# Competitive Markets for Personal Data

Simone Galperti Tianhao Liu Jacopo Perego UC, San Diego Columbia University Columbia University

April 6, 2024

#### ABSTRACT

We study competitive data markets in which consumers own their personal data and can trade it with intermediaries, such as e-commerce platforms. Intermediaries use this data to provide services to the consumers, such as targeted offers from third-party merchants. Our main results identify a novel inefficiency, resulting in equilibrium data allocations that fail to maximize welfare. This inefficiency hinges on the role that intermediaries play as information gatekeepers, a hallmark of the digital economy. We provide three solutions to this market failure: establishing data unions, which manage consumers' data on their behalf; taxing the trade of data; and letting the price of data depend on its intended use.

**JEL Classification Numbers:** C72, D82, D83

**Keywords:** Consumer Data, Information Design, Data Externalities, Data Union.

We are thankful to Dirk Bergemann, Alessandro Bonatti, Laura Doval, Nima Haghpanah, Kevin He, Dana Foarta, Matthew Gentzkow, Alessandro Lizzeri, Alessandro Pavan, Marco Ottaviani, Paolo Siconolfi, Philipp Strack, Laura Veldkamp, Jidong Zhou, and seminar participants at numerous universities for helpful comments. This research is supported by grants from the NSF (Galperti: SES-2149289; Perego: SES-2149315). Galperti also gratefully acknowledges financial support from the UPenn's CTIC and the Warren Center for Network & Data Sciences. *Contacts:* Simone Galperti (sgalperti@ucsd.edu), Tianhao Liu (tl3014@columbia.edu), Jacopo Perego (corresponding author, jacopo.perego@columbia.edu).

## 1 Introduction

Consumer data has become a crucial input of the modern economy and contributes to the success of many large industries, such as online advertisement and digital marketplaces. Firms collect consumers' data to learn their tastes and offer targeted advertisements or personalized products and services. While consumers themselves are the primary suppliers of this data, they typically have limited control over who uses it and how, and only rarely are they financially compensated in return (Federal Trade Commission, 2014). Such an arrangement could lead to distortions, inefficiencies, and greater inequality (Bergemann et al., 2023). Accordingly, new legislation has been recently introduced around the world to give consumers more control over their data. This creates the legal framework upon which *data markets* can emerge, where consumers own their data and firms compete to acquire and use it. What properties would such markets have? And which institutions (if any) should be designed to ensure they promote efficient outcomes?

This paper contributes to answering these questions by studying a stylized model of competitive data markets. We present two sets of results. First, we identify a novel inefficiency that causes the equilibria of this economy to induce data allocations that fail to maximize welfare. This inefficiency stems from an externality that consumers exert on each other when selling their data to intermediaries, such as e-commerce platforms, which use this data to provide services to the consumers. Second, we discuss three solutions to this market failure: establishing data unions, taxing the trade of data, or letting the price of data depend on its intended use.

Our model features a heterogeneous population of consumers, an e-commerce platform, and a third-party merchant. Each consumer owns her data and can sell it to the platform. When this happens, the platform gains access to this consumer and learns her tastes (i.e., her "type"). In exchange, the platform pays this consumer the market price for her data and, in addition, offers her a service. The service consists of intermediating this consumer with the merchant, from whom she can buy a product. As an intermediary, the platform uses its database of consumers' data to provide information to the merchant about their willingness to pay for the product, in the spirit of Bergemann et al. (2015). The main innovation of our model is that the platform's

<sup>&</sup>lt;sup>1</sup>Most notably, the European Union's General Data Protection Regulation (GDPR) grants consumers the right to object to how firms use their data, to request it to be transferred to other firms, or to be deleted. In the United States, a growing list of states have passed bills with a similar scope.

database is determined endogenously as an equilibrium of the markets in which the consumers and the platform trade the data. We assume that such markets are perfectly competitive: The consumers and the platform take data prices as given, and these are pinned down by market clearing. This assumption shuts down distortions that may stem from market power—which is not specific to data—and allows us to focus on novel distortions that are specific to data and can persist even in a competitive setting.

Our main goal is to study the properties of competitive data markets and their ability to promote efficient data allocations. We identify necessary and sufficient conditions for efficiency and show how they depend on the platform's objective. When this objective is sufficiently aligned with that of the merchant, data markets are efficient and consumers' welfare is maximized. By contrast—and perhaps counterintuitively—when the platform's objective is sufficiently aligned with that of the consumers, the equilibrium data allocation can be inefficient. In some cases, markets can entirely unravel, resulting in no data being traded, which leads to the lowest possible welfare. Finally, we sharpen this negative finding by identifying conditions under which *all* equilibria are inefficient.

Inefficiency results because consumers exert an externality on each other when selling their data to the platform. This externality arises endogenously from how the platform uses this data to pursue its objective. When it cares relatively more about creating a surplus for the consumers, it withholds some information from the merchant to curb surplus extraction. To do so, it pools consumers of different types to prevent the merchant from learning their willingness to pay. We show that such pooling introduces a wedge between a consumer's *private* benefit from selling her data and joining the pool and the *social* benefit that this consumer generates for the other consumers in the pool. This wedge leads consumers to make socially inefficient decisions regarding the sale of their data. Thus, the inefficiency originates from the platform's practice of withholding information to harness the conflicting objectives of the constituencies it intermediates.<sup>2</sup> Building on Galperti, Levkun, and Perego (2023), our results show that this practice not only creates interdependencies in how the platform values the data records but effectively sets the stage for market failure. Therefore, our inefficiency does not originate from

<sup>&</sup>lt;sup>2</sup>This practice is common in the marketplace. For instance, Google Search withholds information about users' characteristics from advertisers to increase competition, and ridesharing platforms similarly withhold information from drivers about riders' destinations to discourage selection.

exogenous correlation in consumers' data nor from the fact that the merchant is a monopolist.<sup>3</sup>

We then analyze three institutions that correct the aforementioned market failure. The first is a new intermediary, called a data union, which manages consumers' data on their behalf. Our data union collects data from participating consumers, sells some of the data to the platform, and distributes the proceeds back to the consumers as compensation. We show that a data union can act as a benevolent planner, helping consumers coordinate the decisions to sell their data and internalize the aforementioned externality. Indeed, any equilibrium of the economy with the data union is efficient and consumers' welfare is maximized, regardless of the platform's objective. This result offers theoretical support to recent policy proposals that discuss the potential benefits of data unions for the data economy (e.g., see Posner and Weyl, 2018; Bergemann et al., 2023).

The second solution we consider introduces data taxes that are levied when consumers sell their data. These taxes are "simple" in that they depend only on the type of data that is traded—not on the consumers' identity or how the data is used. When properly designed, data taxes make each consumer internalize the effects that selling her data creates on the rest of the economy. These taxes increase (decrease) the effective data price paid to consumers who exert a positive (negative) externality—typically, these are consumers whose underlying type is low (high). Consequently, we show that any efficient allocation can be implemented by an equilibrium of the competitive economy with budget-balanced data taxes.

The third solution consists in letting the market price of data depend not only on its type but also on how the platform intends to use it. In other words, there are different markets in which to trade the same type of data, depending on its intended use. This solution is inspired by classic approaches to competitive economies with externalities (Arrow, 1969; Laffont, 1976). It is also broadly in the spirit of legislation like the aforementioned GDPR, which requires that the intended use for which consumer data is collected should be determined at the time of its collection (see, GDPR 2016/679 (39)). Since the externalities originate from how the platform uses the data, this richer price system lets the market assign a value to these externalities and convey it to all market participants. We show that this guarantees equilibrium efficiency, regardless of the platform's objective.

<sup>&</sup>lt;sup>3</sup>Indeed, in our model consumers' data is uncorrelated. Moreover, it is known that a platform can find it optimal to withhold information even in the presence of competing merchants (see, e.g., Bergemann et al., 2022b; Elliott et al., 2022), .

**Related Literature**. Our approach is rooted in a general-equilibrium tradition, but leverages the recent progress of the information-design literature (for reviews of this literature, see, Bergemann and Morris, 2019; Kamenica, 2019). This allows us to offer a principled microfoundation of some of the key components of a data economy. For instance, how platforms *use* the data traces back to Bergemann et al. (2015), how platforms *value* the data builds on Galperti et al. (2023), and how competitive markets *price* the data—a novel component of this paper—continues this line of research.

Our paper contributes to a recent literature that studies data markets and their properties. Particularly close to our work is a set of papers that identifies a data externality that can lead to inefficient outcomes. In Choi et al. (2019), Ichihashi (2021), Acemoglu et al. (2022), and Bergemann et al. (2022a), a consumer's decision to sell her data can create externalities on other consumers because their data are correlated. Relative to these papers, we emphasize a different inefficiency, which instead arises from how platforms endogenously use data. To emphasize this difference, we assume throughout that consumers' data are uncorrelated; our platform learns nothing about a consumer when acquiring the data of another. Our inefficiency builds on previous work by Galperti et al. (2023), which characterizes how much a platform values the data record of each consumer in a database. This value has two components: The payoff a platform earns from the consumer associated with the record and the indirect payoff it earns from pooling that record with those of other consumers to achieve better outcomes. They refer to the second component as a "pooling" externality. With the present paper, we contribute to this agenda in several ways. First, in a model of competitive data markets, we endogenize consumers' decisions to sell their data records and participate in the platform's database. Because of the competitive nature of these markets, the price of data records equals the value the platform assigns to them, and thus inherently accounts for these "pooling" externalities. Nevertheless, we show that pooling misaligns the consumers' private and social benefits from selling their data records, which leads to inefficiencies. Finally, we present solutions that are designed to circumvent this market failure.

More broadly, our paper contributes to a growing literature on the economics of digitization, data, and privacy (for reviews, see Acquisti et al., 2016; Bergemann and Bonatti, 2019; Bergemann and Ottaviani, 2021; Goldfarb and Tucker, 2023). Hidir and Vellodi (2021) and Galperti and Perego (2023) study a model a la Bergemann et al. (2015) with consumer participation constraints but without transfers or prices. Chen (2022) studies a model in which plat-

forms compete for consumers' attention and profit from targeted advertising. Since consumers are ex-ante identical and are not paid by the platforms, Chen (2022) does not analyze markets for data but focuses on other important aspects of the data economy, such as the effects of privacy policies in the short run and in the long run. Jones and Tonetti (2020) study the incentives firms have to hoard data and the distortions this creates due to their non-rivalous nature (see, also, Varian, 2009; Farboodi et al., 2019); Finally, Taylor (2004), Acquisti and Varian (2005), Calzolari and Pavan (2006), Dworczak (2020), and Doval and Skreta (2023) study the implications of the repeated nature of online interactions and the learning externalities it may create. Our model is intentionally stylized and abstracts from these other important aspects of the data economy in order to convey our message more clearly and concisely.

## 2 The Model

We present a stylized model of a data economy. It features a platform (*it*), a merchant (*he*), and a unit mass of consumers (*she*). The consumers can sell their personal data to the platform. The platform uses this data to provide information to the merchant about the consumers' preferences. Finally, the merchant charges a fee to each consumer in exchange for the product he produces. A discussion of the main modeling assumption appears in Section 2.1.

