricing and product recommendations and shows that a seller benefits from committing not to use the consumer's information to set prices. Our result on optimal anonymization and market-level pricing has similar implications, but is entirely driven by the data externality that appears when multiple consumers are present.

Finally, Liang and Madsen (2020) investigate how data policies can provide incentives in principal-agent relationships. They emphasize the structure of individual data and how the substitute or complement nature of individual signals determines the impact of data on incentives.

Ichihashi (2020a) considers a single data intermediary and asks how complement or substitute consumer signals affect the equilibrium price of the individual data.

2. Model

■ We consider an idealized trading environment with many consumers, a single intermediary in the data market, and a single producer in the product market.

□ **Product market.** There are finitely many consumers, labeled $i = 1, \dots, N$. In the product market, each consumer (she) chooses a quantity level q_i to maximize her net utility given a unit price p_i offered by the producer (he):

$$u_i(w_i, q_i, p_i) \triangleq w_i q_i - p_i q_i - \frac{1}{2} q_i^2.$$

Each consumer i has a baseline willingness to pay for the product $w_i \in \mathbb{R}$.

The producer sets the unit price p_i at which he offers his product to each consumer i . The producer has a linear production cost

$$c(q) \triangleq c \cdot q, \text{ for some } c \geq 0.$$

The producer's profits are given by

$$\pi(p_i, q_i) \triangleq \sum_i (p_i - c) q_i.$$

□ **Data environment.** The vector of willingness to pay, $w = (\dots, w_i, \dots) \in \mathbb{R}^N$, is distributed according to a joint distribution F_w :

$$w \sim F_w. \quad (1)$$

Initially, each consumer may have only imperfect information about her willingness to pay. In particular, consumer i observes a signal

$$s_i \triangleq w_i + \sigma \cdot e_i, \quad (2)$$

where $\sigma > 0$ and e_i is consumer i 's error term. The error terms $e = (\dots, e_i, \dots) \in \mathbb{R}^N$ are independent of the willingness to pay w , and they are distributed according to a joint distribution F_e :

$$e \sim F_e. \quad (3)$$

We denote by S the information structure generated by the complete vector of consumer signals $s = (\dots, s_i, \dots) \in \mathbb{R}^N$. We allow for arbitrary distributions of fundamentals w and errors e , and hence arbitrary correlation structures across consumers, under the restriction that the (F_w, F_e) are symmetric across individuals. We view the richness in the data structure as represented by (1) and (3) as the defining feature of social data in the digital economy. In particular, the noise in the signal s_i of each individual given by (2) reflects the importance of social learning as enabled by recommender, rating and search engines.

Without loss of generality, we assume that (i) each individual willingness to pay w_i has mean μ and variance 1; (ii) individual errors e_i have mean 0 and variance 1 (which is scaled by the parameter σ).

The producer knows the structure of demand and thus the common prior distribution of consumers' willingness to pay. However, absent any additional information, the producer does not know the realized willingness to pay w_i of any consumer (nor her signal s_i) prior to setting prices. This data environment has two important features. First, any demand information beyond the common prior comes from the signals of the individual consumers. Second, with any amount

of noise in the signals (i.e., if $\sigma > 0$), each consumer can learn more about her own demand from the signals of the other consumers.

The following leading examples illustrate two ways in which data sharing can help each consumer learn more about her individual willingness to pay. In Example 1, a new product has a common value that consumers are imperfectly informed about.

Example 1 (Common preferences). Fundamentals w_i are perfectly correlated and errors e_i are independent: $s_i = w + \sigma \cdot e_i$.

In this case, data sharing helps to filter out the idiosyncratic error terms: As N becomes large, the average signal across all consumers identifies the common willingness to pay.

In Example 2, individual consumers have independent values for a new therapy but are exposed to a common health shock.

Example 2 (Common experience). Errors e_i are perfectly correlated, and fundamentals w_i are independent: $s_i = w_i + \sigma \cdot e$.

Under this structure, the average signal identifies the common error component e as N becomes large. All market participants can then precisely estimate each w_i from the difference between individual and average signal.

As we shall see, information sharing enables learning in both examples. However, the actions of consumers $-i$ impact the surplus of consumer i quite differently in the two cases, which has implications for the equilibrium price of data. More generally, the data structure will determine how to separate the individual and the aggregate information.

□ **Data market.** The data market is run by a single data intermediary (it). As a monopolist market maker, the data intermediary decides how to collect the available information (s_i) from each consumer and how to share it with the other consumers and the producer. Thus, the data intermediary faces both an information design problem and an information pricing problem.

We consider bilateral contracts between the individual consumers and the intermediary and between the producer and the intermediary. The data intermediary offers these bilateral contracts *ex ante*, that is, before the realization of any demand shocks. Each bilateral contract defines a *data policy* and a *data price*.

The data contract with consumer i specifies a *data inflow* policy X_i and a fee $m_i \in \mathbb{R}$ paid to the consumer. The data inflow policy describes how each signal s_i enters the database of the intermediary. We restrict attention to the following two policies: (i) the *complete (identity-revealing)* data policy $X = S$, where the intermediary collects each consumer's signal s_i ; and (ii) the *anonymized* data policy $X = A$, where the intermediary collects individual signals without identifying information. We model the anonymized data policy as

$$A : S \rightarrow \delta(S),$$

for a random permutation of the consumers' indices $i \rightarrow \delta(i)$. In both cases, as discussed in Section 2 below, there are no further incentive constraints, that is, consumers transmit their information truthfully.

In our product market model, where the consumer's demand is linear in her signal, the anonymized data policy A is equivalent to an *aggregate* data policy that conveys information about the average willingness to pay. Intuitively, the anonymized data policy prevents the producer from matching signals to consumers, that is, from profitably charging personalized prices. In Section B, we enrich the intermediary's strategy space by allowing for data policies that collect partial information about the consumers' signals.

A data contract with the producer specifies a *data outflow* policy Y and a fee $m_0 \in \mathbb{R}$ paid by the producer. The data outflow policy determines how each consumer's collected signal is

transmitted to the producer and to other consumers. In particular, letting X denote the intermediary's *realized* data inflow, a data outflow policy $Y = (Y_0, Y_1, \dots, Y_N)$ describes how the collected data are released to the seller,

$$Y_0 : X \rightarrow \Delta(\mathbb{R}^N),$$

and to each consumer,

$$Y_i : X \rightarrow \Delta(\mathbb{R}^N).$$

Sharing data with other consumers is a critical design element because doing so allows each consumer to adjust her quantity demanded at any price. Therefore, the information received by consumers also impacts the *producer's* willingness to pay for the intermediary's data.

The data intermediary maximizes the net revenue

$$R \triangleq m_0 - \sum_{i=1}^N m_i. \quad (4)$$

□ **Equilibrium and timing.** The game proceeds sequentially. First, the terms of trade on the data market and then the terms of trade on the product market are established. The timing of the game is as follows:

1. The data intermediary offers a data inflow policy (m_i, X_i) to each consumer i . Consumers simultaneously accept or reject the intermediary's offer.
2. The data intermediary offers a data outflow policy (m_0, Y) to the producer, based on the (known) number of consumers who have accepted. The producer accepts or rejects the offer.
3. Consumers' signals s are realized, and the information flows (x, y) are transmitted according to the terms of the data policies.
4. The producer sets a unit price p_i for each consumer i who makes a purchase decision q_i , given her available information about w_i .

We analyze the perfect Bayesian equilibria of the game. Under the timing described above, information is imperfect but symmetric at the contracting stage. Furthermore, when the consumer receives the intermediary's offer, she must anticipate the intermediary's choice of data outflow policy, which determines what data are shared with her, as well as with the producer. We denote by $a_0, a_i \in \{0, 1\}$ the participation decisions by the producer and by consumer i , respectively. A perfect Bayesian equilibrium is then a tuple of inflow and outflow data policies, data and product pricing policies, and participation decisions:

$$\{(X^*, Y^*, m^*); p^*, q^*; a^*\},$$

where

$$a_0^* : X \times Y \times \mathbb{R} \rightarrow \{0, 1\}, \quad a_i^* : X_i \times \mathbb{R} \rightarrow \{0, 1\},$$

such that (i) the producer maximizes his expected profits, (ii) the intermediary maximizes its expected revenue, and (iii) each consumer maximizes her net utility. In our baseline analysis, we focus on the best equilibrium for the data intermediary; in the best equilibrium, every consumer accepts the offer from the data intermediary. We discuss a unique implementation in Section 4.

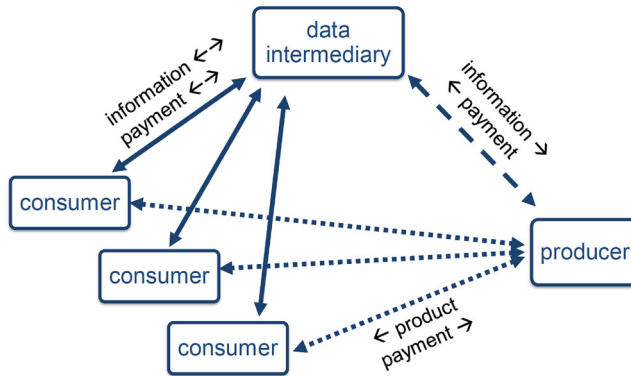
Figure 1 summarizes the information and value flow in the data and product markets.

□ Discussion of model features.

Participation constraints. The participation constraints of every consumer and of the producer are required to hold at the *ex ante* level. Thus, the consumers agree to the data policy before the realization of their signals. The choice of *ex ante* participation constraints captures the prevailing conditions under which users interact with large digital platforms. For instance, users of services

FIGURE 1

DATA AND VALUE FLOWS



on Amazon, Facebook, or Google typically establish an account and accept the “terms of service” before making any specific query or post. Through the lens of our model, the consumer requires a level of compensation that allows her to profitably share the information in expectation. Upon agreeing to participate, there are no further incentive compatibility constraints on the transmission of her information.

Lack of commitment. As we mentioned above, targeted advertising is the primary source of revenue for digital platforms. Consequently, the data intermediary in our sequential game sells the consumers’ data only to the producer, and cannot commit to withholding any information from him. Similarly, the intermediary’s choice of data outflow policy occurs after consumers are enlisted but before their data are realized. This assumption captures the limited ability of a platform to write advertising contracts contingent on, say, the volume of activity taking place at any point in time.

Linear price discrimination. The producer in our baseline model uses the consumers’ data to charge a unit price to each consumer. Whether richer pricing instruments enable online merchants to implement more sophisticated forms of price discrimination is largely an empirical question. However, as we discuss in Section 5, the general anonymization result of Proposition 8 guides the intermediary’s optimal data policy even if the producer can extract more of the surplus generated through better information. Indeed, even in the case of perfect price discrimination, the presence of the data externality would continue to drive a gap between equilibrium and socially efficient allocation.

3. Value of social data

□ **Data sharing and product market.** In our extensive form game, given the realized data inflow, the intermediary offers a data outflow policy to the producer. This policy specifies both the fee m_0 and the flow of information to all market participants, including the consumers. The data outflow policy thus determines both how well informed the producer is and how well informed his customers are.

Because the information is being sold to the producer, the intermediary will choose a data outflow policy that maximizes the producer’s profits, and then extracts the value of this information through the fee m_0 . The intermediary’s choice of outflow policy is simplified considerably by the following result, which allows us to restrict attention to policies in which every consumer is weakly better informed than the producer.

Lemma 1 (Data outflow policy). For any data inflow X , it is without loss of generality to consider data outflow policies where $Y_0(X)$ is measurable with respect to each $Y_i(X)$.

Interestingly, the result of Lemma 1 does not rely on the specific nature of the product market interaction. The key observation is that, if the producer were better informed, the prices charged would convey information to consumers about their own willingness to pay. The ensuing signaling incentives impose a cost on the producer because he will need to deviate from the prices that maximize his profits, holding fixed the consumers' beliefs. The intermediary can then increase the producer's profits (weakly) by revealing any information contained in the equilibrium prices directly to the consumers. Furthermore, this improvement is strict if, when the consumers receive this additional information, the producer modifies his choice of price.

Under the data outflow policies of Lemma 1, we can easily quantify the value of information for consumers and producers. The shared data help each consumer estimate her own willingness to pay. For the producer, the shared data enable a more informed pricing policy. In particular, given the realized data outflow $y \in Y$, the optimal pricing policy for the producer consists of a vector of (potentially personalized) prices p^* , thus resulting in a vector of individual quantities q_i^* purchased. We denote the predicted value of consumer i 's willingness to pay, given the signals (s_i, y_i) by:

$$\widehat{w}_i(s_i, y_i) \triangleq \mathbb{E}[w_i | s_i, y_i].$$

The realized demand of consumer i is given by

$$q_i(s_i, y_i, p) = \widehat{w}_i(s_i, y_i) - p.$$

As Y_0 is (weakly) less informative than Y_i , the producer chooses the optimal price

$$p_i^*(y_0) = \frac{\mathbb{E}[\widehat{w}_i(s_i, y_i) | Y_0] + c}{2} = \frac{\widehat{w}_i(y_0) + c}{2},$$

which results in the equilibrium quantity:

$$q_i^*(s_i, y_i, y_0) \triangleq q_i(s_i, y_i, p_i^*(y_0)).$$

The *ex ante* expected profit of the producer from interacting with consumer i is given by

$$\Pi_i((S_i, Y_i), Y_0) \triangleq \mathbb{E}[\pi(p_i^*(y_0), q_i^*(s_i, y_i, y_0)) | Y_0] = \frac{1}{4} \mathbb{E}[(\widehat{w}_i(y_0) - c)^2 | Y_0].$$

The first argument in $\Pi_i(\cdot, \cdot)$ refers to consumer i 's information structure (S_i, Y_i) and the second argument refers to the producer's information structure Y_0 . Similarly, we denote the indirect utility of consumer i as

$$U_i((S_i, Y_i), Y_0) \triangleq \mathbb{E}[u_i(w_i, q_i^*(s_i, y_i, y_0), p_i^*(y_0)) | S_i, Y_i] = \frac{1}{2} \mathbb{E}[(\widehat{w}_i(s_i, y_i) - p_i^*(y_0))^2 | S_i, Y_i].$$

□ **Welfare effects of data sharing.** The model with quadratic payoffs (regardless of the prior distribution of consumers' willingness to pay) yields explicit expressions for the value of information for product market participants. In particular, because prices and quantities are linear functions of the posterior mean \widehat{w}_i , the *ex ante* average prices and quantities $\mathbb{E}[p^*]$ and $\mathbb{E}[q^*]$ are constant across all information structures. Consequently, all surplus levels depend only on the *ex ante* variance of the posterior mean \widehat{w}_i under the consumers' and the merchant's information structures.