Formally, each consumer has a unit demand for the product sold by the merchant. We denote her willingness to pay by  $\omega \in \Omega \subset \mathbb{R}_{++}$ . Let  $\bar{q} \in \Delta(\Omega)$  be the distribution of  $\omega$  in the population and assume  $\Omega$  is finite with  $|\Omega| \geq 2$ . Each consumer owns a *data record* that fully reveals her corresponding  $\omega$ .<sup>4</sup>

The model has two periods. In the first period, the data markets are open. The platform and the consumers trade the data records at prices  $p=(p(\omega))_{\omega\in\Omega}\in\mathbb{R}^\Omega$ , which they take as given. On the demand side of these markets, the platform chooses how many records of each type to demand. Let  $q=(q(\omega))_{\omega\in\Omega}\in\mathbb{R}^\Omega_+$  denote the composition of the *database* demanded by the platform, for which it pays a total of  $\sum_{\omega\in\Omega}q(\omega)p(\omega)$ . On the supply side, each consumer chooses whether to sell her record to the platform. If a type- $\omega$  consumer sells

<sup>&</sup>lt;sup>4</sup>As in Galperti et al. (2023), a data record is a list of identifiers (e.g., IP address, telephone number, etc.) and personal characteristics (e.g., gender, age, etc.) of a consumer. The former grants access to this consumer and, thus, the ability to intermediate her with the merchant. The latter provides information about her type  $\omega$ .

her record, the platform pays her the price  $p(\omega)$  and later intermediates her with the merchant, as described below. Without loss of generality, we assume that consumers of the same type sell their records with the same probability, denoted by  $z(\omega) \in [0,1]$ . Conversely, if a type- $\omega$  consumer does not sell her record to the platform, she forgoes the opportunity to interact with the merchant and obtains a reservation utility of  $r(\omega) \geq 0$ .

In the second period, the product market is open. The platform uses the acquired database q—whose composition is publicly known—to mediate the interaction between the merchant and the subset of consumers who have sold their records. In particular, the platform acts as an information intermediary: It provides the merchant with information about the consumers in the database. Formally, the platform solves a standard information-design problem where the relative frequency of consumers' types is given by q. The platform commits to an information structure that maps the record of each consumer in its database into random signals. Given the signal received, the merchant sets fee  $a \in A$  for the consumer, who then purchases the product if and only if the merchant's fee a is lower than her willingness to pay  $\omega$ . Therefore, given  $\omega$  and a, the consumer's and the merchant's second-period payoffs can be written as  $u(a, \omega) = \max\{\omega - a, 0\}$  and  $\pi(a, \omega) = a\mathbb{1}(\omega \ge a)$ , respectively. Finally, the platform's payoff is a linear combination of the consumer's trading surplus and the merchant's profits, that is,  $v(a, \omega) = \gamma_u u(a, \omega) + \gamma_\pi \pi(a, \omega)$ . We assume  $\gamma_u, \gamma_\pi \ge 0$ , with at least one strict inequality.

By standard results from the information-design literature (e.g., see Bergemann and Morris, 2016), the platform's problem in the second period can be formulated as choosing a recommendation mechanism  $x: \Omega \to \Delta(A)$  that solves:

$$\begin{split} V(q) &= \max_{x:\Omega \to \Delta(A)} \quad \sum_{a,\omega} v(a,\omega) x(a|\omega) q(\omega) \\ &\text{such that} \quad \sum_{\omega} \left( \pi(a,\omega) - \pi(a',\omega) \right) x(a|\omega) q(\omega) \geq 0 \qquad \forall \ a,a' \in A. \end{split} \tag{$\mathcal{P}_q$}$$

Without loss of generality, we let  $A = \Omega$ .

To summarize, we have introduced four endogenous variables: prices p for the data records; the consumers' decisions z to supply their records; the platform's demanded database q; and the platform's mechanism x for problem  $\mathcal{P}_q$ . We define an equilibrium of this economy as follows:

**Definition 1.** A profile  $(p^*, z^*, q^*, x^*)$  is an equilibrium of the competitive economy if

(a). Given  $p^*$ ,  $q^*$  solves the platform's problem in the first period, that is,

$$q^* \in \arg\max_{q \in \mathbb{R}^{\Omega}_+} V(q) - \sum p^*(\omega) q(\omega).$$
 (1)

- (b). Given  $q^*$ ,  $x^*$  solves the platform's problem  $\mathcal{P}_{q^*}$  in the second period.
- (c). Given  $x^*$  and  $p^*$ ,  $z^*$  solves the consumers' problem in the first period. That is, for all  $\omega$ ,

$$z^*(\omega) \in \arg\max_{\zeta \in [0,1]} \zeta \Big( p^*(\omega) + \sum_a x^*(a|\omega) u(a,\omega) \Big) + (1-\zeta) r(\omega).$$

(d). Data markets clear. That is, for all  $\omega$ ,  $q^*(\omega) = z^*(\omega)\bar{q}(\omega)$ .

Conditions (a) and (b) require that the platform acquire a database that maximizes its payoff taking prices as given, while anticipating it will use its data optimally in the second period. Condition (c) requires that each type- $\omega$  consumer choose  $z(\omega)$  optimally, again taking prices as given and anticipating that the platform will acquire a database  $q^*$  and use it to implement mechanism  $x^*$ . Therefore, she sells her record at price  $p(\omega)$  only if  $p(\omega) + \sum_a u(a,\omega)x^*(a|\omega) \geq r(\omega)$ , where  $\sum_a u(a,\omega)x^*(a|\omega)$  captures her expected trading surplus. Finally, condition (d) requires that the demand of each type of record equals its supply. This last condition pins down data prices, in the spirit of a traditional competitive equilibrium. Proposition B.1 in the Online Appendix shows that an equilibrium exists.

Hereafter, we will refer to the two-period model just described as the "competitive economy." In the next sections, we will study the equilibria of this economy, their inefficiency, and possible remedies to it.

## 2.1 Discussion of Modeling Assumptions

Before proceeding, we briefly discuss our main modeling assumptions. Our economy features a single platform taking the prices of data records as given. This assumption has a substantive component and an expositional one. The substantive component is that the platform is a price taker, and thus the data market is competitive. This is a distinguishing feature of this paper, and it allows us to shut down other potential sources of inefficiency—such as the platform's market power—that have less to do with data as an input. The expositional component is that we focus on a single platform rather than a finite number of identical ones. This substantially

simplifies notation at little cost to generality. Galperti and Perego (2022) show how to model a competitive economy with multiple competing platforms.

The platform's objective is assumed to be linear in the consumers' trading surplus and the merchant's profits. This specification is especially tractable, while capturing key features of real-world two-sided markets (Xu and Yang, 2023, offer a dynamic microfoundation for such an objective). The results of Section 4 do not depend on this assumption.

Three aspects of the consumer's problem have been simplified. First, the reservation utility  $r(\omega)$  is exogenous. This assumption rules out settings in which the consumer can bypass the platform and trade directly with the merchant. While not focusing on data markets, Bergemann and Bonatti (2023) study the interplay of online and offline interactions. Second, the consumer cannot participate in the platform's mechanism without revealing her type. That is, the data record bundles "access" to the consumer and information about her willingness to pay, which can be restrictive in some applications. With different goals than ours, Hidir and Vellodi (2021) and Ali et al. (2022) study models in which these two aspects are unbundled. Third, selling the data record fully reveals the underlying consumer's type. This simplifies notation since the type of a consumer and that of her record are the same. Our analysis can be extended to records that are only partially informative of the consumer's type, a model of which is proposed by Galperti et al. (2023).

## 2.2 Efficiency Benchmark

This paper aims to shed light on how efficiently the data markets allocate records between the consumers and the platform. To answer this, we first need to introduce an efficiency benchmark. Let us refer to the pair (q, x) as an *allocation*. For each allocation, denote the sum of the payoff of the platform and the consumers as

$$\mathcal{W}(q,x) \triangleq \sum_{a,\omega} \left( v(a,\omega) + u(a,\omega) \right) x(a|\omega) q(\omega) + \sum_{\omega} \left( \bar{q}(\omega) - q(\omega) \right) r(\omega). \tag{2}$$

Note that since data prices p only affect the distribution of payoffs between the platform and the consumers, this notion of welfare only depends on the allocation (q, x).

**Definition 2.** An allocation  $(q^{\circ}, x^{\circ})$  is **constrained efficient** if it solves

$$W^{\circ} = \max_{q,x} \quad \mathcal{W}(q,x)$$
 such that  $q \leq \bar{q}$ ,  $(\mathcal{SB})$  and  $x \text{ solves } \mathcal{P}_q$ .

According to this efficiency benchmark, an allocation is constrained efficient if it maximizes the welfare function W(q, x) subject to two constraints. The first requires that the database q is feasible, i.e., it does not allocate to the platform more records than those that exist in the economy. The second constraint requires that the mechanism x is sequentially optimal for the platform given q.<sup>5</sup>

We briefly motivate the efficiency benchmark just introduced. Since the main goal of the paper is to demonstrate that equilibria of this economy can be inefficient, a less demanding efficiency benchmark is more desirable, as it makes such a negative result starker. Additionally, a less demanding benchmark allows us to ignore other potential sources of inefficiency that are rather standard and not special to data as an input. In particular, two aspects of Definition 2 make it less demanding. First, it focuses on "constrained" efficiency, as it requires that the mechanism x be sequentially optimal for the platform given q. Dropping this constraint, i.e., focusing on unconstrained efficiency, would lead us to detect an additional inefficiency that is merely driven by the fact that, in the first period, the platform cannot commit to a mechanism for the second period. Second, we exclude the merchant's payoff from the welfare function W. Including the merchant's payoff would lead us to detect an additional inefficiency that is merely driven by the fact that the platform does not "sell" information to the merchant and, thus, cannot transfer the merchant's profit to the consumer. Therefore, the consumer does not internalize the effects that selling her data create on the profit of the merchant.

An additional useful feature of our welfare notion is that in any equilibrium of the competitive economy,  $W(q^*, x^*)$  coincides with the consumers' welfare because the platform must earn a payoff of zero.<sup>7</sup> Therefore, any constrained efficient equilibrium also maximizes con-

 $<sup>^5</sup>$ A constrained efficient outcome always exists. This follows from the fact that  $\mathcal{W}$  is continuous and, by Lemma B.1, the feasible set of outcomes in the planner's problem is nonempty and compact.

<sup>&</sup>lt;sup>6</sup>Online Appendix D shows that our main qualitative results are unchanged when we consider "unconstrained" efficiency or a welfare function that includes the merchant profits.

<sup>&</sup>lt;sup>7</sup>Indeed, notice that V(q) is homogeneous of degree one. If the platform earned a strictly positive payoff at  $q^*$ , it could profitably deviate by acquiring database  $q' = \alpha q^*$ , with  $\alpha > 1$ , which earns a payoff V(q') –

sumers' welfare.

## 3 The Inefficiency of the Data Market

In this section, we present a series of results that identify necessary and sufficient conditions for equilibrium efficiency and, by doing so, we uncover the key driver of inefficiency in our competitive economy.

We begin by asking two questions that are instrumental to our analysis. What is the social cost—that is, the decrease in the welfare function  $\mathcal{W}$ —that results from allocating an additional  $\omega$ -record to the platform's database? Clearly, it is  $r(\omega)$ , the reservation utility that is lost by the corresponding consumer. What is the corresponding social benefit? To compute it, fix an arbitrary database q and consider the following maximization problem:

$$W(q) \triangleq \max_{x:\Omega \to \Delta(A)} \sum_{a,\omega} (v(a,\omega) + u(a,\omega))x(a|\omega)q(\omega)$$
such that  $x$  solves  $\mathcal{P}_q$ . (3)

We can think of the above as the problem of a planner who chooses a mechanism x to maximize the welfare of the platform and the consumers in database q. This planner is constrained to choose a mechanism that, given q, the platform would also be willing to implement. We denote by  $\Psi_q$  the set of supergradients of W(q). In other words, just like a derivative, each  $\psi_q(\omega)$  captures how W(q) changes when we add an additional  $\omega$ -record to database q. For this reason,  $\psi_q(\omega)$  identifies the *social benefit* of allocating an additional  $\omega$ -record into the platform's database.<sup>8</sup>

Our first result demonstrates how the social benefit and cost of data records can be used to characterize which allocations are constrained efficient.

**Proposition 1.** An allocation (q, x) is constrained efficient if and only if x solves  $\mathcal{P}_q$  and there exists a  $\psi_q \in \Psi_q$  such that, for all  $\omega$ ,

$$\overline{\sum_{\omega} p^*(\omega) q'(\omega) = \alpha \big( V(q^*) - \sum_{\omega} p^*(\omega) q^*(\omega) \big)} > V(q^*) - \sum_{\omega} p^*(\omega) q^*(\omega).$$

<sup>&</sup>lt;sup>8</sup>In Appendix A.1 we explicitly compute W(q) (Equation (A.2)) and show it is concave and, therefore,  $\Psi_q$  is well-defined. Based on this, Lemma A.1 provides an analytical characterization of  $\Psi_q$ , which shows that  $\Psi_q$  is generically a singleton and is easy to compute. This result generalizes Proposition 2 of Galperti et al. (2023) to the case when the supergradient is not unique.

$$-\psi_q(\omega) \ge r(\omega) \text{ if } q(\omega) > 0,$$

$$- \psi_q(\omega) \le r(\omega) \text{ if } q(\omega) < \bar{q}(\omega).$$

It is clear that under any constrained-efficient allocation, a record allocated to the platform's database should have a social benefit that exceeds its social cost. Perhaps less intuitively, these conditions are also sufficient for constrained efficiency. This will be key to characterizing equilibrium efficiency in terms of the model primitives. To see why, fix an equilibrium  $(p^*, z^*, q^*, x^*)$  and denote by

$$U^*(\omega) \triangleq p^*(\omega) + \sum_{a} x^*(a|\omega)u(a,\omega) \tag{4}$$

the *private benefit* that a type- $\omega$  consumer obtains when selling her record to the platform. Notice that the equilibrium conditions require that  $U^*(\omega) \geq r(\omega)$  if  $q^*(\omega) > 0$ , and that  $U^*(\omega) \leq r(\omega)$  if  $q^*(\omega) < \bar{q}(\omega)$ . Therefore, in light of Proposition 1, this equilibrium is constrained efficient if and only if the private and social benefits of data records are sufficiently "aligned." That is, if there is a  $\psi_{q^*} \in \Psi_{q^*}$  such that  $\psi_{q^*}(\omega) \geq r(\omega)$  if  $U^*(\omega) \geq r(\omega)$  and, conversely,  $\psi_{q^*}(\omega) \leq r(\omega)$  if  $U^*(\omega) \leq r(\omega)$ .

The key question is then, under what conditions on the model primitives, are  $U^*$  and  $\psi_{q^*}$  aligned. To address this question, it is useful to define  $\sigma^*(\omega) \triangleq \psi_{q^*}(\omega) - p^*(\omega)$  and write

$$\psi_{q^*}(\omega) = p^*(\omega) + \sigma^*(\omega). \tag{5}$$

This decomposition of the social benefit, which is analogous to the definition of  $U^*$  in Equation (4), has a particularly useful economic interpretation:  $p^*(\omega)$  and  $\sigma^*(\omega)$  capture the marginal change in the platform's payoff and the consumers' trading surplus, respectively, that result from adding an  $\omega$ -record to  $q^*$ . The following result formalizes this interpretation.

**Lemma 1.** In any equilibrium  $(p^*, z^*, q^*, x^*)$ ,  $p^*$  is a supergradient of  $V(q^*)$ .

This result demonstrates that, due to the competitive nature of the data markets, the marginal change in the platform's payoff from acquiring an additional  $\omega$ -records—formally, the supergradient of  $V(q^*)$ —must equal its cost  $p^*(\omega)$ . Intuitively, if that was not the case, the platform would strictly prefer to change her demand of  $\omega$ -records. As a consequence of Lemma 1, the

<sup>&</sup>lt;sup>9</sup>In Appendix A.1, we show V(q) is concave and provide an analytical characterization of its supergradients in Lemma A.1. In the language of Galperti et al. (2023), the supergradients are the "values of an  $\omega$ -record" for the platform. This implies that equilibrium data prices  $p^*$  inherently account for the "pooling" externalities identified by Galperti et al. (2023).

remaining component of the social benefit  $\psi_{q^*}(\omega)$ —namely,  $\sigma^*(\omega)$ —captures the marginal change in the trading surplus of all consumers that results from adding a  $\omega$ -record to the database  $q^*$ .

Comparing Equations (4) and (5) reveals that equilibrium efficiency hinges on the alignment between the trading surplus that a type- $\omega$  consumer expects to receive when selling her record—namely,  $\sum_a x^*(a,\omega)u(a,\omega)$ —and the effect that this sale has on the trading surplus of *all* consumers—namely,  $\sigma^*(\omega)$ . Since the consumer internalizes the former but not the latter, her decision to sell her record may exert an externality on other consumers, thus introducing inefficiency in the economy.

The following result shows how the presence of these externalities, and thus the efficiency of the economy, depends on the model primitives and, in particular, hinges on the objective of the platform.

**Proposition 2.** If  $\gamma_{\pi} > \gamma_{u}$ , all equilibrium allocations of the competitive economy are constrained efficient and, therefore, maximize consumers' welfare. Conversely, if  $\gamma_{\pi} \leq \gamma_{u}$ , equilibrium allocations can be inefficient.

Perhaps counterintuitively, if the platform cares relatively more about the merchant's profits, the social and private benefits of data records are aligned. Thus, equilibria are constrained efficient and consumers' welfare is maximal. Vice versa, if the platform cares relatively more about the consumers' surplus, this alignment can break, leading to inefficiencies.

To gain intuition, suppose  $\gamma_{\pi} > \gamma_u$  and consider an arbitrary database  $q \neq 0$ . In this case, the platform finds it optimal to reveal the  $\omega$  of each consumer in the database to the merchant, allowing the merchant to extract all their surplus. To see this, notice that, by Bergemann et al. (2015, Theorem 1), the platform's problem  $\mathcal{P}_q$  is equivalent to choosing a point in the triangle of Figure 1, which plots the set of pairs of merchant's profit and consumers' surplus that can be induced by any mechanism. Since the platform's payoff v is linear in  $\pi(a,\omega)$  and  $u(a,\omega)$ , when  $\gamma_{\pi} > \gamma_u$ , the optimal mechanism maximizes the merchant's profits, leaving consumers with no surplus. As a consequence,  $\sum_a x^*(a,\omega)u(a,\omega) = \sigma^*(\omega) = 0$ . Therefore, consumers do not exert externalities on each other when selling their records, and all equilibria are constrained efficient.

Suppose instead  $\gamma_{\pi} \leq \gamma_{u}$ . In this case, for any q, the optimal mechanism  $x^{*}$  typically involves withholding some information from the merchant to prevent excessive surplus extraction

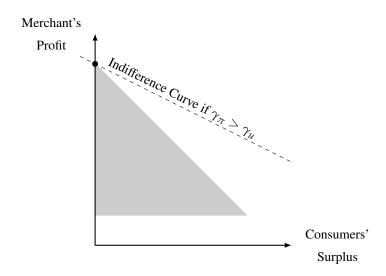


Figure 1: The trading surplus triangle. The dashed line depicts the platform's indifference curve when  $\gamma_{\pi} > \gamma_{u}$ .

(again, see Figure 1). To do so, the platform may pool a consumer's record with those of others. In this case, their payoffs become interdependent, as the presence of this record in the pool can affect the merchant's beliefs in a nontrivial way. Consequently,  $\sigma^*(\omega)$  and  $\sum_a x^*(a|\omega)u(a,\omega)$  can differ, leading to inefficiencies. The externality exerted by low-type consumers is typically positive. To see why, consider a consumer of the lowest type,  $\underline{\omega} \triangleq \min_{\omega} \Omega$ . Her expected trading surplus is zero since the merchant would never charge a fee a lower than  $\underline{\omega}$ . Yet,  $\sigma^*(\underline{\omega})$  can be strictly positive, because when this consumer is pooled with higher-type consumers, she helps them receive a lower fee, and thus earn a higher surplus. For opposite reasons, the externality exerted by high-type consumers is typically negative. Failure to internalize these externalities leads to inefficiencies.

This is the sense in which the market failure we emphasize is enabled by the role played by the platform as an information intermediary. It originates from the platform's incentive to withhold information to harness the conflicting objectives of the constituencies it intermediates. The latter is a distinguishing feature of multi-sided platforms, which explains why information withholding is so commonly observed in the marketplace. <sup>10</sup> Importantly, the market failure we emphasize does not conceptually originate from the absence of competition for the

<sup>&</sup>lt;sup>10</sup>For instance, Google restricts advertisers' access to users' characteristics during ad auctions, aiming to foster competition and optimize ad revenues. Similarly, ridesharing giants such as Uber conceal riders' destination details from drivers, discouraging cherry-picking practices.

merchant and, instead, could arise even in a richer model with multiple competing merchants. For example, in a model in which merchants compete in a second-price auction, Bergemann et al. (2022b) show that a platform has incentives to withhold information from them. This scenario would result in inefficiencies akin to those highlighted in this paper.

We conclude this discussion by finding sufficient conditions under which *all* equilibria are inefficient, thus sharpening the negative message of Proposition 2. To avoid trivial cases, let us focus on economies in which  $W^{\circ} > R := \sum_{\omega} \bar{q}(\omega) r(\omega)$ , that is, the constrained-efficient allocation involves some trade.<sup>11</sup>

**Corollary 1.** Let  $\gamma_{\pi} \leq \gamma_{u}$  and suppose  $W^{\circ} > R$ . If  $\gamma_{u}\underline{\omega} < r(\underline{\omega}) < (1 + \gamma_{u})\underline{\omega}$ , then all equilibria are inefficient.