We therefore quantify the players' information gains under any information structure Y through the variance of posterior expectation:

$$G(Y) \triangleq \text{var}[\widehat{w}_i(y)], \quad (5)$$

and refer to it as the (informational) *gain function*. As we normalized the variance of the fundamental w_i to 1, the gain function $G(Y)$ represents the fraction of the variance of w_i explained by the signal y .

We now turn to the consequences of data sharing relative to no information sharing. Without information, the producer charges a constant price for all consumers based on the prior mean, denoted by \bar{p} . In contrast, the consumer already has an initial signal s_i , according to which she can adjust her quantity. The producer's net revenue and the consumer's expected utility under no information sharing are then given by

$$\begin{aligned}\Pi_i(S_i, \emptyset) &\triangleq \mathbb{E}[\pi(\bar{p}, q_i^*(s_i))], \\ U_i(S_i, \emptyset) &\triangleq \mathbb{E}[u_i(w_i, q_i^*(s_i), \bar{p})|S_i].\end{aligned}$$

We can now express the value of data sharing for the consumers and the producer in terms of the respective information gains.

Proposition 1 (Value of data outflow).

1. The value of data outflow Y for the producer is

$$\Pi_i((S_i, Y_i), Y_0) - \Pi_i(S_i, \emptyset) = \frac{1}{4}G(Y_0). \quad (6)$$

2. The value of data outflow Y for consumer i is

$$U_i((S_i, Y_i), Y_0) - U_i(S_i, \emptyset) = \frac{1}{2}(G(S_i, Y_i) - G(S_i)) - \frac{3}{8}G(Y_0). \quad (7)$$

3. The social value of data outflow Y is

$$W_i((S_i, Y_i), Y_0) - W_i(S_i, \emptyset) = \frac{1}{2}(G(S_i, Y_i) - G(S_i)) - \frac{1}{8}G(Y_0). \quad (8)$$

Thus, consumer and social surplus increase with the additional information learned by the consumers $G(S_i, Y_i) - G(S_i)$, and decrease with the information learned by the producer $G(Y_0)$. Intuitively, the welfare consequences of data sharing operate through two channels. First, with more information about her own preferences, the demand of each consumer is more responsive to her willingness to pay; this responsiveness is beneficial for consumers and (weakly) for the producer. Second, with access to better data, the producer adapts his pricing policy to the estimate of each consumer's willingness to pay \hat{w}_i . This price responsiveness dampens some of the quantity responsiveness. Hence, this second channel reduces both consumer surplus and total welfare.

Whether the first or the second channel dominates depends on the informativeness of the consumers' initial signals $G(S_i)$, and on the degree of correlation in the fundamental and error terms, which jointly determine any consumer's ability to learn from others' information. Corollary 1 formalizes this intuition by deriving the implications of Proposition 1 in several special cases.

Corollary 1 (Welfare effects).

1. If consumers cannot learn from each other's signals ($G(S) = G(S_i)$), any data sharing reduces consumer and social surplus.
2. If individual signals s_i are uninformative ($G(S_i) = 0$), any data sharing improves consumer and social surplus.
3. Social surplus is maximized by collecting and sharing all signals with every consumer ($Y_i = S$), and sharing no data with the producer ($Y_0 = \emptyset$).

Data sharing is detrimental to consumer and social surplus when consumers observe their willingness to pay perfectly ($\sigma = 0$), or when both fundamentals (w_i, w_j) and errors (e_i, e_j) are

independent. In those cases, any data sharing only enables price discrimination.⁴ Conversely, if individual signals become arbitrarily uninformative, but the entire vector s remains informative, then even *symmetric* information gains (i.e., data outflow policies where $Y_0 = Y_i$) yield Pareto improvements in the product market. In this case, the producer and the consumer share the additional gains from trade associated with better informed consumption and pricing decisions.

Finally, an immediate implication of the two channels highlighted by Proposition 1 is that the *first best* allocation of information consists of collecting and sharing all data among the consumers and none with the producer. However, a data intermediary with market power will not implement the socially optimal allocation of information. In particular, because the producer is paying for the data, he will always receive some information. To fully describe the outcome of the game, we then turn to the price of social data.

□ **Price of social data.** We first derive the total payment m_0 charged to the producer and the compensation m_i owed to each consumer. For the producer, the gains from data acquisition have to at least offset the price of the data. At the same time, the intermediary can charge up to the entire value of the information outflow Y_0 to the producer. From expression (6) in Proposition 1, we can then write the payment m_0 as

$$m_0(Y) = N(\Pi_i((S_i, Y_i), Y_0) - \Pi_i(S_i, \emptyset)) = \frac{N}{4}G(Y_0). \quad (9)$$

Not surprisingly, the intermediary's profits are increasing in the amount of information sold to the producer. However, we also know from Lemma 1 that every consumer i receives at least as much information as the producer. These two observations establish the optimality of the *complete data outflow* policy. Under this policy, the entire realized data inflow X is reported to the producer and to all consumers, including those who did not accept the intermediary's offer.

Lemma 2 (Optimal data outflow). Given any realized data inflow X , the complete data outflow policy, $Y_0^*(X) = Y_i^*(X) = X$ for all i , maximizes the gross revenue of the producer among all feasible data-outflow policies.

A critical driver of the consumer's decision to share data is her ability to anticipate the intermediary's use of the information thus gained. By Lemma 2, every consumer knows that all product market participants will receive the same information from the intermediary. In particular, consumer i knows that, by rejecting her contract, she prevents the producer from accessing her data X_i , but she does not forego the opportunity to learn from other consumers' data X_{-i} . In other words, consumer i can learn X_{-i} for free. In practice, this corresponds to searching for the same content on a digital platform without "logging in" or else agreeing to sharing her own data.

For any data inflow policy X and any underlying signal structure, the data intermediary offers positive payments to consumers. This occurs because the intermediary must compensate consumers *on the margin*: Consumer i requires compensation for revealing her data X_i to the producer, given that the other $N - 1$ consumers already share theirs.⁵ Therefore, the payments satisfy

$$m_i \geq U_i((S_i, X_{-i}), X_{-i}) - U_i((S_i, X), X). \quad (10)$$

Intuitively, consumer i is not compensated for the (positive or negative) effect of other consumers' data inflow X_{-i} on her surplus. To see this more formally, suppose consumer i 's participation

⁴ The special case in which each consumer knows her willingness to pay (i.e., signals are noiseless in our model's language) is closely related to the model of third-degree price discrimination in Robinson (1933) and Schmalensee (1981). In our setting, data sharing enables the producer to offer personalized prices; thus, price discrimination occurs across different *realizations* of the willingness to pay. In contrast, in Robinson (1933) and Schmalensee (1981), price discrimination occurs across different market segments. In both settings, the central result is that average demand does not change (with all markets served), but social welfare is lower under finer market segmentation.

⁵ We will revisit this property when we discuss the intermediary's commitment power in Section 5.

constraint (10) binds, and rewrite compensation m_i as

$$\begin{aligned} m_i^*(X) &= U_i((S_i, X_{-i}), X_{-i}) - U_i((S_i, X), X) \\ &= -\underbrace{(U_i((S_i, X), X) - U_i(S_i, \emptyset))}_{\triangleq \Delta U_i(X)} + \underbrace{U_i((S_i, X_{-i}), X_{-i}) - U_i(S_i, \emptyset)}_{\triangleq DE_i(X)}. \end{aligned} \quad (11)$$

The first term in (11), denoted by $\Delta U_i(X)$, is the total impact on consumer i 's surplus associated with data inflow X . The second term is the *data externality* (12) imposed on i by consumers $j \neq i$. It reflects the change in utility when $j \neq i$ sell their data X_j to the intermediary who then shares the data with the producer. As it is central to our analysis, we now examine the latter term in detail.

□ **Data externality and intermediation.** Our notion of *data externality* isolates the effect on consumer i 's surplus of the decision by the other consumers to share their data with all market participants.

Definition 1 (Data externality). The data externality imposed by consumers $-i$ on consumer i is given by

$$DE_i(X) \triangleq U_i((S_i, X_{-i}), X_{-i}) - U_i(S_i, \emptyset). \quad (12)$$

Using expression (7) in Proposition 1, the data externality can be written as

$$DE_i(X) = \frac{1}{2}(G(S_i, X_{-i}) - G(S_i)) - \frac{3}{8}G(X_{-i}). \quad (13)$$

To provide some intuition as to what determines the sign of the data externality, we evaluate expression (13) under the two special information structures in Examples 1 and 2. In both cases, we consider what would happen if consumer i held back her signal, given that the remaining $N - 1$ consumers share theirs with the producer and with consumer i .

Example 1 illustrates that if consumer i does not learn much from the signals of the other consumers, but those signals help predict w_i , then the data externality is negative.

Example 1 (Common preferences). Let $s_i = w + \sigma e_i$, and suppose the intermediary collects all consumer data, $X_i = S_i$. The producer can use $N - 1$ signals to estimate the common preference parameter w , that is, $G(S_{-i}) > 0$. If the individual signals are sufficiently precise so that $G(S) \approx G(S_i)$, then consumer i is worse off when other consumers share their signals, that is, the data externality is negative.

Example 2 illustrates that if the producer cannot learn anything about w_i from signals s_{-i} , then the data externality is unambiguously positive.

Example 2 (Common experience). Let $s_i = w_i + \sigma e$, and suppose the intermediary collects all consumer data, $X_i = S_i$. Because all w_i are independent, $G(S_{-i}) = 0$, that is, the producer cannot learn about w_i from signals s_{-i} only. However, consumer i can use signals s_{-i} to filter out the common error in her own signal s_i , that is, $G(S) > G(S_i)$. Therefore, other consumers' signals help consumer i , that is, the data externality is positive.

Thus, the overall effect of data sharing on consumer surplus (Proposition 1) depends largely on the informativeness of individual signals s_i . Conversely, the impact of other consumers' sharing decisions varies significantly with the data structure, particularly the ability of the producer to learn about w_i from signals s_{-i} .

4. Optimal data intermediation

■ The data externality has direct implications for the intermediary's profit (4). Combining the expressions for payments in (9) and (11), we can write the intermediary's profit R as

$$R(X) = \sum_{i=1}^N (\Delta W_i(X) - DE_i(X)), \quad (14)$$

where $\Delta W(X)$ denotes the effect of sharing data policy X on total surplus, as in (8). The intermediary's profits are equal to the effect of data sharing on social surplus, net of the data externalities across all consumers. The sign of the data externality is therefore critical for the profitability of data intermediation. In particular, if consumers impose negative data externalities on each other, this imposition directly reduces the compensation owed to each one, and conversely if the data externalities are positive. The revenue formula (14) clarifies how the intermediary's objective differs from the social planner's. In particular, if the data externality is negative, then intermediation can be profitable but welfare reducing. Conversely, if the data externality is positive, welfare-enhancing intermediation might not be profitable.

Having characterized the two terms ΔU_i and DE_i in (7) and (13), we can rewrite the payments to consumers in (11) as

$$m_i^*(X) = \frac{3}{8}(G(X) - G(X_{-i})). \quad (15)$$

Finally, combining the terms $m_0^*(X)$ in (9) and $m_i^*(X)$ in (15), we obtain

$$R(X) = \frac{N}{8}(3G(X_{-i}) - G(X)). \quad (16)$$

This yields a necessary and sufficient condition for profitable data intermediation.

Proposition 2 (Profitable data intermediation). Data intermediation with inflow policy X is profitable if and only if

$$3G(X_{-i}) \geq G(X).$$

Proposition 2 considers the amount of information learned by the producer. Specifically, it requires that the signals x_{-i} generate at least $1/3$ of the variance of w_i explained by the entire vector x in order for the data inflow policy X to be profitable. Intuitively, it is cheaper to acquire each signal x_i if the other consumers' signals are substitutes, and this is more likely to occur when the underlying fundamentals w_i and w_{-i} are correlated. Conversely, for independent fundamentals, $G(X_{-i}) = 0$ for any X , and intermediation is not profitable.

We now draw the implications for the intermediary's profits in our two leading examples.

Corollary 2 (Common preference). Suppose fundamentals w_i are perfectly correlated and errors e_i are independent. When N is large, data intermediation is always profitable. However, for sufficiently small σ , the data externality is negative, and data sharing reduces social surplus.

These results echo the findings of Acemoglu et al. (2021), who considered signals with diminishing marginal informativeness, and found socially excessive data intermediation. The information structure in Corollary 2 satisfies this submodularity property. In our model, however, socially insufficient intermediation can also occur. In particular, the intermediary may be unable to generate positive profits from socially efficient information with complementary signals, such as those in Corollary 3.

Corollary 3 (Common experience). Suppose fundamentals w_i are independent and errors e_i are perfectly correlated. For sufficiently large σ , data sharing increases social surplus. However, because the fundamentals are independent, data intermediation is never profitable.

The conclusions of Corollaries 2 and 3 do not depend on whether the intermediary collects complete data or anonymized data. However, as we discuss in the next section, it is always optimal for the intermediary to anonymize the data collected.

□ **Data anonymization.** We now explore the intermediary's decision to anonymize the individual consumers' demand data. We focus on two maximally different policies along this dimension. At one extreme, the intermediary can collect and transmit *complete (identity-revealing)* data about individual consumers ($X = S$), thereby enabling the producer to charge personalized prices. At the other extreme, the intermediary can collect *anonymized* data ($X = A$).

Under anonymized information intermediation, the producer charges the same price to all consumers who participate in the intermediary's data policy. In other words, from the point of view of the producer, anonymized data is equivalent to aggregate demand data. These data still allow the producer to perform third-degree price discrimination across realizations of the total market demand but limit his ability to extract surplus from individual consumers.⁶

Certainly, for the producer, the value of market demand data is lower than the value of individual demand data. However, the cost of acquiring such fine-grained data from consumers is also correspondingly higher. We now show that anonymizing the consumers' information profitably reduces the intermediary's data acquisition costs.

Proposition 3 (Optimality of data anonymization). The intermediary obtains strictly greater profits by collecting anonymized consumer data.

Within the confines of our policies, but independent of the distributions of fundamental and noise terms, the data intermediary finds it advantageous to not elicit the identity of the consumer. Therefore, the producer will not offer personalized prices but variable prices that adjust to the realized information about market demand. In other words, a monopolist intermediary might cause socially inefficient information transmission, but the equilibrium contractual outcome preserves privacy over the personal identity of the consumer.

In Section 5, we explore the boundaries of the anonymization result, under both heterogeneous consumers and alternative product-market specifications. In particular, we will generalize our result to show that anonymization is optimal when consumers are homogeneous *ex ante* and transmitting data to the producer is socially inefficient. Therefore, if information is used to target advertisements, and more precise messages increase social surplus, then we shall predict that the intermediary shares complete data.⁷

For the case of surplus extraction (prices), the finding in Proposition 3 suggests why we might see personalized prices in fewer settings than initially anticipated. In the context of direct sales of information, for example, Nielsen does not sell individual households' data to merchants. Instead, Nielsen aggregates its panel data at the local market level. Similarly, in the context of indirect sales of information, merchants on the retail platform Amazon very rarely engage in personalized pricing. However, the price of every single good or service is subject to substantial variation across both geographic markets and over time.

In light of the above result, we might interpret the restraint on the use of personalized pricing as the optimal resolution of the intermediary's trade-off in acquiring sensitive consumer

⁶ More formally, under the anonymized data policy A , the producer has access to the vector $\delta(s)$, that is, to a uniformly random permutation of the consumers' signals. Because the producer faces a prediction problem for each w_i with a convex loss, he chooses to charge a uniform price that is optimal for the sample average of the consumers' signals.

⁷ If, in addition, responsiveness to ads is heterogeneous in the consumer population, coarser information transmission (i.e., not fully anonymous) can also be optimal, as we show in Proposition 9.

information. Stretching the interpretation, the pushback against Amazon's introduction of personalized pricing in the early 2000s can also be viewed as a tightening of the consumers' participation constraints. Under any degree of competition, the net value of the Amazon platform under perfect personalized pricing would not have exceeded the consumers' outside options.

The data externality is, once again, the key to gaining intuition for why the intermediary chooses data anonymization. Suppose consumers $-i$ reveal their signals, and consumer i does not. With access to identifying information, the producer optimally aggregates the available data to form the best predictor of the missing data point. In this case, the producer charges a personalized price $p_i^*(X_{-i})$ to each nonparticipating consumer i . With anonymous data, the producer charges two prices: a single price for all participating consumers and another price for the deviating, nonparticipating consumers. Because the distribution of consumer willingness to pay and signals is symmetric, however, the producer's inference on w_i is invariant to permutations of the other consumers' signals, that is,

$$\mathbb{E}[w_i | a_{-i}] = \mathbb{E}[w_i | s_{-i}].$$

Therefore, a nonparticipating consumer faces identical prices under both data policies:⁸

$$p_i^*(S_{-i}) = p_i^*(A_{-i}).$$

Likewise, consumer i 's posterior distribution over her own w_i does not depend on the identity of the other consumers. Therefore, removing identity information through the anonymized policy $X = A$ does not have any implications for consumers' learning either:

$$\mathbb{E}[w_i | s_i, a_{-i}] = \mathbb{E}[w_i | s].$$

Because the amount of information available to consumer i and to the producer *off the path of play* is independent of $X \in \{A, S\}$, it follows that

$$U_i((S_i, S_{-i}), S_{-i}) = U_i((S_i, A_{-i}), A_{-i}).$$

In turn, this implies $DE_i(S) = DE_i(A)$. Thus, the data externality term $DE_i(X)$ in the intermediary's profits (14) is not impacted by the choice of inflow $X \in \{A, S\}$.

Along the path of play, however, the two data inflow policies yield different outcomes. In particular, the anonymized data inflow policy reduces the amount of information conveyed to the producer in equilibrium. Crucially, this reduction does not occur at the expense of the consumers' own learning. Therefore, the shift to anonymized data increases the total surplus terms $\Delta W_i(X)$ and the intermediary's profits. Put differently, anonymization reduces the cost of procuring the information, relative to the loss in revenue.

We now show that data anonymization is the key to the "explosive" profitability of data intermediation when the number of consumers becomes large. We then revisit the applications and limitations of our anonymization result in Section 5.

□ **Large markets.** Thus far, we have considered the optimal data policy for a given finite number of consumers, each of whom transmits a single signal. Perhaps, *the* defining feature of data markets is the multitude of (potential) participants, data sources, and services. We now pursue the implications of having many participants (i.e., of many data sources) for the social efficiency of data markets and the price of data.

Each additional consumer presents an additional opportunity for trade in the product market. Thus, the feasible social surplus is linear in the number of consumers. In addition, with every additional consumer, the intermediary obtains additional information about market demand.

⁸ Anonymization remains optimal if we force the producer to charge a single price to all consumers on and off the equilibrium. With this interpretation, we intend to capture the idea that the producer offers one price "on the platform" to the participating consumers while interacting with the deviating consumer "offline." The producer then uses the available market data to tailor the offline price.

These two effects suggest that intermediation becomes increasingly profitable in larger markets, wherein the potential revenue increases without bound, whereas individual consumers make a small marginal contribution to the precision of aggregate data.

For this comparative statics analysis, we adopt the following *additive data structure*. Specifically, we assume the willingness to pay of consumer i is the sum of two components:

$$w_i = \theta + \theta_i. \quad (17)$$

The term θ is *common* to all consumers in the market, whereas the term θ_i is *idiosyncratic* to consumer i . Similarly, the error term of consumer i is given by

$$e_i \triangleq \varepsilon + \varepsilon_i, \quad (18)$$

where the terms ε and ε_i refer to a common and an idiosyncratic error, respectively. We also refer to the willingness to pay w_i as the *fundamental* as opposed to the error term e_i .

As we vary the number of consumers N , the additive data structure allows us to hold the pairwise correlation between any two consumers' fundamentals and noise terms constant. In particular, let α denote the correlation coefficient of any two (w_i, w_j) , and let β denote the correlation coefficient of (e_i, e_j) .

We first establish a sufficient condition for the profitability of complete data sharing as the number of consumers becomes large, and then we analyze the data intermediary's revenue and total cost separately.

Proposition 4 (Profitable intermediation of anonymized data). For any $\alpha > 0$, there exists N^* such that anonymized data sharing is profitable if $N > N^*$.

We already know from Corollary 2 that a high degree of correlation in the consumers' willingness to pay allows the intermediary to profit from data sharing with sufficiently precise signals. Under the optimal data-sharing policy, *any* degree of correlation in the consumers' willingness to pay makes the anonymized signals sufficiently close substitutes that intermediation is profitable when N is large.

In Proposition 5, we assume that error terms are independent. This allows us to use the sample average to establish a lower bound on learning from $N - 1$ signals. We suspect that similar results hold more generally under correlated errors.⁹

Proposition 5 (Large markets). Consider the additive data structure and assume that errors are independent across consumers. As $N \rightarrow \infty$:

1. Each consumer's compensation m_i^* converges to zero.
2. Total consumer compensation is bounded by a constant,

$$Nm_i^* \leq \frac{9}{8}(\text{var}[\theta_i] + \text{var}[\varepsilon_i]), \quad \forall N.$$

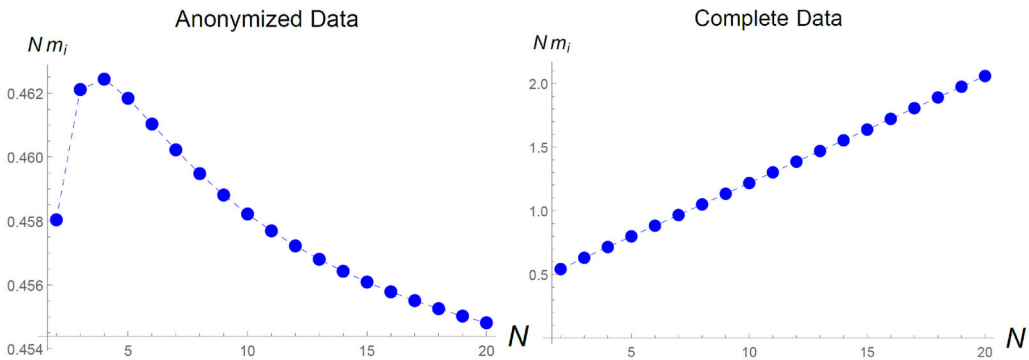
3. The intermediary's revenue and profit grow linearly in N .

As the optimal data policy aggregates the consumers' signals, each additional consumer has a rapidly decreasing marginal value. Furthermore, each consumer is paid only for her marginal contribution; this explains why the total payments Nm_i converge to a finite number. Strikingly, this convergence can occur from above: when the consumers' willingness to pay is sufficiently correlated, the decrease in each i 's marginal contribution can be sufficiently strong to offset the increase in N .

Whereas total costs converge to a constant, the revenue that the data intermediary can extract from the producer is linear in the number of consumers. Our model therefore implies that, as the

⁹ This result holds, for example, when both fundamentals and errors have Gaussian distributions.

FIGURE 2

TOTAL CONSUMER COMPENSATION ($\sigma_w = 1, \sigma_e = 0$)

market size grows without bound, the per capita profit of the data intermediary converges to the per capita profit when the (anonymized) data are freely available. Conversely, the impact on consumer surplus depends on the degree of correlation in the underlying fundamentals and on the precision of the consumers' initial signals.¹⁰

Finally, we show that data anonymization is crucial for the large N properties of the intermediary's profits. Recall that, with complete data intermediation, individual consumer payments are proportional to $G(S) - G(S_{-i})$. As long as fundamentals w_i are not perfectly correlated, these payments are then bounded away from zero for any finite N . Proposition 6 shows that this property also holds in the limit.

Proposition 6 (Asymptotics with complete sharing). Consider the additive data structure with $\text{var}[\theta_i] > 0$. Under complete (identity-revealing) data sharing, the asymptotic individual compensation is bounded away from 0:

$$\liminf_{N \rightarrow \infty} m_i^* \geq \frac{3}{8} \frac{\text{var}^2[\theta_i]}{1 + \text{var}[e_i]} > 0.$$

An immediate consequence of Proposition 6 is that, with complete data sharing, total payments to consumers grow linearly in N . Thus, anonymization is critical to achieving increasing returns to scale in data intermediation: even in settings where complete data intermediation $X = S$ is profitable, the per capita profits are bounded away from the full value of information.

Figure 2 illustrates an example with normally distributed fundamentals and errors, in which it can be less expensive for the intermediary to acquire a larger *anonymized* dataset than a smaller one, but not a larger *complete* dataset.

□ **Unique implementation.** Our analysis thus far has characterized the intermediary's most preferred equilibrium. An ensuing question is whether the qualitative insights and the asymptotic properties discussed above would hold across all equilibria, particularly in the intermediary's least preferred equilibrium. A seminal result in the literature on contracting with externalities (see Segal (1999)) is the “divide-and-conquer” scheme that guarantees a unique equilibrium outcome (see Segal and Whinston (2000) and Miklos-Thal and Shaffer (2016)). Under this scheme, the intermediary can sequentially approach consumers and offer compensation conditional on all earlier consumers having accepted an offer. In this scheme, the first consumer receives

¹⁰ In a recent contribution, Loertscher and Marx (2020) study large digital monopoly markets, where data have the countervailing effects of improving consumer valuations and increasing monopoly prices.

compensation equal to her entire surplus loss, thereby guaranteeing her acceptance regardless of the other consumers' decisions. More generally, consumer i receives the optimal compensation level in the baseline equilibrium when $N = i$.

The cost of acquiring the consumers' data is strictly higher under "divide and conquer" than in the intermediary's most preferred equilibrium. Nonetheless, the impact of the ensuring unique implementation on per capita profits vanishes in the limit.

Proposition 7 ("Divide and conquer"). Consider the additive data structure with independent errors. Under the "Divide and Conquer" scheme, total consumer compensation satisfies

$$Nm_i^* \leq \frac{3}{4}(1 + \log N)(\text{var}[\theta_i] + \text{var}[\varepsilon_i]).$$

Under divide and conquer, the total payments to the consumers do not converge to a finite constant as N grows without bound. However, the growth rate of these payments is far smaller than the rate at which the producer's willingness to pay for data diverges. Therefore, regardless of the equilibrium-selection criterion, the intermediary's per capita profits converge to the benchmark level when anonymized consumer data are made available.

5. Implications for consumer privacy

■ In this section, we enrich our model along several dimensions to characterize the implications of the optimal data intermediation policy for consumer privacy. In particular, we consider richer pricing instruments in the product market, heterogeneous consumers, heterogeneous product varieties, noisy information collection, and commitment power in the data market.

□ **Data anonymization and social efficiency.** In our baseline setting, data anonymization is optimal independent of the model parameters, such as the number of consumers or the distribution of fundamentals and error terms. This result relies on two crucial assumptions: (i) consumer are homogeneous, by which we mean that the distributions of fundamentals w and errors e are symmetric, and (ii) data sharing has unambiguous welfare effects on product market participants. Indeed, in the model of linear price discrimination, transmitting X_i anonymously improves consumer and social surplus, relative to complete data intermediation.

We can generalize this insight to *any arbitrary* product market interaction beyond the linear pricing model of the previous section. Proposition 8 generalizes the intuition behind the optimal anonymization result in Proposition 3: It establishes that *social surplus* is the criterion guiding the intermediary's decision to optimally collect anonymized data.

Proposition 8 (Social optimality of data anonymization). With homogeneous consumers, anonymized data intermediation is more profitable than complete data intermediation if and only if anonymization increases social surplus.