Corollary 1 gives a sufficient condition under which the positive externality discussed above causes all equilibria of the economy to be inefficient. First, we show that if  $\gamma_u\underline{\omega} < r(\underline{\omega})$ , the platform is unwilling to pay a price higher than  $r(\underline{\omega})$  for  $\underline{\omega}$ -records. This implies that  $U^*(\underline{\omega}) < r(\underline{\omega})$ , since a consumer of type  $\underline{\omega}$  necessarily earns a zero trading surplus when she sells her record. Thus, no such consumer sells her record. Second, we additionally show that, if  $r(\underline{\omega}) < (1 + \gamma_u)\underline{\omega}$ , the social benefit of  $\underline{\omega}$ -records,  $\psi_{q^*}(\underline{\omega})$ , exceed its private cost,  $r(\underline{\omega})$ . As a result, a trade that would be socially beneficial does not occur in equilibrium, generating an inefficiency.

Under the sufficient condition of Corollary 1, the  $\underline{\omega}$ -type consumers exert a positive externality on other consumers, which they fail to internalize, thus leading to inefficiencies. Conversely, Corollary A.1 in Appendix A.1 provides alternative sufficient conditions under which the inefficiency in the economy is caused by a negative externality exerted by higher-type consumers. In general, both positive and negative externalities exist, as illustrated by the next example.

## 3.1 An Example

We now provide an example to illustrate why the data economy can be inefficient. Suppose there are two types of consumers,  $\Omega=\{1,2\}$ , and  $\bar{q}(2)>\bar{q}(1)$ . All consumers have the same reservation utility:  $r(\omega)=\bar{r}\in(0,1)$  for all  $\omega$ . The platform only cares about consumers' trading surplus:  $\gamma_u>\gamma_\pi=0$ . To avoid uninteresting cases, we assume that  $\bar{r}<\frac{1+\gamma_u}{2}$  so that

When  $W^{\circ}=R$ , all equilibria are constrained efficient. To see this notice that in any equilibrium  $(p^*,z^*,q^*,x^*)$ , it must be that  $R \leq \mathcal{W}(q^*,x^*) \leq W^{\circ}$ . Therefore, if  $W^{\circ}=R$ , we have  $\mathcal{W}(q^*,x^*)=W^{\circ}$ .

some trade is required to achieve constrained efficiency.<sup>12</sup> We will show that all equilibria of this economy are inefficient.

We first characterize the constrained-efficient allocations. Let  $(q^{\circ}, x^{\circ})$  be such that  $q^{\circ}(1) = q^{\circ}(2) = \bar{q}(1)$  and  $x^{\circ}(1|\omega) = 1$  for all  $\omega$ . In other words, the platform is given the records of all the low-type consumers and an equal amount of high-type ones. The platform then pools all these records in the same segment, inducing the merchant to charge the lowest fee (i.e., a=1) to all consumers in the database. We argue that  $(q^{\circ}, x^{\circ})$  is the unique constrained-efficient allocation. To see why, consider any other database q>0 and notice that an optimal mechanism  $x_q$  given q is to set  $x_q(1|\omega)=\min\{q(1),q(2)\}/q(\omega)$  for all  $\omega$ . That is, the platform creates the largest possible segment with an equal quantity of low- and high-type consumers, who are then charged the lowest fee. Thus, the allocation  $(q,x_q)$  induces a welfare of  $\mathcal{W}(q,x_q)=(1+\gamma_u)\min\{q(1),q(2)\}+(1-q(1)-q(2))\bar{r}$ . To maximize  $\mathcal{W}(q,x_q)$ , any constrained-efficient allocation  $(q^{\circ},x^{\circ})$  must satisfy  $q^{\circ}(1)=q^{\circ}(2)$ . Since by assumption  $\bar{q}(1)<\bar{q}(2)$  and  $\bar{r}<\frac{1+\gamma_u}{2}$ , setting  $q^{\circ}(1)=q^{\circ}(2)=\bar{q}(1)$  uniquely maximizes  $\mathcal{W}(q,x_q)$ . For future reference, note that welfare under the constrained-efficient allocation is  $W^{\circ}=\bar{r}+\bar{q}(1)(1+\gamma_u-2\bar{r})$ .

We now characterize all equilibria and show that they are inefficient. Let  $(p^*, z^*, q^*, x^*)$  be an equilibrium. Let  $a_{q^*}$  be the fee the merchant would charge if the platform did not provide him with any information besides the database composition. Using Lemma A.1 in Appendix A, we can compute the marginal change in the consumers' trading surplus that results from adding an  $\omega$ -record to  $q^*$ , which is

$$\sigma^*(\omega) = \omega - a_{q^*} \mathbb{1}(\omega \ge a_{q^*}), \tag{6}$$

and the equilibrium price, which is

$$p^*(\omega) = \gamma_u \sigma^*(\omega) = \gamma_u \Big( \omega - a_{q^*} \mathbb{1}(\omega \ge a_{q^*}) \Big). \tag{7}$$

Since  $1 \le a_{q^*} \le 2$ , we have  $0 \le \sigma^*(\omega) \le 1$  and  $0 \le p^*(\omega) \le \gamma_u$ , for all  $\omega$ .

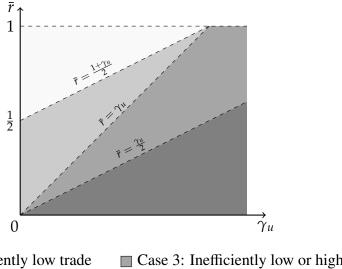
Case 1:  $\bar{r} > \gamma_u$ —Inefficiently Low Trade. The only equilibrium involves no trading and—since Corollary 1 applies to this case—is inefficient. To see why, we first argue that any equilibrium involves no trade of type-1 records:  $q^*(1) = 0$ . Indeed, type-1 consumers always earn

When  $\bar{r} \ge \frac{1+\gamma_u}{2}$ , the absence of trade—i.e.,  $q_0 = (0,0)$ —is constrained efficient. That is,  $W^\circ = R = \bar{r}$ . In this case, any equilibrium allocation  $(q^*, x^*)$  is constrained efficient since  $W^\circ \ge W(q^*, x^*) \ge \bar{r}$ .

a zero trading surplus when they sell their records. Moreover, by Equation (7), the price  $p^*(1)$  can be no higher than  $\gamma_u$ . Therefore, their net payoff is no higher than  $\gamma_u$ , which is strictly smaller than  $\bar{r}$  by assumption. Thus, type-1 consumers do not sell their data in any equilibrium:  $q^*(1) = 0$ . Next, we argue that this implies  $q^*(2) = 0$ . To see why, suppose  $q^*(2) > 0$ . Since  $q^*(1) = 0$ , we must have  $a_{q^*} = 2$ , and thus Equation (7) implies that  $p^*(2) = 0$ . Moreover, since  $q^*(1) = 0$ , type-2 consumers will be perfectly discriminated against and earn a zero trading surplus. Since type-2 consumers get a zero net payoff when they sell their data, they must be unwilling to do so, contradicting  $q^*(2) > 0$ . Therefore,  $q^* = (0,0)$  is the only database compatible with equilibrium. Under this complete market unraveling, any equilibrium allocation  $(q^*, x^*)$  must yield  $\mathcal{W}(q^*, x^*) = \bar{r} < W^\circ$ , and thus the inefficiency is as severe as it could be.

Why are these equilibria inefficient? By selling her record, a low-type consumer could create a positive externality: The platform would pool this consumer with a high-type consumer, thus creating a social benefit of  $1 + \gamma_u$ , which by assumption is larger than  $2\bar{r}$ , that is, the sum of the reservation utilities of these two consumers. For this trade to happen, however, the price  $p^*(1)$  has to exceed the reservation utility  $\bar{r}$ . Unfortunately, the platform is unwilling to pay such a price, since its value for a type-1 record is, at most,  $\gamma_u < \bar{r}$ .

Case 2:  $\bar{r} < \frac{\gamma_u}{2}$ —Inefficiently High Trade. When  $\bar{r} < \frac{\gamma_u}{2}$ , the unique equilibrium involves  $q^*(1) = \bar{q}(1)$  and  $q^*(2) = \min\{\bar{q}(2), \frac{\bar{q}(1)}{\bar{r}}\}$ . To see this, note first that by the no-profit condition, the equilibrium prices must satisfy  $p^*(1) + p^*(2) \ge \gamma_u$ . If not, the platform could acquire a pair of low- and high-type records, pool them in the same segment, and generate a payoff of  $\gamma_u$ , hence generating a positive profit. Therefore, either  $p^*(1)$  or  $p^*(2)$  must exceed  $\frac{\gamma_u}{2}$ . We argue that  $p^*(2) \le \frac{\gamma_u}{2}$ . Otherwise, all high-type consumers would strictly prefer to sell, leading to an uninformed merchant's price of  $a_{q^*} = 2$ , which by Equation (7) implies  $p^*(2) = 0$ , a contradiction. Since then  $p^*(1) \ge \frac{\gamma_u}{2}$ , the low-type consumers strictly prefer to sell, and thus  $q^*(1) = \bar{q}(1)$ . As in the constrained-efficient allocation, the optimal mechanism involves  $x^*(1|2) = \min\{\bar{q}(1), q^*(2)\}/q^*(2)$ : The platform creates a segment that includes all the low-type consumers and as many high-type consumers as possible subject to inducing the fee a = 1. Since  $\bar{r} < 1$ , we must have  $q^*(2) > q^*(1)$ . Otherwise, the expected trading surplus of a high-type consumer selling her record would be 1, and thus all such consumers would sell their records, leading to a contradiction. Given such  $q^*$ , note that  $a_{q^*} = 2$  and, by Equation (7),  $p^*(2) = 0$ . Thus,  $q^*(2) = \min\{\bar{q}(2), \frac{\bar{q}(1)}{\bar{p}}\}$ . We conclude that the unique equilibrium is



- ☐ Case 1: Inefficiently low trade ■ Case 3: Inefficiently low or high trade
- Case 2: Inefficiently high trade ☐ No trade

Figure 2: Equilibrium Inefficiency in the Example of Section 3.1

inefficient because 
$$\mathcal{W}(q^*, x^*) = (1 + \gamma_u)\bar{q}(1) + \max\{0, \bar{r} - \bar{q}(1)(1 + \bar{r})\} < W^\circ$$
.

In this equilibrium, too many high-type consumers inefficiently sell their records to the platform:  $q^*(2) > q^{\circ}(2)$ . They are attracted by the possibility of buying the merchant's product at the lowest fee. However, when an additional high-type consumer sells her record, she exerts a negative externality on other consumers, by undermining their chances of buying the product at the lowest fee. On the one hand, such a consumer gains individually from selling her record, since  $\sum_a x^*(a|2)u(a,2) = q^*(1)/q^*(2) \ge \bar{r}$ . On the other hand, she does not help increase the aggregate trading surplus for all consumers. Indeed,  $\sigma^*(2) = 0$  (by Equation (6)).<sup>13</sup>

We cover the case of  $\frac{\gamma_u}{2} \leq \bar{r} \leq \gamma_u$  in Appendix C. It shares similar intuitions to the previous two cases and features multiple equilibria with either inefficiently low or inefficiently high trade. In summary, all equilibria of this simple economy are inefficient as represented in Figure 2.

<sup>&</sup>lt;sup>13</sup>One may wonder whether a negative price for type-2 records could correct this inefficiency, as it would disincentivize selling these records. Such a price is incompatible with equilibrium behavior, as the platform would then demand an infinite quantity of type-2 records. More generally, Corollary A.1 in Appendix A provides sufficient conditions in an arbitrary economy for similar inefficiencies.

### 4 Remedies to the Market Failure

This section discusses three market designs that provide a remedy to the inefficiency of the competitive economy. The first establishes a "data union" that manages consumers' data on their behalf. The second introduces data taxes in the competitive economy. The third renders data markets "more complete."

### 4.1 Data Unions

We first introduce a new intermediary in the economy, called a *data union*.<sup>14</sup> A data union manages consumers' data records on their behalf to maximize their welfare: It collects records from participating consumers, sells some or all of them to the platform, and distributes the proceeds back to the consumers to incentivize participation. Thus, the union coordinates consumers' actions by unilaterally deciding which records should be sold to the platform and how consumers should be compensated. We show that a data union implements allocations that are constrained efficient. This result offers theoretical support to recent policy proposals that discuss the potential role that a data union could play in the data economy (e.g., see Posner and Weyl, 2018; Bergemann et al., 2023).

More formally, our data union operates as follows. Consumers voluntarily decide whether to join the union and relinquish their records to it. Given the collected database  $\hat{q}$ , the union sells part of it,  $q \leq \hat{q}$ , to the platform at a price that extracts all the platform's expected payoff V(q). Finally, the union distributes the proceeds V(q) to its members through transfers  $p \in \mathbb{R}^{\Omega}$  such that  $\sum_{\omega} p(\omega)\bar{q}(\omega) \leq V(q)$ . The union's objective is to maximize the welfare of its members. We assume that, if a consumer joins the union, she retains her reservation utility  $r(\omega)$  unless her record is sold to the platform. This assumption is critical for our result. A by-product of this assumption is that, without loss of generality, all consumers join the union, i.e.,  $\hat{q} = \bar{q}$ . Thus, the problem of the data union can be written as follows:

$$\max_{(p,q,x)} \quad \sum_{\omega} p(\omega) \bar{q}(\omega) + \sum_{a,\omega} u(a,\omega) x(a|\omega) q(\omega) + \sum_{\omega} (\bar{q}(\omega) - q(\omega)) r(\omega)$$
 such that  $q \leq \bar{q}$ , and  $x$  solves  $\mathcal{P}_a$ ,

<sup>&</sup>lt;sup>14</sup>Data unions have also been discussed in Posner and Weyl (2018) and Bergemann et al. (2023).

and 
$$\sum_{\omega} p(\omega) \bar{q}(\omega) \leq V(q),$$
 and 
$$p(\omega) + \frac{q(\omega)}{\bar{q}(\omega)} \sum_{a} u(a,\omega) x(a|\omega) + \left(1 - \frac{q(\omega)}{\bar{q}(\omega)}\right) r(\omega) \geq r(\omega).$$

The first two constraints are familiar from Definition 2. The third constraint requires that the data union does not run a deficit. The last constraint guarantees consumers' participation. By participating, a type- $\omega$  consumer receives compensation  $p(\omega)$ ; with probability  $\frac{q(\omega)}{\bar{q}(\omega)}$ , her record is sold to the platform, in which case she also earns a trading surplus of  $\sum_a u(a,\omega)x(a|\omega)$ ; with remaining probability, her record is not sold and she gets her reservation utility  $r(\omega)$ . Notice that, in this specification, all consumers are entitled to a payment  $p(\omega)$ , regardless of whether their data records are used. This assumption is immaterial for the result.

**Proposition 3.** Let  $(p^*, q^*, x^*)$  be a solution to the data union's problem. The allocation  $(q^*, x^*)$  is constrained efficient. Conversely, if  $(q^{\circ}, x^{\circ})$  is a constrained-efficient allocation, there exists  $p^{\circ}$  such that  $(p^{\circ}, q^{\circ}, x^{\circ})$  is a solution to the data union's problem.

There are two main differences from the competitive data economy in Section 3. First, while consumers can decide to leave the union, they have no say in whether their data records are sold to the platform. This allows the union to coordinate consumers in a way that competitive markets cannot. Second, the union has bargaining power vis-à-vis the platform and the consumers: The prices p are chosen by the union rather than being determined by market clearing. As a result, the union both internalizes the externalities discussed in the previous section and can properly compensate consumers for their participation.

It is worth noting that the data union's objective of maximizing consumers' welfare is not essential for its ability to induce constrained efficiency allocations. If the union maximized the platform's payoff—or, equivalently, if the platform was the union—then it can be shown that the induced allocations would still be constrained efficient. Clearly, such a union would have no incentive to distribute the proceeds of its activity back to the consumers. Indeed, while allocations would be constrained efficient, consumers' welfare would be minimized—i.e., equal to  $R = \sum_{\omega} \bar{q}(\omega) r(\omega)$ , with all consumers earning their reservation utility.

### 4.2 Data Taxes

While the data union helps achieve constrained-efficient allocations, it involves a concentration of bargaining power in a single institution. One may wonder if constrained-efficient allocations can be decentralized in our original competitive economy. This section provides a positive answer by introducing data taxes levied on consumers.<sup>15</sup>

These taxes work as follows. Whenever a type- $\omega$  consumer sells her record to the platform, she pays a "data tax"  $\tau(\omega) \in \mathbb{R}$  to the government. We can interpret  $\tau(\omega) \leq 0$  as a subsidy paid by the government. To make the problem interesting, suppose the government cannot run a deficit, i.e.,  $\sum_{\omega} q(\omega)\tau(\omega) \geq 0$ . Besides taxes, all other components of the model are as in Section 2. The next result shows that any constrained-efficient allocation can be supported as an equilibrium of the economy with taxation.

**Proposition 4.** Let  $(q^{\circ}, x^{\circ})$  be a constrained-efficient allocation. There exists a profile of taxes  $\tau^*$ , of prices  $p^*$ , and of consumer choices  $z^*$ , such that  $(p^*, z^*, q^{\circ}, x^{\circ})$  is an equilibrium of the economy with taxation  $\tau^*$  and the government does not run a deficit.

We can explicitly characterize the taxes and the equilibrium that supports any constrained-efficient allocation  $(q^{\circ}, x^{\circ})$ . Let  $p^*$  be a supergradient of  $V(q^{\circ})$  and define

$$\tau^*(\omega) \triangleq p^*(\omega) + \sum_{a} x^{\circ}(a|\omega)u(a,\omega) - r(\omega). \tag{8}$$

Additionally, define  $z^*(\omega) \triangleq q^\circ(\omega)/\bar{q}(\omega)$  for all  $\omega$ . It is easy to check that  $(p^*,z^*,q^\circ,x^\circ)$  is an equilibrium of the economy with taxation  $\tau^*$ . First, since  $p^*$  is a supergradient of  $V(q^\circ)$ ,  $q^\circ$  must solve the platform's problem in the first period. Moreover, since  $(q^\circ,x^\circ)$  is constrained efficient,  $x^\circ$  must solve  $\mathcal{P}_{q^\circ}$ , i.e., the platform's problem in the second period. Third, all consumers are indifferent between selling or not selling their data records to the platform since, if they sell, they earn  $p^*(\omega) + \sum_a x^\circ(a|\omega)u(a,\omega) - \tau^*(\omega) = r(\omega)$ . Finally, by the definition of  $z^*$ , data markets clear. In this equilibrium, consumer welfare equals  $R = \sum_{\omega} \bar{q}(\omega)r(\omega)$  while the platform's payoff is zero. Since the allocation is constrained efficient, it must be that the government runs a budget surplus of  $W^\circ - R$ . If these proceeds are distributed to the consumers (in a way that does not affect their behavior, e.g., in a lump-sum manner), then consumer surplus is maximized in equilibrium.

<sup>15</sup> As usual, the data tax could also be levied on the platform and, in equilibrium, passed on to the consumers.

The data taxes correct the inefficiency of the competitive economy by making consumers indifferent between selling or not selling their data records. In particular, the taxes increase the effective data price paid to the consumers who exert a positive externality (typically, the lowtype ones) and decrease the price paid to the consumers who exert a negative externality (typically, the high-type ones). As a result, any constrained-efficient allocation can be supported in equilibrium. It is instructive to see this at play in the context of a concrete example.

**Example of a Data Tax.** Consider Case 1 in Section 3.1. We argued that low-type consumers, whose records would be socially beneficial to sell, are not sufficiently remunerated by the competitive market, leading to inefficiency. Our taxation subsidizes these consumers just enough to make them indifferent between selling or not:  $\tau^*(1) = \gamma_u - \bar{r} < 0$ . This subsidy is financed by taxing the high-type consumers:  $\tau^*(2) = 1 - \bar{r}$ . The equilibrium prices are  $p^*(1) = \gamma_u$  and  $p^*(2) = 0$ . Under these taxes and prices, the constrained-efficient allocation  $(q^{\circ}, x^{\circ})$  can be supported in equilibrium, and the government's tax revenue is  $\bar{q}(1)(1 + \gamma_u - 2\bar{r}) > 0$ .

### 4.3 Lindahl Pricing

Finally, we show how the inefficiency of the competitive economy can be corrected by enriching the price system for data records. We do so by allowing that the terms of trade between a consumer and the platform involve not only whether she sells her record, but also how the platform intends to use it.<sup>16</sup> That is, the price of a record can depend not only on its type  $\omega$ , but also on the fee a that the platform will recommend to the merchant. The basic logic behind this approach is that the aforementioned inefficiency stems from the externalities generated by how data records are used, so market participants should have ways to take those externalities into account in their trades. As such, our approach is similar to classic ways of modeling competitive economies with externalities (e.g., Arrow (1969) and Laffont (1976)).<sup>17</sup> For this reason, we refer to this setting as the Lindahl economy.

More formally, this economy features one market for each pair  $(a, \omega)$  with an associated

<sup>&</sup>lt;sup>16</sup>This is reminiscent of the European Union's GDPR, which requires that "the specific purposes for which personal data are processed should be explicit and legitimate and determined at the time of the collection of the personal data" (see, GDPR 2016/679 (39)).

<sup>&</sup>lt;sup>17</sup>In particular, we follow Bonnisseau et al. (2023).

price  $p(a,\omega)$  at which  $\omega$ -records can be traded for use a. A type- $\omega$  consumer decides in which market to sell her  $\omega$ -record, if any. That is, for all a, she chooses the probability of selling her record to the platform for use a, denoted by  $z(a,\omega) \in [0,1]$ . As in our baseline economy, the platform takes prices as given and chooses a database q and a mechanism x, with  $x(a|\omega)q(\omega)$  representing its demand of  $\omega$ -records in market  $(a,\omega)$ . It is implicit in the trade agreement between the platform and the consumers that, if the platform acquires a record for use a, it has to deliver on this promise. The platform's problem can then be written as

$$\max_{q,x} \quad \sum_{a,\omega} \Big( v(a,\omega) - p(a,\omega) \Big) x(a|\omega) q(\omega)$$
such that 
$$\sum_{\omega} \Big( \pi(a,\omega) - \pi(a',\omega) \Big) x(a|\omega) q(\omega) \ge 0 \qquad \forall \ a,a' \in A$$

It is instructive to compare the platform's problem in the Lindahl economy with the one in the baseline economy (conditions (a) and (b) in Definition 1). They only differ insofar as the Lindahl economy has richer markets, with prices that can depend on a and not just on  $\omega$ . <sup>18</sup>

The equilibrium definition in the Lindahl economy follows from Definition 1.