We note two important aspects of this result. First, it establishes the congruence between the intermediary private objective and the social welfare with respect to the anonymization decision *only*. It does not claim that the equilibrium information flow itself is socially efficient. Second, the argument does not require any specific feature of the product market interaction. As the decision between anonymization and de-anonymization pertains precisely to the marginal value of the private information of i for the prediction of the willingness to pay w_i , intermediary and consumer i can attain a socially efficient arrangement.

The result has immediate implications for how equilibrium data sharing policies depend on the nature of the product market interaction. In particular, Proposition 8 allows us to examine the role of richer pricing instruments. Bergemann, Brooks, and Morris (2015) have shown that every feasible combination of consumer and producer surplus is consistent with *some* form of

price discrimination. Proposition 8 shows that, if the producer had the ability to extract all the expected surplus (given the consumers' information), then the intermediary would find it more profitable to collect complete, identifying data.

A canonical example where this prediction is relevant is the case of unit demand by the consumers, where our model would predict the prevalence of perfect price discrimination. However, to the extent to which consumers have options to retain some surplus *ex post*, such as by scaling down their purchase level as in our baseline model, then full surplus extraction would require the producer to have access to more complex pricing mechanisms.

□ **Market segmentation and data.** The assumption of *ex ante* homogeneity among consumers has enabled us to produce some of the central implications of social data. A more complete description of consumer demand should introduce heterogeneity across groups of consumers along characteristics such as location, demographics, income, and wealth.

We now explore how these additional characteristics influence information policy and the profits of the data intermediary. To this end, we augment the description of consumer demand by splitting the population into J homogeneous groups:

$$w_{ij} \sim F_{w,j}, e_{i,j} \sim F_{e,j}, i = 1, 2, \dots, N_j, j = 1, \dots, J.$$

The intermediary's data inflow policy must now specify whether to anonymize the consumers' signals across groups and within each group. However, Proposition 8 establishes that it is always more profitable to anonymize all signals within each group, rather than revealing the consumers' identities.

Corollary 4 (No discrimination within groups). The data policy that anonymizes all signals within each group $j = 1, \dots, J$ and only reveals the group identity of each consumer i is more profitable than the complete data-sharing policy.

By further specifying the model, we can identify conditions under which the data intermediary will collect and transmit group characteristics. By collecting information about the group characteristics, the intermediary influences the extent of price discrimination. For example, the intermediary could anonymize all signals across groups, thus forcing the producer to offer only a single price. Alternatively, the intermediary could allow the producer to discriminate between two groups of consumers by recording and transmitting the group identities. As intuition would suggest, enabling price discrimination across groups not only allows the intermediary to charge a higher fee to the producer but also increases the compensation owed to consumers.

Proposition 9 below sheds light on the optimal resolution of this trade-off. In this result, we restrict attention to the case of symmetric groups ($N_j = N$ for all j), with the additive data structure $w_i = \theta + \theta_i$, and independent noise terms in the consumers' signals.

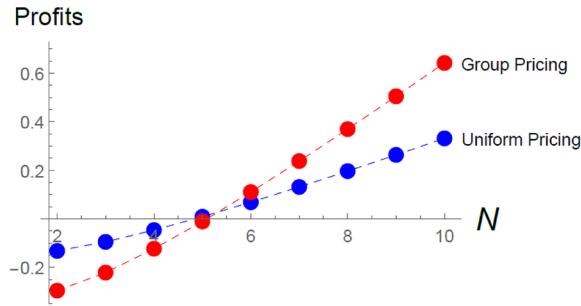
Proposition 9 (Segmentation). If N is large enough, inducing group-level pricing is more profitable for the intermediary than inducing uniform pricing.

Whereas Proposition 3 stated that the intermediary will not reveal any information about consumer identity, Proposition 9 refines that result: If the market is sufficiently large, then the intermediary will convey limited identity information, that is, each consumer's group identity. This policy allows the producer to price discriminate across, but not within, groups. Conversely, if the producer faces few consumers and their willingness to pay is not highly correlated, then pooling all signals reduces the cost of sourcing the data.

The limited amount of price discrimination, which operates optimally at the group level rather than the individual level, can explain the behavior of many platforms. For example, Uber and Amazon claim that they do not discriminate at the individual level, but they condition prices on location, time, and other dimensions that capture group characteristics.

FIGURE 3

MARGINAL VALUE OF AN ADDITIONAL CONSUMER



The result in Proposition 9 is perhaps the sharpest manifestation of the value of big data. By enabling the producer to adopt a richer pricing model, a larger database allows the intermediary to extract more surplus. Our result also clarifies the appetite of the platforms for large datasets: Because having more consumers allows the platform to profitably segment the market more precisely, the value of the marginal consumer $i = N$ to the intermediary remains large even as N grows. In other words, allowing the producer to segment the market is akin to paying a fixed cost (i.e., higher compensation to the current consumers) to access a better technology (i.e., one that scales more easily with N). Figure 3 illustrates this result for an example with normally distributed fundamentals and errors.

The optimality of using a richer pricing model when larger datasets are available is reminiscent of model selection criteria under overfitting concerns, for example, the Akaike information criterion. In our setting, however, the optimality of inducing segmentation is not driven by econometric considerations. Instead, it is entirely driven by the intermediary's cost–benefit analysis in acquiring more precise information from consumers. As the data externality grows sufficiently strong, acquiring the data becomes cheaper as the intermediary exploits the richer structure of consumer demand.¹¹

Finally, the welfare ranking of the two pricing schemes (group vs. uniform) is *a priori* ambiguous. In particular, group pricing can yield lower total surplus than uniform pricing, for example, when consumers know their willingness to pay and group pricing results in inefficient price discrimination. In other words, outside of the conditions of Proposition 8 (e.g., with heterogeneous consumers), market segmentation can be driven purely by the data externality.

□ **Recommender system.** In our baseline model, the data shared by the intermediary are used by the producer to set prices and by consumers to learn about their own preferences. The first assumption is, in a sense, the worst-case scenario for the intermediary: consider the case in which consumers' initial signals are very precise. As price discrimination reduces total surplus, no intermediation would be profitable without a strong, negative data externality. Consequently, data aggregation is an essential part of the optimal data intermediation policy in this case. In practice, however, consumer data can also be used by the producer in surplus-enhancing ways, for example, to facilitate targeting quality levels and other product characteristics to the consumer's tastes.

In this section, we develop a generalization of our framework; this generalization allows the producer to charge a unit price p_i and to offer a product of characteristic k_i to each consumer.

¹¹ Montiel Olea et al. (2019) offer a demand-side explanation of a similar phenomenon: They showed that data buyers who employ a richer pricing model are willing to pay more for larger datasets.

Consumers differ both in their vertical willingness to pay and in their horizontal taste for the product's characteristics. Consumer i 's utility function is given by

$$u_i(w_i, q_i, p_i, k_i, \ell_i) = (w_i - (k_i - \ell_i)^2 - p_i)q_i - q_i^2/2,$$

with w_i denoting the consumer's willingness to pay and ℓ_i denoting the consumer's ideal location or product characteristic. Both the willingness to pay $w \in \mathbb{R}^N$ and the locations $\ell \in \mathbb{R}^N$ of different consumers are potentially correlated. The producer has a constant marginal cost of quantity provision that we normalize to zero and can freely set the product's characteristic. Therefore, the case of a common location $\ell_i \equiv \ell$ for all consumers yields the baseline model of price discrimination.

We examine the data intermediary's optimal data inflow policy, which allows for separate aggregation policies for willingness to pay and location information. We impose the following assumptions: (i) the gains from trade under no information sharing are sufficiently large; (ii) the consumers' fundamentals have a joint Gaussian distribution; (iii) consumer i perfectly observes (w_i, ℓ_i) . The extension to noisy Gaussian signals is immediate. We then obtain another application of Proposition 8.

Proposition 10 (Optimal aggregation by a recommender system). The intermediary's optimal policy collects anonymized data on the vertical component w_i and complete data on the horizontal component ℓ_i .

Therefore, the recommender system enables the producer to offer targeted product characteristics that match k_i to ℓ_i as closely as possible. However, the system does not allow for personalized pricing. The logic is once again given by the intermediary's sources of profits, that is, the contribution to social welfare ΔW and the data externality DE . Because the data externalities do not depend on the level of data anonymization, the intermediary chooses to aggregate the vertical dimension of consumer data, thereby reducing the total surplus if transmitted to the producer. Conversely, because the distance between a consumer's ideal product and the firm's offer $(k_i - \ell_i)^2$ shifts the consumer's demand function down, the intermediary allows for the personalization of product characteristics.

□ **Intermediary with commitment.** We have assumed thus far that the intermediary cannot refrain from selling information to the producer and cannot sell any acquired data inflow back to the consumers. The latter assumption entails no loss: consumers know that the intermediary sells the data to the producer, and therefore expect to receive all available information regardless of their participation decisions (Proposition 1), because they know that revealing this information maximizes the producer's fee.

The no commitment assumption reflects the substantial control that large online platforms have over the use of the data and the opacity with which the data outflow is linked to the data inflow. In other words, it is difficult to ascertain how any given data input informs an intermediary's data output. Nonetheless, it is useful to consider the implications of the data intermediary's ability to commit to a certain data policy, especially in light of the welfare properties of data sharing discussed above. To that end, suppose the data intermediary offered the consumers contracts that specify a data inflow *and* a data outflow policy.

Through richer contracts, the data intermediary can offer consumers privacy guarantees. In particular, the intermediary can implement the socially efficient data-sharing policy, which consists of sharing all signals among all the consumers who accept the contract and not sharing any data at all with the producer (Corollary 1). In exchange for this commitment, the data intermediary requests compensation from the consumers. In turn, consumers are willing to pay a positive price for these data, and hence the socially optimal data sharing is always profitable.¹²

¹² This environment with commitment is related to the analysis in Lizzeri (1999) but has a number of distinct features. First, in Lizzeri (1999), the private information is held by a single agent, and multiple downstream firms compete

However, the equilibrium outcome under these stronger commitment assumptions fails to capture the role of large online platforms. Even though there are examples in which consumers pay a positive price to access tailored, non-sponsored recommendations, data intermediaries choose to monetize the producers' side of their platform much more frequently. Moreover, the socially efficient policy need not maximize the intermediary's profits. For example, if fundamentals are perfectly correlated and signals are arbitrarily precise, the intermediary's profits from the first-best policy are nil. Under these conditions, an intermediary with commitment (or even an intermediary who sells its own products) would not monetize their information by selling it to consumers.

It is beyond the scope of this article to characterize the optimal commitment policy for any initial data structure, but the data externality clearly remains a key driver of the equilibrium allocation of information even under stronger commitment assumptions.

6. Conclusion

■ We have explored the trading of information between data intermediaries with market power and multiple consumers with correlated preferences. The data externality that we have uncovered strongly suggests that levels of compensation close to zero can induce an individual consumer in a large market to relinquish precise information about her preferences. This finding holds even if the consumer's data are later sold to a firm that seeks to extract their surplus. Thus, giving consumers control rights over their data (a pillar of privacy regulation such as the EU General Data Protection Regulation or the California Privacy Rights Act) is insufficient to bring about the efficient use of information.

Our results regarding the aggregation of consumer information further suggest that privacy regulations must move away from concerns over personalized prices at the individual level. Most often, firms do not set prices in response to individual-level characteristics. Instead, segmentation of consumers occurs at the group level (e.g., as in the case of Uber) or at the temporal and spatial levels (e.g., as in the case of Staples and Amazon). Thus, our analysis points to the significant welfare effects of group-level and market-level dynamic prices that react in real time to changes in demand.

A possible mitigator of the consequences of data externalities—echoed in Posner and Weyl (2018)—consists of facilitating the formation of consumer groups or unions to internalize the data externality when bargaining with powerful intermediaries, such as large online platforms.¹³ A different regulatory solution is based on *privacy managers*, such as internet browsers with heterogeneous privacy settings that compete for consumers' default choice. Yet another solution—suggested by Romer (2019)—consists of making the data outflow costly for the intermediary by, for example, taxing targeted advertising. In our model, taxing the data outflow will limit efficient and inefficient intermediation alike but will affect the intermediary's choice of data policy under the assumptions of Section 5.

Finally, our data intermediary collected and redistributed the consumer data but played no role in the interaction between the consumers and the producer. In contrast, a consumer can often access a given producer only through a data platform.¹⁴ Many platforms can then be thought of as auctioning access to the consumer. The data platform provides the bidding producers with additional information that they can use to tailor their interactions with consumers. Social data

for the information and for the object offered by the agent. Second, the privately informed agent enters the contract after she has observed her private information. The shared insight is that the intermediary *with* commitment power might be able to extract a rent without any influence on the efficiency of the allocation.

¹³ This result echoes the claim in Zuboff (2019) that “privacy is a public issue.”

¹⁴ Product data platforms, such as Amazon, Uber and Lyft, acquire individual data from the consumer through the consumers' purchase of services and products. Social data platforms, such as Google and Facebook, offer data services to individual users and sell the information to third parties, who mostly purchase the information in the form of targeted advertising space. In terms of our model, a product data platform combines the roles of data intermediation and product pricing.

platforms thus trade individual consumer information for services rather than money. In these markets, the data externality manifests itself in the quality of the services offered and in the extent of the consumers' engagement.

Appendix A

Proof of Lemma 1. Under an arbitrary data-inflow policy X , each consumer i observes a noisy signal S_i of her own willingness to pay and sends a potentially noisier signal X_i to the intermediary.¹⁵ Consumer i observes both S_i and X_i . Given the data inflow X , the intermediary chooses an outflow policy, namely, the signal $Y_0 = Y_0(X)$ sent to the producer and the signal $Y_i = Y_i(X)$ sent to each consumer i . The intermediary then chooses a policy Y that maximizes the producer's *ex ante* expected payoff, which it fully extracts through the fee m_0 . We let the intermediary select their favorite equilibrium in the ensuing game.

For any outflow policy $Y = (Y_0, Y_i)$ denote an induced signaling equilibrium as $\bar{\gamma} = (\bar{q}_i, \bar{p})$, where $\bar{p} : Y_0 \rightarrow R^+$ is the pricing strategy of the producer and $\bar{q}_i : Y_i \times S_i \times X_i \times R^+ \rightarrow R^+$ is the demand function of consumer i . We first argue that there exists an equilibrium γ^* under the outflow policy $(\bar{p} \circ Y_0, (Y_i, \bar{p} \circ Y_0))$ that yields a weakly higher *ex ante* payoff for the producer. In this new outflow policy, instead of revealing Y_0 to the producer, the intermediary directly recommends the price $\bar{p}(Y)$ which coincides with the equilibrium pricing strategy in $\bar{\gamma}$ and reveals to consumer i both Y_i and the price recommendation.

On the equilibrium path of $\bar{\gamma}$, consumer i updates her posterior $\mu(Y_i, S_i, \bar{p}_i(Y))$ using Y_i , her own private signal S_i , the report X_i , and the observed price p_i . We denote the consumer's demand as a function of her posterior beliefs and the price as

$$q_i(\mu(Y_i, S_i, X_i, p_i), p_i).$$

The *ex ante* profit of the producer from consumer i is given by

$$\mathbb{E}[\bar{p}_i q_i(\mu(Y_i, S_i, X_i, \bar{p}_i), \bar{p}_i)].$$

Now consider the new outflow policy $(\bar{p} \circ Y_0, (Y_i, \bar{p} \circ Y_0))$. Under this policy, there exists an equilibrium where consumer i forms her demand using the data outflow $(Y_i, \bar{p} \circ Y_0)$ from the intermediary as well as her own signal S_i and the data inflow X_i . Because consumer i knows everything that the producer knows, the price charged by the producer no longer influences the consumer's posterior, which therefore coincides with the on-path beliefs in the original equilibrium $\bar{\gamma}$, that is,

$$\mu(Y_i, S_i, X_i, \bar{p}_i(Y_0)).$$

Knowing this, the producer maximizes his *ex ante* payoff by choosing a pricing strategy $\hat{p}(\cdot)$ as a function of his signal $\bar{p} \circ Y_0$. Thus the producer's equilibrium profit is given by

$$\max_{\hat{p}} \hat{p}(\bar{p} \circ Y_0) q_i(\mu(Y_i, S_i, X_i, \bar{p}_i(Y_0)), \hat{p}(\bar{p} \circ Y_0)).$$

Clearly "following the intermediary's recommendation," that is, setting $\hat{p}(p) = p$ is a feasible strategy that yields the same payoff as in the old equilibrium $\bar{\gamma}$. Consequently, the producer's equilibrium payoff is weakly higher than in $\bar{\gamma}$. \square

Proof of Proposition 1. For any offered price p_i , consumer i demands the quantity

$$q_i = \mathbb{E}[w_i | (S_i, Y_i)] - p_i.$$

The producer finds it optimal to set the following price

$$p_i = \frac{\mathbb{E}[w_i | Y_0]}{2}.$$

Recall that the consumer always has superior information so that Y_0 is measurable with respect to Y_i . The profit of the producer is given by

$$\begin{aligned} \Pi_i((S_i, Y_i), Y_0) &= \mathbb{E} \left[\frac{\mathbb{E}[w_i | Y_0]}{2} \left(\mathbb{E}[w_i | (S_i, Y_i)] - \frac{\mathbb{E}[w_i | Y_0]}{2} \right) \right] \\ &= \frac{\mathbb{E}[(\mathbb{E}[w_i | Y_0])^2]}{4} = \frac{\text{var}[\mathbb{E}[w_i | Y_0]] + \mathbb{E}[w_i]^2}{4} = \frac{1}{4} G(Y_0) + \Pi_i(S_i, \emptyset), \end{aligned}$$

¹⁵ Under complete data sharing, for example, the consumer either reports $X_i = S_i$ or refuses to participate, so that X_i has infinite variance (or the corresponding σ -algebra is the empty set).

where the outside expectation represents integration over the whole probability space. The expected consumer surplus is given by

$$U_i((S_i, Y_i), Y_0) = \mathbb{E} \left[\left(w_i - \frac{\mathbb{E}[w_i|Y_0]}{2} \right) \left(\mathbb{E}[w_i|(S_i, Y_i)] - \frac{\mathbb{E}[w_i|Y_0]}{2} \right) \right] - \frac{1}{2} \mathbb{E} \left[\left(\mathbb{E}[w_i|(S_i, Y_i)] - \frac{\mathbb{E}[w_i|Y_0]}{2} \right)^2 \right] \quad (\text{A1})$$

$$\begin{aligned} &= \frac{1}{2} \mathbb{E}[(\mathbb{E}[w_i|(S_i, Y_i)])^2] - \frac{3}{4} (\mathbb{E}[w_i|Y_0])^2, \\ &= \frac{1}{2} (G((S_i, Y_i)) - G(S_i)) - \frac{3}{8} G(Y_0) + U_i(S_i, \emptyset) \end{aligned} \quad (\text{A2})$$

Finally, the impact on total surplus is given by the sum of the two effects:

$$W_i((S_i, Y_i), Y_0) - W_i(S_i, \emptyset) = \frac{1}{2} (G((S_i, Y_i)) - G(S_i)) - \frac{1}{8} G(Y_0),$$

which completes the proof. \square

Proof of Proposition 2. By Lemma 1, it is without loss of generality to assume the producer receives a signal Y , and the consumer receives a signal (Y_i, Y) . Thus, we can focus on equilibria where prices have no signaling effect. These equilibria coincide with those described in Proposition 1. As we have shown there, the profit of the producer is:

$$\mathbb{E} \left[\frac{\mathbb{E}[w_i|Y]}{2} \left(\mathbb{E}[w_i|Y \cup Y_i, S_i, X_i] - \frac{\mathbb{E}[w_i|Y]}{2} \right) \right] = \frac{\mathbb{E}[(\mathbb{E}[w_i|Y])^2]}{4} = \frac{\text{var}[\mathbb{E}[w_i|Y]] + \mathbb{E}[w_i]^2}{4}.$$

Therefore it is optimal to maximize $\text{var}[\mathbb{E}[w_i|Y]]$, which is achieved by setting $Y = X$. Hence, the intermediary reveals all information collected ($Y = X$) both to the producer and to consumer i . \square

Proof of Corollary 2. When fundamentals w_i are perfectly correlated,

$$\mathbb{E}[w_i|S] = \mathbb{E}[\theta|S] = \mathbb{E}[\theta|S_1, \dots, S_N],$$

$$\mathbb{E}[w_i|S_{-i}] = \mathbb{E}[\theta|S_{-i}],$$

$$\text{var}[\mathbb{E}[\theta|S_{-i}]] = \text{var}[\mathbb{E}[\theta|S_1, \dots, S_{N-1}]].$$

Under our symmetry assumption, the variance of the posterior expectation of the common willingness to pay $\text{var}[\mathbb{E}[\theta|S_{1,\dots,N}]]$ can be written as a function of N . Now we argue that $\text{var}[\mathbb{E}[\theta|S_{1,\dots,N}]]$ is increasing in N . We first define $g(S_{1,\dots,N-1}) \triangleq \mathbb{E}[\theta|S_{1,\dots,N-1}]$. Then, according to Lemma A1 below, we have

$$\begin{aligned} \text{var}[\mathbb{E}[\theta|S_{1,\dots,N}]] &= \max_{f \in L^2} \text{var}[\theta] - \mathbb{E}[(\theta - g(S_{1,\dots,N-1}))^2], \\ &\geq \max_{f \in L^2} \text{var}[\theta] - \mathbb{E}[(\theta - g(S_{1,\dots,N-1}))^2], \\ &= \text{var}[\mathbb{E}[\theta|S_{1,\dots,N-1}]]. \end{aligned}$$

The sequence $\text{var}[\mathbb{E}[\theta|S_{1,\dots,N}]]$ is increasing and bounded. Therefore, it converges:

$$\lim_{N \rightarrow \infty} G(S) = \lim_{N \rightarrow \infty} G(S_{-i}),$$

and intermediation is then profitable:

$$\lim_{N \rightarrow \infty} \frac{R(S)}{N} = \frac{1}{4} \lim_{N \rightarrow \infty} G(S) > 0.$$

In the limit for $N \rightarrow \infty$, the data externality and the consumer surplus are given by

$$\begin{aligned} \lim_{N \rightarrow \infty} U_i(S, S) - U_i(S_i, \emptyset) &= \lim_{N \rightarrow \infty} \mathbb{E} \left[\frac{1}{8} (\mathbb{E}[w_i|S])^2 - \frac{1}{2} (\mathbb{E}[w_i|S_i] - \mathbb{E}[w_i])^2 - \frac{1}{8} \mathbb{E}[w_i]^2 \right] \\ &= \lim_{N \rightarrow \infty} \frac{1}{8} \text{var}[\mathbb{E}[w_i|S]] - \frac{1}{2} \text{var}[\mathbb{E}[w_i|S_i]], \\ \lim_{N \rightarrow \infty} DE_i(S) &= \lim_{N \rightarrow \infty} \frac{1}{2} \text{var}[\mathbb{E}[w_i|S]] - \frac{1}{2} \text{var}[\mathbb{E}[w_i|S_i]] - \frac{3}{8} \text{var}[\mathbb{E}[w_i|S_{-i}]] \\ &= \lim_{N \rightarrow \infty} \frac{1}{8} \text{var}[\mathbb{E}[w_i|S]] - \frac{1}{2} \text{var}[\mathbb{E}[w_i|S_i]]. \end{aligned}$$

Therefore, when the initial noise is sufficiently small (i.e., when $\text{var}[\mathbb{E}[w_i|S_i]]$ is close to $\text{var}[w_i]$), the data externality is negative and data sharing hurts consumers. \square

Proof of Corollary 3. Because w_i is independent from the other consumers' signals, we have $\text{var}[\mathbb{E}[w_i|S_{-i}]] = 0$. Thus, intermediation is always unprofitable, and the data externality is always positive,

$$R(S) = -\frac{N}{8}\text{var}[\mathbb{E}[w_i|S]] < 0,$$

$$DE(S) = \frac{1}{2}(\text{var}[\mathbb{E}[w_i|S]] - \text{var}[\mathbb{E}[w_i|S_i]]) \geq 0.$$

Finally, for the results on consumer surplus, we turn to Lemma A1. In particular, we know

$$\begin{aligned}\text{var}[\mathbb{E}[w_i|S]] &= \text{var}[w_i] - \mathbb{E}[(w_i - \mathbb{E}[w_i|S])^2], \\ &\geq \text{var}[w_i] - \mathbb{E}\left[\left(w_i - \left(s_i - \frac{1}{N-1}\sum_{j \neq i} s_j\right)\right)^2\right], \\ &= \text{var}[\theta_i] - \mathbb{E}\left[\left(\theta_i - \left(\theta_i + \varepsilon - \frac{1}{N-1}\sum_{j \neq i} \theta_j - \varepsilon\right)\right)^2\right], \\ &= \text{var}[\theta_i] - \mathbb{E}\left[\left(\frac{1}{N-1}\sum_{j \neq i} \theta_j\right)^2\right] = \frac{N-2}{N-1}\text{var}[\theta_i] \rightarrow \text{var}[w_i].\end{aligned}$$

Thus, we obtain

$$\lim_{N \rightarrow \infty} U_i(S, S) - U_i(S_i, \emptyset) = \frac{1}{8}\text{var}[w_i] - \frac{1}{2}\text{var}[\mathbb{E}[w_i|S_i]].$$

When σ is sufficiently large, so that $\text{var}[\mathbb{E}[w_i|S_i]]$ is close to 0, intermediation increases consumer surplus. \square

The proof of Proposition 2 follows from expressions (15) and (16) in the text.

Proof of Proposition 3. In the main text, the data inflow from consumer i is given by $X_i = S_i$ (under complete sharing) and we compare it with $X_i^* = \delta(S_i)$ (under anonymization). Note that $(S_i, X_i)_i$ in this case is symmetrically distributed, that is, its joint density is unchanged under permutations of indices. Here, we prove a slightly more general version of the result by allowing an arbitrary information inflow X such that $(S_i, X_i)_i$ is symmetrically distributed.¹⁶ We assume that apart from the private signal S_i and information outflow Y_i provided by the intermediary, consumer i also observes her own data inflow X_i . This assumption is needed because information set (S_i, X_i, X_{-i}^*) and (S_i, X_i, X^*) are equivalent by construction, but (S_i, X_{-i}^*) and (S_i, X^*) maybe not. Note that when we restrict to the less general case (when $X_i = S_i$), the latter holds automatically, so we do not need this assumption.

For any fixed inflow policy X , we refer to p_{-i} as the off-path price charged to consumer i when she does not accept the intermediary's contract, and to p_i as the on-path price charged to consumer i . Now consider another inflow policy X^* identical to X up to a random permutation of the consumers' identities. Under this scheme, we refer to p_{-i}^* as the off-path price for consumer i , and to p_i^* as the on-path price for consumer i .

We first argue that $p_{-i} = p_{-i}^*$ for any realization of W, S, X . To do so, let us calculate consumer i 's posterior about W_i under each inflow policy. Under the non-anonymized scheme, the posterior distribution of consumer i 's willingness to pay is given by

$$\begin{aligned}f_i(W_i = w_i | S_i = s_i, X_i = x_i, X = x) \\ = \frac{\int f(W_i = w_i, W_{-i} = w'_{-i}, S_i = s_i, S_{-i} = s'_{-i}, X_i = x_i, X_{-i} = x_{-i}) ds'_{-i} dw'_{-i}}{\int f(W = w', S_i = s_i, S_{-i} = s'_{-i}, X_i = x_i, X_{-i} = x_{-i}) ds'_{-i} dw'}.\end{aligned}$$

Recall from Proposition 2 that the intermediary's optimal data outflow policy consists of revealing to the consumers all the available information, even if the consumer refuses to participate. When the data is anonymized, because consumer i knows her own report X_i , the data outflow reveals to her the vector of reports X_{-i} without knowing who generated each one. We now define $\delta \in S^{n-1}$ as permutation of consumer indices. Consumer i 's posterior distribution over her willingness to pay w_i is now given by

$$f_i(W_i = w_i | S_i = s_i, X_i = x_i, X_{-i}^* = x_{-i}).$$

¹⁶ For example, in Section B, X_i might be a noisier signal of S_i .