**Definition 3.** A profile  $(p^*, z^*, q^*, x^*)$  is an equilibrium of the Lindahl economy if

- (a) Given  $p^*$ ,  $(q^*, x^*)$  solves the platform's problem in Equation (9).
- (b) Given  $p^*$ ,  $z^*$  solves the consumers problem: for all  $\omega$ ,

$$z^*(\cdot,\omega) \in \arg\max_{z \in \mathbb{R}^A_+ \text{ s.t. } \sum_a z(a) \le 1} \sum_a z(a) \left( p^*(a,\omega) + u(a,\omega) \right) + \left(1 - \sum_a z(a)\right) r(\omega).$$

(c) Markets clear:  $x^*(a|\omega)q^*(\omega) = z^*(a,\omega)\bar{q}(\omega)$  for all  $\omega$  and a.

For the Lindhal economy, we will use the following notion of effciency.

**Definition 4.** An allocation  $(q^{\dagger}, x^{\dagger})$  is unconstrained efficient if it solves

$$W^{\dagger} = \max_{q,x} \quad W(q,x)$$
such that  $q \leq \bar{q}$ ,  $(\mathcal{FB})$ 
and  $\sum_{\omega} (\pi(a,\omega) - \pi(a',\omega)) x(a|\omega) q(\omega) \geq 0 \quad \forall a,a' \in A$ 

<sup>&</sup>lt;sup>18</sup>In particular, the timing of the two economies is the same. To see this, note that we could have equivalently written conditions (a) and (b) in Definition 1 as Equation (9), with the additional restriction that  $p(a, \omega) = p(\omega)$  for all  $(a, \omega)$ .

Compared to the notion of constrained efficiency (Definition 2), an unconstrained-efficient allocation (q, x) requires x to be optimal for the planner and not necessarily for the platform. Besides this, the planner's problem is the same. Indeed, notice that, in Definition 2, x was already required to be obedient since x solved  $\mathcal{P}_q$ . Therefore, the welfare of an unconstrained-efficient allocation is weakly higher than that of a constrained-efficient allocation:  $W^{\dagger} \geq W^{\circ}$ .

The next result shows that the equilibria of the Lindahl economy are unconstrained efficient.

**Proposition 5.** Let  $(p^*, z^*, q^*, x^*)$  be an equilibrium of the Lindahl economy. The allocation  $(q^*, x^*)$  is unconstrained efficient. Moreover, consumer welfare equals  $W^{\dagger}$ . Conversely, any unconstrained-efficient allocation  $(q^{\dagger}, x^{\dagger})$  can be supported as an equilibrium of the Lindahl economy.<sup>19</sup>

The richness of the price system allows the equilibria of this economy not only to avoid the inefficiency highlighted in Section 3, but to achieve unconstrained efficiency. For the same reasons as before, the platform still earns zero profit, so  $W^{\dagger}$  is also the consumer welfare. The following example illustrates how the richer price system helps induce efficient allocations.

Example of a Lindahl economy. Consider again the example of Section 3.1. Since  $\gamma_{\pi}=0$ , a mechanism x is optimal for the platform if and only if it is also optimal for the planner. Therefore, the constrained- and unconstrained-efficient allocations coincide, leading to a welfare of  $W^{\circ}=W^{\dagger}=\bar{r}+\bar{q}(1)(1+\gamma_{u}-2\bar{r})$ . Moreover, as in Section 3.1, the unconstrained-efficient allocation  $(q^{\dagger},x^{\dagger})$  is unique and given by  $q^{\dagger}(\omega)=\bar{q}(1)$  and  $x^{\dagger}(1|\omega)=1$  for all  $\omega$ . Recall that in the baseline economy all equilibria are inefficient. By contrast, we now discuss an equilibrium of the Lindahl economy and show it is unconstrained efficient. Let  $(p^{*},z^{*},q^{*},x^{*})$  be defined as follows. First, let  $(q^{*},x^{*})=(q^{\dagger},x^{\dagger})$ , i.e., the candidate equilibrium supports the unconstrained-efficient allocation. Second, for all  $\omega$ , let  $z^{*}(1,\omega)=\frac{\bar{q}(1)}{\bar{q}(\omega)}$  and  $z^{*}(2,\omega)=0$ , so that  $z^{*}$  and  $(q^{*},x^{*})$  clear the data markets. Finally, let prices be  $p(a=2,\omega)=0$ , for all  $\omega$ ,  $p^{*}(a=1,\omega=1)=\gamma_{u}+(1-\bar{r})$ , and  $p^{*}(a=1,\omega=2)=-(1-\bar{r})$ . To see that this is an equilibrium of the Lindahl economy, note that given prices  $p^{*}$ , type-1 consumers strictly prefer to sell their record in market  $(a=1,\omega=1)$ . Type-2 consumers, instead, are indifferent between not selling and selling in market  $(a=1,\omega=2)$ . Finally, the platform maximizes

<sup>&</sup>lt;sup>19</sup>Since an unconstrained-efficient allocation always exists, this result implies the existence of an equilibrium of the Lindahl economy.

 $(\gamma_u + 1 - \bar{r})(x(1|2)q(2) - x(1|1)q(1))$  subject to  $x(1|1)q(1) \ge x(1|2)q(2)$ . Therefore, the platform cannot make a positive payoff, and  $(q^*, x^*)$  achieves the maximum of 0 given  $p^*$ .

It is worth noting the crucial role of the richer price system. The price  $p^*(a=1,\omega=1)$  incorporates the positive externality that a low-type consumer generates when selling her record. This high price paid by the platform is financed by the high-type consumers, who have to pay to participate in the platform's mechanism. The platform uses their payments to acquire low-type records. That is, it is as if high-type consumers who participate in the platform's mechanism subsidize the participation of low-type consumers. Notice that the equilibrium exists even if  $p^*(a=1,\omega=2) < 0$ . Despite the negative price, the platform does not have an incentive to acquire an arbitrary quantity of such records. This is because, to fulfill the terms of trade for this market, the platform has to guarantee the merchant is willing to charge a low fee a=1 to the type-2 records it acquires.

### References

- Acemoglu, D., A. Makhdoumi, A. Malekian, and A. Ozdaglar (2022): "Too much data: Prices and inefficiencies in data markets," *American Economic Journal: Microeconomics*, 14, 218–256.
- Acquisti, A., C. Taylor, and L. Wagman (2016): "The Economics of Privacy," *Journal of Economic Literature*, 54, 442–92.
- Acquisti, A. and H. R. Varian (2005): "Conditioning prices on purchase history," *Marketing Science*, 24, 367–381.
- ALI, S. N., G. Lewis, and S. Vasserman (2022): "Voluntary Disclosure and Personalized Pricing," forthcoming, Review of Economic Studies.
- Arrow, K. J. (1969): "The Organization of Economic Activity: Issues Pertinent to the Choice of Market versus Non-Market Allocation," *The Analysis and Evaluation of Public Expenditures: the PPB System*, Joint Economic Committee, Congress of the United States, Washington, D.C., 47–64.
- Bergemann, D. and A. Bonatti (2019): "Markets for Information: An Introduction," *Annual Review of Economics*, 11, 85–107.
- ——— (2023): "Data, competition, and digital platforms," arXiv preprint arXiv:2304.07653.

- Bergemann, D., A. Bonatti, and T. Gan (2022a): "The economics of social data," *The RAND Journal of Economics*, 53, 263–296.
- Bergemann, D., B. Brooks, and S. Morris (2015): "The Limits of Price Discrimination," *American Economic Review*, 105 (3).
- BERGEMANN, D., J. CRÉMER, D. DINIELLI, C.-C. GROH, P. HEIDHUES, M. SCHAFER, M. SCHNITZER, F. M. S. MORTON, K. SEIM, AND M. SULLIVAN (2023): "Market design for personal data," Yale J. on Reg., 40, 1056–1120.
- BERGEMANN, D., T. HEUMANN, S. MORRIS, C. SOROKIN, AND E. WINTER (2022b): "Optimal Information Disclosure in Classic Auctions," *American Economic Review: Insights*, 4, 371–88.
- Bergemann, D. and S. Morris (2016): "Bayes Correlated Equilibrium and the Comparison of Information Structures in Games," *Theoretical Economics*, 11, 487–522.
- ——— (2019): "Information Design: A Unified Perspective," *Journal of Economic Literature*, 57(1), pp. 44-95).
- BERGEMANN, D. AND M. OTTAVIANI (2021): "Information markets and nonmarkets," in *Handbook of industrial organization*, Elsevier, vol. 4(1), 593–672.
- Bertsimas, D. and J. N. Tsitsiklis (1997): *Introduction to linear optimization*, vol. 6, Athena scientific Belmont, MA.
- Böнм, V. (1975): "On the continuity of the optimal policy set for linear programs," *SIAM Journal on Applied Mathematics*, 28, 303–306.
- Bonnisseau, J.-M., E. L. del Mercato, and P. Siconolfi (2023): "Existence of an equilibrium in arrowian markets for consumption externalities," *Journal of Economic Theory*, 209, 105638.
- Calzolari, G. and A. Pavan (2006): "On the Optimality of Privacy in Sequential Contracting," *Journal of Economic Theory*, 130, 168–204.
- CHEN, D. (2022): "The market for attention," Available at SSRN 4024597.
- CHOI, J. P., D.-S. JEON, AND B.-C. KIM (2019): "Privacy and personal data collection with information externalities," *Journal of Public Economics*, 173, 113–124.
- DOVAL, L. AND V. SKRETA (2023): "Purchase history and product personalization," *RAND Journal of Economics*, Forthcoming.
- Dworczak, P. (2020): "Mechanism design with aftermarkets: Cutoff mechanisms," *Econometrica*, 88, 2629–2661.
- Elliott, M., A. Galeotti, A. Koh, and W. Li (2022): "Market Segmentation through Information," *Working Paper*.

- FARBOODI, M., R. MIHET, T. PHILIPPON, AND L. VELDKAMP (2019): "Big data and firm dynamics," in *AEA papers and proceedings*, American Economic Association 2014 Broadway, Suite 305, Nashville, TN 37203, vol. 109, 38–42.
- Federal Trade Commission (2014): *Data Brokers: A Call for Transparency and Accountability*, A Report by the Federal Trade Commission, May.
- Galperti, S., A. Levkun, and J. Perego (2023): "The Value of Data Records," *Review of Economic Studies*, rdad044.
- Galperti, S. and J. Perego (2022): "Competitive Markets for Personal Data," *Available at SSRN 4309966*.
- ——— (2023): "Privacy and the Value of Data," AEA Papers and Proceedings, 113, 197–203.
- GOLDFARB, A. AND C. TUCKER (2023): The Economics of Privacy, NBER Conference Volume.
- HIDIR, S. AND N. VELLODI (2021): "Privacy, Personalization and Price Discrimination," *Journal of the European Economic Association*, 19, 1342–1363.
- Icнінаsні, S. (2021): "The economics of data externalities," *Journal of Economic Theory*, 196, 105316.
- Jones, C. I. and C. Tonetti (2020): "Nonrivalry and the Economics of Data," *American Economic Review*, 110, 2819–58.
- Kamenica, E. (2019): "Bayesian persuasion and information design," *Annual Review of Economics*, 11, 249–272.
- LAFFONT, J. J. (1976): "Decentralization with Externalities," *European Economic Review*, 359–375.
- Posner, E. and E. G. Weyl (2018): Radical Markets: Uprooting Capitalism and Democracy for a Just Society, Princeton University Press.
- ROCKAFELLAR, R. T. AND R. J.-B. Wets (2009): *Variational analysis*, vol. 317, Springer Science & Business Media.
- Taylor, C. R. (2004): "Consumer privacy and the market for customer information," *Rand Journal of Economics*.
- Varian, H. R. (2009): "Economic Aspects of Personal Privacy," in *Internet Policy and Economics*, ed. by W. H. Lehr and L. M. Pupillo, New York: Springer.
- Xu, W. And K. H. Yang (2023): "Informational Intermediation, Market Feedback, and Welfare Losses," *Working Paper*.

# **Appendix**

## **A** Proofs

### A.1 Proofs for Section 3

As a result of the linear specification of the platform's payoff,  $v(a, \omega) = \gamma_u u(a, \omega) + \gamma_\pi \pi(a, \omega)$ , we can explicitly compute V(q) and W(q) for all q, and characterize their supergradients. To do so, fix an arbitrary database  $q \geq 0$  and let  $w(q) \triangleq \sum_{\omega} \omega q(\omega)$  be the gains from trade and  $\pi_m(q) \triangleq \max_a \sum_{\omega} \pi(a, \omega) q(\omega)$  be the payoff of the uninformed merchant who only knows the composition of q.

As discussed in the second paragraph after Proposition 2 and shown in Figure 1, when the platform's payoff v has the linear specification, the optimal mechanism  $x^*$  either maximizes merchant's profit (when  $\gamma_{\pi} > \gamma_{u}$ ) or maximizes consumers' surplus (when  $\gamma_{u} \geq \gamma_{\pi}$ ). Therefore, we can explicitly compute the platform's value V and the planner's value W, depending on the relative magnitude of  $\gamma_{u}$  and  $\gamma_{\pi}$ :

$$V(q) = \begin{cases} \gamma_{\pi}w(q), & \gamma_{\pi} > \gamma_{u} \\ \gamma_{\pi}\pi_{m}(q) + \gamma_{u}(w(q) - \pi_{m}(q)), & \gamma_{\pi} \leq \gamma_{u} \end{cases}$$
(A.1)

$$W(q) = \begin{cases} \gamma_{\pi}w(q), & \gamma_{\pi} > \gamma_{u} \\ \gamma_{\pi}\pi_{m}(q) + (1 + \gamma_{u})(w(q) - \pi_{m}(q)), & \gamma_{\pi} \leq \gamma_{u} \end{cases}$$
(A.2)

Since w(q) is linear and  $\pi_m(q)$  is convex, V and W are concave functions and their supergradients are well-defined. Next, we characterize these supergradients. Let  $A_q \triangleq \arg\max_a \sum_{\omega} \pi(a,\omega) q(\omega)$  be the set of optimal monopoly prices under q. Let  $G_q \triangleq \{g \in \mathbb{R}^{\Omega} : g(\omega) = 0 \text{ if } q(\omega) > 0 \text{ and } g(\omega) \geq 0 \text{ if } q(\omega) = 0\}$ . Define  $\Phi_q(\gamma_u, \gamma_\pi) \subset \mathbb{R}^{\Omega}$  as

$$\Phi_{q}(\gamma_{u},\gamma_{\pi}) \triangleq \begin{cases} (\gamma_{\pi}\omega)_{\omega \in \Omega}, & \gamma_{\pi} > \gamma_{u} \\ cov\{(\gamma_{u}\omega + (\gamma_{\pi} - \gamma_{u})a\mathbb{1}(\omega \geq a))_{\omega \in \Omega} : a \in A_{q}\}, & \gamma_{\pi} \leq \gamma_{u} \end{cases} + G_{q},$$

where *cov* refers to convex hull and summation between two sets A and B is defined in the sense that  $A + B \triangleq \{a + b : a \in A, b \in B\}$ .

### Lemma A.1.

- 1. The set of supergradients of V(q) at q is  $\Phi_q(\gamma_u, \gamma_\pi)$ .
- 2. The set of supergradients of W(q) at q is  $\Psi_q = \Phi_q(\gamma_u, \gamma_\pi)$  if  $\gamma_\pi > \gamma_u$  and  $\Psi_q = \Phi_q(1 + \gamma_u, \gamma_\pi)$  if  $\gamma_\pi \leq \gamma_u$ .

*Proof.* Note that  $V, W : \mathbb{R}^{\Omega}_{+} \to \mathbb{R}$  defined in Equation (A.1) and (A.2) are linear combinations of functions w and  $-\pi_{m}$ . Therefore, in order to compute the supergradients of V and W, we only need to compute the supergradients of w and  $-\pi_{m}$ . To do this, we first extend w and  $-\pi_{m}$  from  $\mathbb{R}^{\Omega}_{+}$  to  $\mathbb{R}^{\Omega}$  by defining  $\tilde{w}, \tilde{\pi}_{m} : \mathbb{R}^{\Omega} \to \mathbb{R}$  as

$$\tilde{w}(q) = \sum_{\omega} \omega q(\omega) + \delta(q), \ -\tilde{\pi}_m(q) = \min_{a \in \Omega} -\pi(a, \omega)q(\omega) + \delta(q),$$

where  $\bar{\mathbb{R}} \triangleq [-\infty, \infty)$ ,  $\delta(q) : \mathbb{R}^{\Omega} \to \bar{\mathbb{R}}$  is defined by  $\delta(q) = 0$  if  $q \in \mathbb{R}_+^{\Omega}$  and  $\delta(q) = -\infty$  if  $q \notin \mathbb{R}_+^{\Omega}$ . It immediately follows that  $\partial w(q) = \partial \tilde{w}(q)$  and  $\partial(-\pi_m)(q) = \partial(-\tilde{\pi}_m)(q)$  for all  $q \in \mathbb{R}_+^{\Omega}$ .

Next, we compute the supergradients of  $\tilde{w}$  and  $-\tilde{\pi}_m$  on  $\mathbb{R}^{\Omega}_+$ . Note that  $\sum_{\omega} \omega q(\omega)$ ,  $\min_{a \in \Omega} -\pi(a, \omega)q(\omega)$ ,  $\delta(q)$  are upper-semicontinuous in q on  $\mathbb{R}^{\Omega}$  and finite on  $\mathbb{R}^{\Omega}_+$ . Therefore, we can invoke Corollary 10.9 and Example 7.27 in Rockafellar and Wets (2009) to derive that

$$\partial \tilde{w}(q) = \partial (\sum_{\omega} \omega q(\omega)) + \partial \delta(q), \ \partial (-\tilde{\pi}_m)(q) = \partial (\min_{a \in \Omega} -\pi(a, \omega)q(\omega)) + \partial \delta(q)$$

on  $q \in \mathbb{R}^{\Omega}_+$ . We have  $\partial(\sum_{\omega} \omega q(\omega)) = (\omega)_{\omega \in \Omega}$ . Moreover, Exercise 8.31 in Rockafellar and Wets (2009) implies that  $\partial(\min_{a \in \Omega} -\pi(a,\omega)q(\omega)) = cov\{(-\pi(a,\omega) = -a\mathbb{1}(\omega \geq a)_{\omega \in \Omega} : a \in A_q\}$ . It is left to compute  $\partial \delta(q)$ . By definition,  $g \in \mathbb{R}^{\Omega}$  is a supergradient of  $\delta$  at  $q \in \mathbb{R}^{\Omega}_+$  if and only if for all  $q' \in \mathbb{R}^{\Omega}_+$ :

$$\delta(q') - \delta(q) \le g \cdot (q' - q) \iff g \cdot (q' - q) \ge 0.$$

Therefore,  $\partial \delta(q) = G_q$  on  $\mathbb{R}^{\Omega}_+$ .

Finally, we compute  $\partial V(q)$  and  $\partial W(q)$ . If  $\gamma_{\pi} > \gamma_{u}$ , Equation (A.1) and (A.2) imply that

$$\partial V(q) = \partial W(q) = \gamma_{\pi} \partial w(q) = \{(\gamma_{\pi} \omega)_{\omega \in \Omega}\} + G_q = \Phi(\gamma_u, \gamma_{\pi}).$$

If  $\gamma_{\pi} \leq \gamma_{u}$ , Equation (A.1) implies that

$$\partial V(q) = \gamma_u \partial w(q) + (\gamma_u - \gamma_\pi) \partial (-\pi_m)(q)$$

$$=cov\{(\gamma_u\omega+(\gamma_\pi-\gamma_u)a\mathbb{1}(\omega\geq a))_{\omega\in\Omega}:a\in A_q\}+G_q=\Phi(\gamma_u,\gamma_\pi);$$

Equation (A.2) implies that

$$\partial W(q) = (1 + \gamma_u)\partial w(q) + (1 + \gamma_u - \gamma_\pi)\partial(-\pi_m)(q)$$
$$= cov\{((1 + \gamma_u)\omega + (\gamma_\pi - 1 - \gamma_u)a\mathbb{1}(\omega \ge a))_{\omega \in \Omega} : a \in A_q\} + G_q = \Phi(1 + \gamma_u, \gamma_\pi).$$

Note that  $\Phi_q$  and  $\Psi_q$  are singletons when  $A_q$  is a singleton and  $q(\omega) > 0$  for all  $\omega$ , so the supergradients are generically unique. Moreover, Theorem 5.2 in Bertsimas and Tsitsiklis (1997) shows that the set of supergradients of the value function of a linear program coincides with the set of solutions of its dual. Therefore,  $\Phi_q$  is also the set of solutions to the dual problem of  $\mathcal{P}_q$ :

$$\begin{split} & \min_{\phi,\lambda} \quad \sum_{\omega} \phi(\omega) q(\omega) \\ & \text{such that} \quad \phi(\omega) \geq v(a,\omega) + \sum_{\hat{a}} (\pi(a,\omega) - \pi(\hat{a},\omega)) \lambda(\hat{a}|a) \qquad \forall a,\omega \qquad (\mathcal{P}_q'(v)) \\ & \text{and} \quad \lambda(\hat{a}|a) \geq 0 \qquad \forall \hat{a},a \end{split}$$

We will use this fact in the following proofs.

It is with noting that Lemma A.1 extends Proposition 2 in Galperti et al. (2023). In their setting, Galperti et al. (2023) characterize the solution to problem  $\mathcal{P}'_q(v)$  when it is unique, i.e., when  $A_q$  is a singleton and  $q(\omega) > 0$  for all  $\omega$ . We complement their results by providing a full characterization of the solutions even when it is not unique, and also by characterizing the supergradients of W, i.e., the marginal value of data to the planner under the incentive constraint of the platform.

With these results in mind, we proceed to prove Proposition 1 and Lemma 1.