For notational simplicity, we use \Pr to denote both probability and the proper marginal density. Then the posterior can be rewritten as

$$\begin{aligned} & \frac{\Pr(W_i = w_i, S_i = s_i, X_i = x_i, X_{-i}^* = x_{-i})}{\Pr(S_i = s_i, X_i = x_i, X_{-i}^* = x_{-i})} \\ &= \frac{\sum_{\delta \in S^{n-1}} \Pr(\delta, W_i = w_i, S_i = s_i, X_i = x_i, X_{-i} = x_{\delta(-i)})}{\sum_{\delta \in S^{n-1}} \Pr(\delta, S_i = s_i, X_i = x_i, X_{-i} = x_{\delta(-i)})} \\ &= \frac{\sum_{\delta \in S^{n-1}} \Pr(\delta) \Pr(W_i = w_i, S_i = s_i, X_i = x_i, X_{-i} = x_{\delta(-i)})}{\sum_{\delta \in S^{n-1}} \Pr(\delta) \Pr(S_i = s_i, X_i = x_i, X_{-i} = x_{\delta(-i)})}. \end{aligned}$$

Because of the symmetry assumption, we know that

$$\begin{aligned} & \Pr(W_i = w_i, S_i = s_i, X_i = x_i, X_{-i} = x_{\delta(-i)}) \\ &= \int f(W_i = w_i, W_{-i} = w'_{-i}, S_i = s_i, S_{-i} = s'_{-i}, X_i = x_i, X_{-i} = x_{\delta(-i)}) ds'_{-i} dw'_{-i} \\ &= \int f(W_i = w_i, W_{-i} = w'_{\delta^{-1}(-i)}, S_i = s_i, S_{-i} = s'_{\delta^{-1}(-i)}, X_i = x_i, X_{-i} = x_{-i}) ds'_{-i} dw'_{-i} \\ &= \int f(W_i = w_i, W_{-i} = w'_{\delta^{-1}(-i)}, S_i = s_i, S_{-i} = s'_{\delta^{-1}(-i)}, X_i = x_i, X_{-i} = x_{-i}) ds'_{\delta^{-1}(-i)} dw'_{\delta^{-1}(-i)} \\ &= \Pr(W_i = w_i, S_i = s_i, X_i = x_i, X_{-i} = x_{-i}). \end{aligned}$$

For the same reason, we also have

$$\Pr(S_i = s_i, X_i = x_i, X_{-i} = x_{\delta(-i)}) = \Pr(S_i = s_i, X_i = x_i, X_{-i} = x_{-i}).$$

Thus, the posterior of consumer i can be simplified as:

$$\begin{aligned} & f_i(W_i = w_i | S_i = s_i, X_i = x_i, X_{-i}^* = x_{-i}) \\ &= \frac{\sum_{\delta \in S^{n-1}} \Pr(\delta) \Pr(\delta, W_i = w_i, S_i = s_i, X_i = x_i, X_{-i} = x_{-i})}{\sum_{\delta \in S^{n-1}} \Pr(\delta) \Pr(S_i = s_i, X_i = x_i, X_{-i} = x_{-i})} \\ &= \frac{\sum_{\delta \in S^{n-1}} \frac{1}{|S^{n-1}|} \Pr(\delta, W_i = w_i, S_i = s_i, X_i = x_i, X_{-i} = x_{-i})}{\sum_{\delta \in S^{n-1}} \frac{1}{|S^{n-1}|} \Pr(S_i = s_i, X_i = x_i, X_{-i} = x_{-i})} \\ &= f_i(W_i = w_i | S_i = s_i, X_i = x_i, X_{-i} = x_{-i}). \end{aligned}$$

We have therefore proved that consumer i has the same posterior about her willingness to pay w_i for any realization of W, S, X irrespective of whether the data are anonymized or not. Furthermore, this holds both on and off the path of play.

Next, we show that the producer also has the same posterior about W_i for any realization of W, S, X when consumer i refuses to report. Under the non-anonymized scheme, the posterior density is given by:

$$f_i(W_i = w_i | X = x) = \frac{\int f(W_i = w_i, W_{-i} = w'_{-i}, S = s', X_i = x_i, X_{-i} = x_{-i}) ds' dw'_{-i}}{\int f(W = w', S = s', X = x_i, X_{-i} = x_{-i}) ds' dw'}.$$

Under the anonymized scheme, the posterior density is given by

$$f_i(W_i = w_i | X^* = x) = \frac{\sum_{\delta \in S^{n-1}} \Pr(\delta) \Pr(W_i = w_i, X = \delta(x))}{\sum_{\delta \in S^{n-1}} \Pr(\delta) \Pr(X = \delta(x))}$$

By the earlier argument, we can simplify it as follows:

$$\frac{\sum_{\delta \in S^{n-1}} \Pr(\delta) \Pr(W_i = w_i, X = x)}{\sum_{\delta \in S^{n-1}} \Pr(\delta) \Pr(X = x)} = f_i(W_i = w_i | X = x)$$

Because the posteriors for both parties are the same for any realization, so is the price, and hence the welfare impact of information

The profit of the intermediary from consumer i 's data under inflow policy X is given by

$$R_i(X) = \Pi(X, X) - \Pi(S_i, \emptyset) - U_i((S_i, X_{-i}), X_{-i}) + U_i((S_i, X), X).$$

We have argued that consumer surplus off the path is the same:

$$U_i((S_i, X_{-i}), X_{-i}) = U_i((S_i, X_{-i}^*), X_{-i}^*).$$

We now turn to the last term—the impact on social welfare on the path of play:

$$\begin{aligned} & \Pi((S_i, X), X) + U_i((S_i, X), X) \\ &= \frac{1}{2} \mathbb{E} \left[(\mathbb{E}[w_i | S_i, X_i, X] - \mathbb{E}[w_i | X])^2 + \frac{1}{4} (\mathbb{E}[w_i | X])^2 \right] + \frac{\text{var}[\mathbb{E}[w_i | X]]}{4} \\ &= \frac{1}{2} \text{var}[\mathbb{E}[w_i | S_i, X_i, X]] - \frac{1}{8} \text{var}[\mathbb{E}[w_i | X]]. \end{aligned}$$

Recall that consumer i has the same on path posterior under two different scheme. Therefore, the difference in the intermediary's profits under the two policies reduces to

$$\begin{aligned} & \frac{1}{2} \text{var}[\mathbb{E}[w_i | S_i, X_i, X]] - \frac{1}{8} \text{var}[\mathbb{E}[w_i | X]] - \frac{1}{2} \text{var}[\mathbb{E}[w_i | S_i, X_i, X^*]] + \frac{1}{8} \text{var}[\mathbb{E}[w_i | X^*]] \\ &= -\frac{1}{8} \text{var}[\mathbb{E}[w_i | X]] + \frac{1}{8} \text{var}[\mathbb{E}[w_i | X^*]] \leq 0. \end{aligned}$$

Therefore, anonymization is more profitable than complete sharing, and strictly so whenever anonymization makes the estimation less precise. \square

In the remainder of the Appendix, we often make use of the following classical result in statistics, which we state as a lemma without proof—the result is an immediate consequence of the fact that $\mathbb{E}[X|Y]$ is the projection of X on $\mathcal{F}(Y)$ in L^2 space.

Lemma A1. Let W and Y be two random variables. Then it holds that

$$\text{var}[\mathbb{E}[W|Y]] = \text{var}[W] - \mathbb{E}[(W - \mathbb{E}[W|Y])^2] \leq \text{var}[W],$$

and

$$\mathbb{E}[(W - \mathbb{E}[W|Y])^2] \leq \mathbb{E}[(W - f(Y))^2], \quad \forall f \in L^2.$$

To prove Proposition 4, we first state a basic property of anonymized data sharing in our symmetric environment.

Lemma A2. When the data is anonymized, the following holds:

$$\mathbb{E}[w_i | A] = \mathbb{E}[w_j | A].$$

Proof of Lemma A2. Denote the joint distribution of W and S as $f(W = w, S = s)$ and the posterior of W_i after observing A as $f(W_i = w_i | A)$. Denote the permutation in S^N as v and especially the swapping between i and j as v_{ij} . For notational simplicity, we use \Pr to denote both probability and the proper marginal density.

$$\begin{aligned} f_i(W_i = w_i | A = s) &= \frac{\Pr(W_i = w_i, A = s)}{\Pr(A = s)} = \frac{\sum_{v \in S^N} \Pr(v) \Pr(W_i = w_i, S_v = s)}{\Pr(A = s)} \\ &= \frac{\sum_{v \in S^N} \frac{1}{|S^N|} \int f(W_i = w_i, W_j = w_j, W_{-ij} = w_{-ij}, S_v = s) dw_j dw_{-ij}}{\Pr(A = s)}. \end{aligned}$$

Because f is unchanged under permutation, we can apply the following transformation:

$$\begin{aligned} f_i(W_i = w_i | A = s) &= \frac{\sum_{v \in S^N} \frac{1}{|S^N|} \int f(W_j = w_i, W_i = w_j, W_{-ij} = w_{-ij}, S_{v_{ij} \circ v} = s) dw_j dw_{-ij}}{\Pr(A = s)}, \\ &= \frac{\sum_{v_{ij} \circ v \in S^N} \frac{1}{|S^N|} \int f(W_j = w_i, W_i = w'_i, W_{-ij} = w_{-ij}, S_{v_{ij} \circ v} = s) dw'_i dw_{-ij}}{\Pr(A = s)} = f_j(W_j = w_i | A = s). \end{aligned}$$

Because the posterior distribution is the same, so is the conditional expectation because

$$\mathbb{E}[w_i | A] = \int w_i f_i(W_i = w_i | A) dw_i,$$

which completes the proof. \square

Proof of Proposition 4. Combining Lemmas A1 and A2, we obtain

$$\mathbb{E}[w_i | A] = \mathbb{E} \left[\frac{1}{N} \sum_i w_i | A \right] = \mathbb{E} \left[\theta + \frac{1}{N} \sum_i \theta_i | A \right];$$

$$G(A) = \text{var} \left[\mathbb{E} \left[\theta + \frac{1}{N} \sum_i \theta_i | A \right] \right] = \text{var} \left[\theta + \frac{1}{N} \sum_i \theta_i \right] - \mathbb{E} \left[\left(\theta + \frac{1}{N} \sum_i \theta_i - \mathbb{E} \left[\theta + \frac{1}{N} \sum_i \theta_i | A \right] \right)^2 \right].$$

We can simplify the last term as follows:

$$\begin{aligned} & \mathbb{E} \left[\left(\theta + \frac{1}{N} \sum_i \theta_i - \mathbb{E} \left[\theta + \frac{1}{N} \sum_i \theta_i | A \right] \right)^2 \right] \\ &= \mathbb{E} \left[(\theta - \mathbb{E}[\theta | A])^2 + \frac{1}{N^2} (\sum_i \theta_i - \sum_i \mathbb{E}[\theta_i | A])^2 - \frac{2}{N} (\theta - \mathbb{E}[\theta | A]) (\sum_i \theta_i - \sum_i \mathbb{E}[\theta_i | A]) \right] \\ &\geq \mathbb{E}[(\theta - \mathbb{E}[\theta | A])^2] - \frac{2}{N} \sqrt{\text{var}[\theta - \mathbb{E}[\theta | A]] \text{var}[\sum_i \theta_i - \sum_i \mathbb{E}[\theta_i | A]]} \\ &\geq \mathbb{E}[(\theta - \mathbb{E}[\theta | A])^2] - \frac{2}{N} \sqrt{N \text{var}[\theta] \text{var}[\theta_i]}, \end{aligned}$$

where the last inequality comes from Lemma A1. The intermediary's profit can be written as

$$\begin{aligned} R &= 3G(A_{-i}) - G(A), \\ &= 3\text{var}[\mathbb{E}[\theta | A_{-i}]] - \text{var}[\theta] - \frac{1}{N} \text{var}[\theta_i] + \mathbb{E} \left[\left(\theta + \frac{1}{N} \sum_i \theta_i - \mathbb{E} \left[\theta + \frac{1}{N} \sum_i \theta_i | A \right] \right)^2 \right], \\ &\geq 3\text{var}[\mathbb{E}[\theta | A_{-i}]] - \text{var}[\theta] - \frac{1}{N} \text{var}[\theta_i] + \mathbb{E}[(\theta - \mathbb{E}[\theta | A])^2] - \frac{2}{N} \sqrt{N \text{var}[\theta] \text{var}[\theta_i]} \\ &= 3\text{var}[\mathbb{E}[\theta | A_{-i}]] - \text{var}[\mathbb{E}[\theta | A]] - \frac{1}{N} \text{var}[\theta_i] - \frac{2}{\sqrt{N}} \sqrt{\text{var}[\theta] \text{var}[\theta_i]}. \end{aligned}$$

Therefore, in the limit we have:

$$\lim_{N \rightarrow \infty} R = 2 \lim_{N \rightarrow \infty} \text{var}[\mathbb{E}[\theta | A]] > 0,$$

which completes the proof. \square

Proof of Proposition 5. We first prove that the total compensation is bounded from above, which immediately implies that the individual compensation goes to 0 as $N \rightarrow \infty$. From Lemma A2, we know that

$$\begin{aligned} G(A) &= \text{var}[\mathbb{E}[w_i | A]] = \text{var} \left[\mathbb{E} \left[\sum_i \frac{w_i}{N} | A \right] \right], \\ &\leq \text{var} \left[\sum_i \frac{w_i}{N} \right] = \text{var}[\theta] + \frac{\text{var}[\theta_i] + \text{var}[\varepsilon_i]}{N}. \end{aligned}$$

On the other hand, we also know

$$G(A_{-i}) = \text{var}[\mathbb{E}[\theta | A_{-i}]] = \text{var}[\theta] - \mathbb{E}[(\theta - \mathbb{E}[\theta | A_{-i}])^2].$$

Because the conditional expectation is the best L^2 approximation, we know it leads to a smaller error than the “sample average estimator,”

$$\mathbb{E}[(\theta - \mathbb{E}[\theta | A_{-i}])^2] \leq \mathbb{E} \left[\theta - \frac{1}{N-1} \sum_{j \neq i} (\theta + \theta_j + \varepsilon_j)^2 \right] = \frac{1}{N-1} (\text{var}[\theta_i] + \text{var}[\varepsilon_j]).$$

Therefore, we have:

$$\begin{aligned} N(G(A) - G(A_{-i})) &\leq N \left(\text{var}[\theta] + \frac{\text{var}[\theta_i] + \text{var}[\varepsilon_i]}{N} - \text{var}[\theta] + \frac{1}{N-1} (\text{var}[\theta_i] + \text{var}[\varepsilon_j]) \right), \\ &= N \left(\frac{\text{var}[\theta_i] + \text{var}[\varepsilon_i]}{N} + \frac{1}{N-1} (\text{var}[\theta_i] + \text{var}[\varepsilon_i]) \right) \\ &\leq 3(\text{var}[\theta_i] + \text{var}[\varepsilon_i]). \end{aligned}$$

The total consumer compensation is then given by

$$\frac{3N}{8} (G(A) - G(A_{-i})) \leq \frac{9}{8} (\text{var}[\theta_i] + \text{var}[\varepsilon_i]).$$

Finally, the intermediary's profit is growing linearly in N because

$$R(S) = \frac{N}{4} G(A) - \frac{3N}{8} (G(A) - G(A_{-i})),$$

$$\lim_{N \rightarrow \infty} \frac{R(S)}{N} = \frac{1}{4} \lim_{N \rightarrow \infty} G(A),$$

which completes the proof. \square

Proof of Proposition 6. When data is not anonymized we have:

$$G(S) - G(S_{-i}) = \text{var}[\mathbb{E}[\theta + \theta_i | S]] - \text{var}[\mathbb{E}[\theta | S_{-i}]].$$

Because of symmetry, we have

$$\text{cov}[\mathbb{E}[\theta | S], \mathbb{E}[\theta_i | S]] = \text{cov}[\mathbb{E}[\theta | S], \mathbb{E}[\theta_j | S]] = \text{cov}[\mathbb{E}[\theta | S], \Sigma_{j=1}^N \mathbb{E}[\theta_j / N | S]].$$

Because the correlation coefficient is always greater than -1 , we obtain

$$\begin{aligned} \text{cov}[\mathbb{E}[\theta | S], \Sigma_{j=1}^N \mathbb{E}[\theta_j / N | S]] &\geq -\sqrt{\text{var}[\theta] \text{var}[\Sigma_{j=1}^N \mathbb{E}[\theta_j / N | S]]}, \\ &\geq -\sqrt{\text{var}[\theta] \text{var}[\Sigma_{j=1}^N \theta_j / N]}. \end{aligned}$$

Therefore, according to Lemma A1 we have:

$$\begin{aligned} G(S) - G(S_{-i}) &= \text{var}[\mathbb{E}[\theta | S]] + 2\text{cov}[\mathbb{E}[\theta | S], \mathbb{E}[\theta_i | S]] + \text{var}[\mathbb{E}[\theta_i | S]] - \text{var}[\mathbb{E}[\theta | S_{-i}]] \\ &\geq \text{var}[\mathbb{E}[\theta | S]] - 2\frac{1}{\sqrt{N}}\sqrt{\text{var}[\theta] \text{var}[\theta_i]} + \text{var}[\mathbb{E}[\theta_i | S]] - \text{var}[\mathbb{E}[\theta | S_{-i}]], \end{aligned}$$

and hence

$$\liminf_{N \rightarrow \infty} G(S) - G(S_{-i}) \geq \text{var}[\mathbb{E}[\theta_i | S]].$$

The last term is strictly positive because

$$\begin{aligned} \text{var}[\mathbb{E}[\theta_i | S]] &= \text{var}[\theta_i] - \mathbb{E}[(\theta_i - \mathbb{E}[\theta_i | S])^2] \\ &\geq \text{var}[\theta_i] - \mathbb{E}\left[\left(\theta_i - \frac{\text{var}[\theta_i]}{\text{var}[\theta_i] + \text{var}[\theta] + \text{var}[e]} S_i\right)^2\right], \\ &= \text{var}[\theta_i] - \left(\text{var}[\theta_i] - \frac{\text{var}^2[\theta_i]}{\text{var}[\theta_i] + \text{var}[\theta] + \text{var}[e]}\right), \\ &= \frac{\text{var}^2[\theta_i]}{\text{var}[\theta_i] + \text{var}[\theta] + \text{var}[e]} > 0, \end{aligned}$$

where the first inequality again uses Lemma A1. \square

Proof of Proposition 7. In the standard “divide and conquer” scheme, the compensation for the i -th consumer is the marginal loss of revealing her information given that $i - 1$ consumers reveal their signals:

$$\frac{3}{8} G(S_{1, \dots, i}) - \frac{3}{8} G(S_{1, \dots, i-1}).$$

Because in general we do not know whether this marginal loss is decreasing in i , we consider the following revised version of divide and conquer, where consumer i receives

$$m_i = \max_{k \geq i} \frac{3}{8} G(S_{1, \dots, k}) - \frac{3}{8} G(S_{1, \dots, k-1}).$$

Under this payment scheme, it is a dominant strategy for consumer 1 to accept the offer. Moreover, it is optimal for consumer i to accept the offer, given that the first $i - 1$ consumers accept. Using an identical proof to Proposition 5, we obtain

$$\begin{aligned} \frac{3}{8} G(S_{1, \dots, i}) - \frac{3}{8} G(S_{1, \dots, i-1}) &\leq \frac{3}{8} \left(\frac{1}{i} + \frac{1}{i-1} \right) (\text{var}[\theta_i] + \text{var}[\varepsilon_i]), \\ &\leq \frac{3}{4} \frac{1}{i-1} (\text{var}[\theta_i] + \text{var}[\varepsilon_i]). \end{aligned}$$

Therefore, we obtain an upper bound on the compensation paid to consumer i :

$$m_i \leq \max_{k \geq i} \frac{3}{4} \frac{1}{k-1} (\text{var}[\theta_i] + \text{var}[\varepsilon_i]) = \frac{3}{4} \frac{1}{i-1} (\text{var}[\theta_i] + \text{var}[\varepsilon_i]).$$

Finally, because we have

$$\begin{aligned}\Sigma_i \frac{3}{8} (G(S_{1,\dots,i}) - G(S_{1,\dots,i-1})) &\leq \frac{3}{4} \left(1 + \sum_{i=3}^N \frac{1}{i-1} \right) (\text{var}[\theta_i] + \text{var}[\varepsilon_i]), \\ &\leq \frac{3}{4} (1 + \log N) (\text{var}[\theta_i] + \text{var}[\varepsilon_i]),\end{aligned}$$

the total compensation grows at a speed less than $\log N$. \square

Proof of Proposition 8. The proof of this proposition is similar to that of Proposition 3. By Lemma 1, we know that the intermediary will transmit whatever information it collected to all consumers. By homogeneity, we know the consumer's posterior about their own fundamental w_{ij} is the same whether the signals are anonymized or not, and the producer's posterior about any deviating consumer's fundamental is also the same under the two schemes.

Denote the broker's revenue under the non-anonymized and anonymized scheme as $R(X)$ and $R(X^*)$, respectively. It then holds that

$$\begin{aligned}R_i(X) &= \Pi(X, X) - \Pi(S_i, \emptyset) - U_i((S_i, X_{-i}), X_{-i}) + U_i((S_i, X), X), \\ R_i(X^*) &= \Pi(X^*, X^*) - \Pi(S_i, \emptyset) - U_i((S_i, X_{-i}^*), X_{-i}^*) + U_i((S_i, X^*), X^*).\end{aligned}$$

Our analysis in previous paragraph implies

$$U_i((S_i, X_{-i}), X_{-i}) = U_i((S_i, X_{-i}^*), X_{-i}^*).$$

Therefore the intermediary prefers anonymization if and only if

$$R_i(X^*) - R_i(X) = W((S_i, X^*), X^*) - W((S_i, X), X) \geq 0,$$

which completes the proof. \square

Proof of Proposition 9. We first consider the case where the intermediary anonymizes all data, including the group identities. Similar to the result in Lemma A2, we know that the producer offers one price to all consumers on the path of play,

$$\mathbb{E}[w_{ij}|A] = \mathbb{E}[w_{i'j'}|A].$$

Denoting $N = \Sigma_j N_j$, we have

$$\mathbb{E}[w_{i'j'}|A] = \frac{1}{\Sigma_j N_j} \Sigma_j \Sigma_i \mathbb{E}[w_{ij}|A] = \frac{1}{\Sigma_j N_j} \Sigma_j \Sigma_i \mathbb{E}[\theta_j + \theta_{ij}|A]$$

Therefore we obtain an upper bound on the revenue per capita

$$\begin{aligned}\frac{R(A)}{N} &= \frac{3}{8} G(A_{-ij}) - \frac{1}{8} G(A) = \frac{3}{8} \text{var}[\mathbb{E}[w_{ij}|A_{-ij}]] - \frac{1}{8} \text{var}[\mathbb{E}[w_{ij}|A]] \\ &\leq \frac{1}{4} \text{var}[\mathbb{E}[w_{ij}|A]] \leq \frac{1}{4} \frac{1}{N^2} \Sigma N_j^2 \text{var}[\theta_j] + \frac{1}{4N} \text{var}[\theta_{ij}].\end{aligned}$$

Next, consider the case where the intermediary reveals the group identity. Instead of A we use A^g to denote the information that the producer receives. By an argument similar to the proof of Lemma A2, we know that (on path) the producer offers one price to all consumers in each group:

$$\begin{aligned}\mathbb{E}[w_{ij}|A^g] &= \mathbb{E}[w_{i'j'}|A^g] \\ &= \frac{1}{N_j} \Sigma_{i'=1}^{N_j} \mathbb{E}[w_{i'j'}|A^g] = \frac{1}{N_j} \Sigma_{i'=1}^{N_j} \mathbb{E}[w_{i'j'}|A^g] = \mathbb{E}\left[\theta_j + \frac{1}{N_j} \Sigma_{i'=1}^{N_j} \theta_{i'j} \middle| A^g\right].\end{aligned}$$

When consumer ij rejects the offer, the intermediary will know the group identity of this deviating consumer and use all the available data to estimate the demand:

$$\mathbb{E}[w_{ij}|A_{-ij}^g] = \mathbb{E}[\theta_j + \theta_{ij}|A_{-ij}^g] = \mathbb{E}[\theta_j|A_{-ij}^g].$$

The revenue that the intermediary obtains from consumer ij 's data is then given by

$$\begin{aligned}&\frac{3}{8} \text{var}[\mathbb{E}[w_{ij}|A_{-ij}^g]] - \frac{1}{8} \text{var}[\mathbb{E}[w_{ij}|A^g]], \\ &= \frac{3}{8} \text{var}[\mathbb{E}[\theta_j|A_{-ij}^g]] - \frac{1}{8} \text{var}\left[\mathbb{E}\left[\theta_j + \frac{1}{N_j} \Sigma_{i'=1}^{N_j} \theta_{i'j} \middle| A^g\right]\right],\end{aligned}$$

$$\geq \frac{3}{8} \text{var}[\mathbb{E}[\theta_j | A_{-ij}^g]] - \frac{1}{8} \text{var}[\theta_j] - \frac{1}{8N_j} \text{var}[\theta_{ij}] - \frac{1}{4} \sqrt{\frac{1}{N_j} \text{var}[\theta_j] \text{var}[\theta_{ij}]}.$$

From Lemma A1, we know

$$\begin{aligned} \mathbb{E}[(\theta_j - \mathbb{E}[\theta_j | A_{-ij}^g])^2] &\leq \mathbb{E}\left[\left(\theta_j - \frac{1}{N_j - 1} \sum_{\ell \neq i, \ell \neq j} \theta_\ell\right)^2\right], \\ &= \frac{1}{N_j - 1} (\text{var}[\theta_{ij}] + \text{var}[\varepsilon_{ij}]); \\ \text{var}[\mathbb{E}[\theta_j | A_{-ij}^g]] &= \text{var}[\theta_j] - \mathbb{E}[(\theta_j - \mathbb{E}[\theta_j | A_{-ij}^g])^2], \\ &\geq \text{var}[\theta_j] - \frac{1}{N_j - 1} (\text{var}[\theta_{ij}] + \text{var}[\varepsilon_{ij}]). \end{aligned}$$

Thus we obtain a lower bound on the revenue from consumer ij :

$$\frac{1}{4} \text{var}[\theta_j] - \frac{3}{8} \frac{1}{N_j - 1} (\text{var}[\theta_{ij}] + \text{var}[\varepsilon_{ij}]) - \frac{1}{8N_j} \text{var}[\theta_{ij}] - \frac{1}{4} \sqrt{\frac{1}{N_j} \text{var}[\theta_j] \text{var}[\theta_{ij}]}.$$

Finally, we can compute the difference in the revenues

$$\begin{aligned} R(A) - R(A^g) &\leq \sum_j \frac{N_j^2}{4N} \text{var}[\theta_j] + \frac{1}{4} \text{var}[\theta_{ij}] \\ &\quad - \sum_j \left(\frac{N_j}{4} \text{var}[\theta_j] - \frac{3}{8} \frac{N_j}{N_j - 1} (\text{var}[\theta_{ij}] + \text{var}[\varepsilon_{ij}]) - \frac{1}{8} \text{var}[\theta_{ij}] - \frac{\sqrt{N_j}}{4} \sqrt{\text{var}[\theta_j] \text{var}[\theta_{ij}]} \right). \end{aligned}$$

As long as $N_j < kN$ where $k < 1$, we know that

$$\begin{aligned} R(A) - R(A^g) &< \frac{1}{4} \text{var}[\theta_{ij}] + \sum_j \left(-\frac{1-k}{4} N_j \text{var}[\theta_j] + \frac{3}{8} \frac{N_j}{N_j - 1} (\text{var}[\theta_{ij}] + \text{var}[\varepsilon_{ij}]) + \frac{1}{8} \text{var}[\theta_{ij}] + \frac{\sqrt{N_j}}{4} \sqrt{\text{var}[\theta_j] \text{var}[\theta_{ij}]} \right). \end{aligned}$$

The dominant linear term is decreasing in N_j , and hence we know that as $N_j \rightarrow \infty$, revealing group identities is more profitable. \square

Proof of Proposition 10. Each consumer's demand function is given by

$$q_i = w_i - (\ell_i - x_i)^2 - p_i.$$

This means the producer's profit is given by

$$\pi = \sum_{i=1}^N p_i (w_i - (\ell_i - x_i)^2 - p_i).$$

Therefore, under any information structure S , the producer offers

$$\begin{aligned} p_i &= (\mathbb{E}[w_i | S] - \mathbb{E}[(\ell_i - \mathbb{E}[\ell_i | S])^2 | S]) / 2, \\ &= (\mathbb{E}[w_i | S] - \mathbb{E}[(\ell_i - \mathbb{E}[\ell_i | S])^2]) / 2, \\ x_i &= \mathbb{E}[\ell_i | S], \end{aligned}$$

where the second line relies on the fact that the underlying random variables are normal so that $\ell_i - \mathbb{E}[\ell_i | S]$ is independent of S .

The consumer's surplus is then given by

$$\begin{aligned} U_i(S) &= \frac{1}{2} \mathbb{E} \left[\left(w_i - (\ell_i - \mathbb{E}[\ell_i | S])^2 - \frac{\mathbb{E}[w_i | S] - \mathbb{E}[(\ell_i - \mathbb{E}[\ell_i | S])^2]}{2} \right)^2 \right] \\ &= \frac{1}{2} \mathbb{E} \left[\left(w_i - \frac{1}{2} \mathbb{E}[w_i | S] \right)^2 \right] + \frac{1}{2} \mathbb{E} \left[\left((\ell_i - \mathbb{E}[\ell_i | S])^2 - \frac{1}{2} \mathbb{E}[(\ell_i - \mathbb{E}[\ell_i | S])^2] \right)^2 \right] \\ &\quad - \mathbb{E} \left[\left(w_i - \frac{1}{2} \mathbb{E}[w_i | S] \right) \right] \mathbb{E} \left[(\ell_i - \mathbb{E}[\ell_i | S])^2 - \frac{1}{2} \mathbb{E}[(\ell_i - \mathbb{E}[\ell_i | S])^2] \right] \end{aligned}$$

$$\begin{aligned}
&= \frac{1}{2} \mathbb{E} \left[w_i^2 - \frac{3}{4} \mathbb{E}[w_i | S]^2 \right] + \frac{1}{2} \mathbb{E} \left[(\ell_i - \mathbb{E}[\ell_i | S])^4 - \frac{3}{4} \mathbb{E}[(\ell_i - \mathbb{E}[\ell_i | S])^2]^2 \right] \\
&\quad - \frac{1}{4} \mathbb{E}[w_i] \mathbb{E}[(\ell_i - \mathbb{E}[\ell_i | S])^2].
\end{aligned}$$

Therefore, the difference is:

$$\begin{aligned}
U_i(S) - U_i(\emptyset) &= -\frac{3}{8} \text{var}[\mathbb{E}[w_i | S]] + \frac{1}{4} \mu \text{var}[\mathbb{E}[\ell_i | S]] + \frac{1}{2} \mathbb{E}[(\ell_i - \mathbb{E}[\ell_i | S])^4] \\
&\quad - \frac{3}{8} \mathbb{E}[(\ell_i - \mathbb{E}[\ell_i | S])^2]^2 - \frac{1}{2} \mathbb{E}[(\ell_i - \mu_\tau)^4] + \frac{3}{8} \mathbb{E}[(\ell_i - \mu_\tau)^2]^2.
\end{aligned}$$

Because every random variable is assumed to be normal, $\ell_i - \mathbb{E}[\ell_i | S]$ is also normal with zero mean. We can further simplify and obtain

$$\begin{aligned}
U_i(S) - U_i(\emptyset) &= -\frac{3}{8} \text{var}[\mathbb{E}[w_i | S]] + \frac{1}{4} \mu \text{var}[\mathbb{E}[\ell_i | S]] + \frac{3}{2} (\text{var}[\ell_i] - \text{var}[\mathbb{E}[\ell_i | S]])^2 \\
&\quad - \frac{3}{8} (\text{var}[\ell_i] - \text{var}[\mathbb{E}[\ell_i | S]])^2 - \frac{3}{2} \text{var}[\ell_i]^2 + \frac{3}{8} \text{var}[\ell_i]^2, \\
&= -\frac{3}{8} \text{var}[\mathbb{E}[w_i | S]] + \frac{1}{4} \mu \text{var}[\mathbb{E}[\ell_i | S]] + \frac{9}{8} (\text{var}[\mathbb{E}[\ell_i | S]]^2 - 2 \text{var}[\mathbb{E}[\ell_i | S]] (\sigma_\tau^2 + \sigma_{\epsilon_i}^2)).
\end{aligned}$$

Similarly we have:

$$\begin{aligned}
\Pi_i(S) &= \frac{1}{4} \mathbb{E}[(\mathbb{E}[w_i | S] - \mathbb{E}[(\ell_i - \mathbb{E}[\ell_i | S])^2])] \\
&= \frac{1}{4} \mathbb{E}[\mathbb{E}[w_i | S]^2 - 2 \mathbb{E}[w_i | S] \mathbb{E}[(\ell_i - \mathbb{E}[\ell_i | S])^2] + \mathbb{E}[(\ell_i - \mathbb{E}[\ell_i | S])^2]^2],
\end{aligned}$$

and hence

$$\begin{aligned}
\Pi_i(S) - \Pi_i(\emptyset) &= \frac{1}{4} \text{var}[\mathbb{E}[w_i | S]] + \frac{1}{2} \mu \text{var}[\mathbb{E}[\ell_i | S]] \\
&\quad + \frac{1}{4} (\text{var}[\mathbb{E}[\ell_i | S]]^2 - 2 \text{var}[\mathbb{E}[\ell_i | S]] (\sigma_\tau^2 + \sigma_{\epsilon_i}^2)).
\end{aligned}$$

To summarize, the impact of data sharing on social surplus is given by

$$\begin{aligned}
W_i(S) - W_i(\emptyset) &= U_i(S) - U_i(\emptyset) + \Pi_i(S) - \Pi_i(\emptyset), \\
&= -\frac{1}{8} \text{var}[\mathbb{E}[w_i | S]] + \frac{3}{4} \mu \text{var}[\mathbb{E}[\ell_i | S]] + \frac{11}{8} (\text{var}[\mathbb{E}[\ell_i | S]]^2 - 2 \text{var}[\mathbb{E}[\ell_i | S]] (\sigma_\tau^2 + \sigma_{\epsilon_i}^2)).
\end{aligned}$$

Therefore the difference $W_i(S) - W_i(\emptyset)$ is a quadratic function of the variance of the conditional expectation $x \triangleq \text{var}[\mathbb{E}[\ell_i | S]]$. In particular, we let

$$g(x) \triangleq \frac{11}{8} x^2 + \left(\frac{3}{4} \mu - \frac{11}{4} (\sigma_\tau^2 + \sigma_{\epsilon_i}^2) \right) x.$$

As long as $3\mu > 11(\sigma_\tau^2 + \sigma_{\epsilon_i}^2)$, this function is positive and increasing in x , which means a higher $\text{var}[\mathbb{E}[\ell_i | S]]$ increases consumer surplus.

Finally, as in the proof of Proposition 3, aggregating w_i increases $W_i(S)$ but keeps $\Pi(\emptyset)$ and $U_i(S_{-i})$ unchanged. Not aggregating ℓ_i increases $W_i(S)$ while keeping $\Pi(\emptyset)$ and $U_i(S_{-i})$ unchanged. Therefore, it is optimal for the intermediary to aggregate w_i but not ℓ_i . \square

References

- ACEMOGLU, D., MAKHDOUNI, A., MALEKIAN, A., and OZDAGLAR, A. “Too Much Data: Prices and Inefficiencies in Data Markets.” *American Economic Journal: Microeconomics*, forthcoming.
- ACQUISTI, A., TAYLOR, C., and WAGMAN, L. “The Economics of Privacy.” *Journal of Economic Literature*, Vol. 54 (2016), pp. 442–492.
- ALI, S.N., LEWIS, G., and VASSERMAN, S. “Voluntary Disclosure and Personalized Pricing.” Technical Report 26592, National Bureau of Economic Research, 2019.
- ARIDOR, G., CHE, Y.K., and SALZ, T. “The Economic Consequences of Data Privacy Regulation: Empirical Evidence from GDPR.” Technical Report 26900, National Bureau of Economic Research, 2020.
- ARRIETA-IBARRA, I., GOFF, L., JIMENEZ-HERNANDEZ, D., LANIER, J., and WEYL, G. “Should We Treat Data as Labor? Moving beyond “Free”.” *American Economic Review Paper and Proceedings*, Vol. 108 (2018), pp. 38–42.

- ATHEY, S., CATALINI, C., and TUCKER, C. "The Digital Privacy Paradox: Small Money, Small Costs, Small Talk." Technical Report 23488, National Bureau of Economic Research, 2017.
- BERGEMANN, D. and BONATTI, A. "Targeting in Advertising Markets: Implications for Offline Vs. Online Media." *Rand Journal of Economics*, Vol. 42 (2011), pp. 417–443.
- . "Markets for Information: An Introduction." *Annual Review of Economics*, Vol. 11 (2019), pp. 85–107.
- BERGEMANN, D., BROOKS, B., and MORRIS, S. "The Limits of Price Discrimination." *American Economic Review*, Vol. 105 (2015), pp. 921–957.
- CHOI, J., JEON, D., and KIM, B. "Privacy and Personal Data Collection with Information Externalities." *Journal of Public Economics*, Vol. 173 (2019), pp. 113–124.
- Committee on the Judiciary. "Investigation of Competition in Digital Markets." Technical Report, United States House of Representatives, 2020.
- Competition & Markets Authority. "Online Platforms and Digital Advertising." Technical Report, UK Government, 2020.
- CREMÉR, J., DE MONTJOYE, Y.A., and SCHWEITZER, H. "Competition Policy for the Digital Era." Technical Report, European Commission, 2019.
- CUMMINGS, R., LIGETT, K., PAI, M., and ROTH, A. "The Strange Case of Privacy in Equilibrium Models." In ACM-EC (Economics and Computation) 2016. New York: Association for Computer Machinery, 2016.
- FAINMESSER, I., GALEOTTI, A., and MOMOT, R. "Digital Privacy." Technical Report, Johns Hopkins University, 2020.
- FURMAN, J., COYLE, D., FLETCHER, A., MCAULES, D., and MARSDEN, P. "Unlocking Digital Competition: Report of the Digital Competition Expert Panel." HM Treasury, United Kingdom.
- GRADWOHL, R. "Information Sharing and Privacy in Networks." In ACM-EC (Economics and Computation) 2017. New York: Association for Computer Machinery, 2017.
- ICHIHASHI, S. "The Economics of Data Externalities." Technical Report, Bank of Canada, 2020a.
- . "Online Privacy and Information Disclosure by Consumers." *American Economic Review*, Vol. 110 (2020b), pp. 569–595.
- JOHNSON, G., SHRIVER, S., and GOLDBERG, S. "Privacy & Market Concentration: Intended & Unintended Consequences of the GDPR." Technical Report, Boston University Questrom School of Business, 2020.
- JULLIEN, B., LEFOULI, Y., and RIORDAN, M. "Privacy Protection, Security, and Consumer Retention." Technical Report, Toulouse School of Economics, 2020.
- LIANG, A. and MADSEN, E. "Data and Incentives." Tech. rep., Northwestern University, 2020.
- LIN, T. "Valuing Intrinsic and Instrumental Preferences for Privacy." Technical Report, Boston University, 2019.
- LIZZERI, A. "Information Revelation and Certification Intermediaries." *RAND Journal of Economics*, Vol. 30 (1999), pp. 214–231.
- LOERTSCHER, S. and MARX, L.M. "Digital Monopolies: Privacy Protection or Price Regulation?" *International Journal of Industrial Organization*, Vol. 71 (2020), pp. 1–13.
- MIKLOS-THAL, J. and SHAFFER, G. "Naked Exclusion with Private Offers." *American Economic Journal: Microeconomics*, Vol. 8 (2016), pp. 174–194.
- MONTIEL OLEA, J.L., ORTOLEVA, P., PAI, M., and PRAT, A. "Competing Models." *arXiv preprint arXiv:1907.03809*.
- POSNER, E.A. and WEYL, E.G. *Radical Markets: Uprooting Capitalism and Democracy for a Just Society*. Princeton University Press, 2018.
- ROBINSON, J. *The Economics of Imperfect Competition*. Macmillan, London, 1933.
- ROMER, P. "A Tax That Could Fix Big Tech." *The New York Times*.
- SCHMALENSEE, R. "Output and Welfare Implications of Monopolistic Third-Degree Price Discrimination." *American Economic Review*, Vol. 71 (1981), pp. 242–247.
- SEGAL, I. "Contracting with Externalities." *Quarterly Journal of Economics*, Vol. 114 (1999), pp. 337–388.
- SEGAL, I. and WHINSTON, M. "Naked Exclusion: Comment." *American Economic Review*, Vol. 90 (2000), pp. 296–309.
- Stigler Committee on Digital Platforms. "Final Report." Technical Report, Stigler Center for the Study of the Economy and the State, 2019.
- TANG, H. "The Value of Privacy: Evidence from Online Borrowers." Technical Report, HEC Paris, 2019.
- TAYLOR, C. "Consumer Privacy and the Market for Customer Information." *RAND Journal of Economics*, Vol. 35 (2004), pp. 631–651.
- ZUBOFF, S. *The Age of Surveillance Capitalism: The Fight for a Human Future at the New Frontier of Power*. Public Affairs, New York, 2019.

Supporting information

Additional supporting information may be found online in the Supporting Information section at the end of the article.