**Proof of Proposition 1**. Suppose  $\gamma_{\pi} > \gamma_{u}$ . In this case, x as a solution to  $\mathcal{P}_{q}$  reveals all records in the database. Given this, it is easy to see that in this case the planner's solution to problem  $\mathcal{SB}$  (see Definition 2) is  $q(\omega) = 0$  if  $\gamma_{\pi}\omega < r(\omega)$ ,  $q(\omega) = \bar{q}(\omega)$  if  $\gamma_{\pi}\omega > r(\omega)$ , and  $q(\omega) \in [0, \bar{q}(\omega)]$  if  $\gamma_{\pi}\omega = r(\omega)$ . Taken together, (q, x) is constrained efficient if and only if  $q(\omega) > 0$  implies  $\gamma_{\pi}\omega \geq r(\omega)$  while  $q(\omega) < \bar{q}(\omega)$  implies  $\gamma_{\pi}\omega \leq r(\omega)$ , and x solves  $\mathcal{P}_{q}$ . From Lemma A.1 we know  $\psi_{q} \in \Psi_{q}$ , where  $\psi_{q}(\omega) = \gamma_{\pi}\omega$ ,  $\forall \omega$ . This proves

that if (q,x) is constrained efficient, then there exists  $\psi_q \in \Psi_q$  such that  $q(\omega) > 0$  implies  $\psi_q \ge r(\omega)$  while  $q(\omega) < \bar{q}(\omega)$  implies  $\psi_q \le r(\omega)$ . Conversely, suppose x solves  $\mathcal{P}_q$  and for some  $\psi_q \in \Psi_q$  we have  $\psi_q(\omega) \ge r(\omega)$  if  $q(\omega) > 0$  and  $\psi_q(\omega) \le r(\omega)$  if  $q(\omega) < \bar{q}(\omega)$ . By Lemma A.1 we know if  $q(\omega) > 0$  then  $\psi_q(\omega) = \gamma_\pi \omega$ ; if  $q(\omega) < \bar{q}(\omega)$  then  $\psi_q(\omega) \ge \gamma_\pi \omega$ . Therefore, this means  $q(\omega) > 0$  implies  $\gamma_\pi \omega \ge r(\omega)$  while  $q(\omega) < \bar{q}(\omega)$  implies  $\gamma_\pi \omega \le \psi_q(\omega) \le r(\omega)$ , which means (q,x) is constrained efficient.

Now, suppose  $\gamma_{\pi} \leq \gamma_{u}$ . In this case, the constraint that requires x to be sequentially optimal for the platform can be relaxed, and the planner's problem  $\mathcal{SB}$  coincides with  $\mathcal{FB}$  (see Definition 4). Thus, a data allocation (q, x) is constrained efficient if and only if it solves problem  $\mathcal{FB}$ . The dual problem of  $\mathcal{FB}$  can be formulated as:

$$\begin{split} & \min_{\mu,\lambda} \quad \sum_{\omega} \mu(\omega) \bar{q}(\omega) \\ & \text{such that} \quad \mu(\omega) \geq \gamma_{\pi} \pi(a,\omega) + (1+\gamma_{u}) u(a,\omega) + \sum_{\hat{a}} (\pi(a,\omega) - \pi(\hat{a},\omega)) \lambda(\hat{a}|a) \qquad \forall a,\omega \\ & \text{and} \quad \mu(\omega) \geq r(\omega) \qquad \forall \omega \\ & \text{and} \quad \lambda(\hat{a}|a) \geq 0 \qquad \forall \hat{a},a \end{split} \tag{$\mathcal{FB}'$}$$

We first show the "only if" direction of the proposition. Take any efficient allocation (q, x). Then the planner's value under (q, x) can also be written as:

$$\begin{split} \sum_{\omega} \Big( \bar{q}(\omega) - q(\omega) \Big) r(\omega) + & \max_{\chi \in \mathbb{R}_{+}^{A \times \Omega}} & \sum_{a, \omega} \Big( v(a, \omega) + u(a, \omega) \Big) \chi(a, \omega) \\ & \text{such that} & \sum_{a} \chi(a, \omega) = q(\omega), \quad \forall \omega \in \Omega \\ & \text{and} & \sum_{\omega} \Big( \pi(a, \omega) - \pi(\hat{a}, \omega) \Big) \chi(a, \omega) \geq 0 \qquad \forall \ a, \hat{a} \in A \end{split}$$

Let  $(\mu, \lambda)$  be a solution to problem  $\mathcal{FB}'$ . Then  $(\mu, \lambda)$  is also feasible for problem  $\mathcal{P}'_q(v + u)$ , which is the dual of the maximization problem above. Moreover, by strong duality, the planner's value can also be written as  $\sum_{\omega} \mu(\omega) \bar{q}(\omega)$ . Therefore, we have:

$$\sum_{\omega} \mu(\omega) \bar{q}(\omega) \leq \sum_{\omega} \left( \bar{q}(\omega) - q(\omega) \right) r(\omega) + \sum_{\omega} \mu(\omega) q(\omega).$$

Meanwhile, since  $\mu(\omega) \ge r(\omega)$ , we must also have the other direction of the inequality, so we conclude:

$$\sum_{\omega} \mu(\omega) \bar{q}(\omega) = \sum_{\omega} \left( \bar{q}(\omega) - q(\omega) \right) r(\omega) + \sum_{\omega} \mu(\omega) q(\omega).$$

This equality has two implications. First,  $(\mu, \lambda)$  achieves the minimum of  $\mathcal{P}_q'(v+u)$ , so  $\mu \in \Psi_q$ . Second,  $\mu(\omega) = r(\omega)$  whenever  $q(\omega) < \bar{q}(\omega)$ . Letting  $\psi_q = \mu$ , we conclude the proof of the "only if" direction, i.e., if  $q(\omega) > 0$  then  $\psi_q(\omega) \ge r(\omega)$  (which is always true by construction) and if  $q(\omega) < \bar{q}(\omega)$  then  $\psi_q(\omega) \le r(\omega)$ .

Next, we prove the "if" direction. Let  $(\psi, \lambda)$  be a solution to  $\mathcal{P}'_q(v+u)$  such that  $\psi(\omega) \geq r(\omega)$  if  $q(\omega) > 0$  and  $\psi(\omega) \leq r(\omega)$  if  $q(\omega) < \bar{q}(\omega)$ . Let  $\mu \triangleq \max\{\psi, r\}$ . Then,  $(\mu, \lambda)$  is feasible for  $\mathcal{FB}'$  and, thus,  $\sum_{\omega} \mu(\omega) \bar{q}(\omega) \geq W^{\circ}$  by weak duality. Next, we argue (q, x) achieves  $W^{\circ}$ , where x is a solution to (3). Since  $(\psi, \lambda)$  is a solution to  $\mathcal{P}'_q(v+u)$  and x is a solution to (3), by strong duality, we know that under (q, x), the planner's value is

$$W(q,x) = \sum_{\omega} q(\omega)\psi(\omega) + \sum_{\omega} r(\omega)(\bar{q}(\omega) - q(\omega)).$$

When  $q(\omega) = \bar{q}(\omega)$ , we have  $\psi(\omega) \geq r(\omega)$ , and thus  $\psi(\omega)q(\omega) = \mu(\omega)\bar{q}(\omega)$ ; when  $q(\omega) = 0$ ,  $\psi(\omega) \leq r(\omega)$  and thus  $r(\omega)(\bar{q}(\omega) - q(\omega)) = \mu(\omega)\bar{q}(\omega)$ ; when  $0 < q(\omega) < \bar{q}(\omega)$ , we have  $\psi(\omega) = r(\omega) = \mu(\omega)$ . Therefore,

$$W(q,x) = \sum_{\omega} q(\omega)\psi(\omega) + \sum_{\omega} r(\omega)(\bar{q}(\omega) - q(\omega)) = \sum_{\omega} \mu(\omega)\bar{q}(\omega).$$

This means that (q, x) achieves  $W^{\circ}$  and thus (q, x) is constrained efficient.

**Lemma A.2.** Fix an arbitrary v.<sup>20</sup>  $q \leq \bar{q}$  solves the platform's problem in the first period (Problem (1)) if and only if  $p \in \Phi_q$ .

*Proof.* We first observe that the platform's problem in the first period (Problem (1)) is equivalent to choosing (q, x) given price p, or in other words, choosing  $\chi(a, \omega) = \chi(a|\omega)q(\omega)$  without any feasibility constraint. This is because in the first period the platform can choose any  $q \ge 0$ . Therefore, its dual problem can be formulated as:

$$\min_{\lambda} \quad 0$$
such that 
$$\sum_{\hat{a}} (\pi(\hat{a}, \omega) - \pi(a, \omega)) \lambda(\hat{a}|a) \ge v(a, \omega) - p(\omega) \qquad \forall a, \omega \qquad (A.3)$$
and 
$$\lambda(\hat{a}|a) > 0 \qquad \forall \hat{a}, a$$

In other words, the platform's optimal payoff is zero if (A.3) is feasible, and it is infinite otherwise.

 $<sup>^{20}</sup>$ We emphasize that in this lemma v does not need to be a linear combination of u and  $\pi$ .

To show the "only if" direction, note that if  $q \leq \bar{q}$  solves (1), then  $V(q) = \sum_{\omega} p(\omega)q(\omega)$ . If instead  $V(q) > \sum_{\omega} p(\omega)q(\omega)$ , the platform can choose 2q to achieve a higher payoff, which is a contradiction. By strong duality,  $V(q) = \sum_{\omega} p(\omega)q(\omega)$  implies that Problem (A.3) is feasible. Take any feasible solution  $\lambda$ , and consider  $(\phi, \lambda)$  where  $\phi = p$ . Next we argue  $(\phi, \lambda)$  is an optimal solution to  $\mathcal{P}'_q(v)$ . Suppose not, then since  $(\phi, \lambda)$  is feasible to  $\mathcal{P}'_q(v)$ , we must have  $V(q) < \sum_{\omega} \phi(\omega)q(\omega)$ , but this contradicts  $V(q) - \sum_{\omega} p(\omega)q(\omega) = 0$ .

To show the "if" direction, suppose  $(p,\lambda)$  is an optimal solution to  $\mathcal{P}_q'(v)$ . This means Problem (A.3) is feasible. Therefore, the platform's optimal payoff is 0 given p. By strong duality, we have  $V(q) = \sum_{\omega} p(\omega)q(\omega)$ . This means that q gives the platform a payoff of 0. Therefore, q is a solution to the platform's problem (1) given price p.

**Proof of Lemma 1.** Since in any equilibrium  $(p^*, z^*, q^*, x^*)$ , we must have  $q^* \leq \bar{q}$  and  $q^*$  solves (1) given  $p^*$ , Lemma A.2 then implies that  $p^* \in \Phi_{q^*}$ , which means  $p^*$  is a supergradient of V at  $q^*$ .

**Proof of Proposition 2.** We prove the case of  $\gamma_{\pi} > \gamma_{u}$  here. The negative results for the case of  $\gamma_{\pi} \leq \gamma_{u}$  is illustrated by Section 3.1, Corollary 1, and Corollary A.1. Consider any equilibrium  $(p^{*}, z^{*}, q^{*}, x^{*})$ . We know that  $q^{*}(\omega) > 0$  implies  $p^{*}(\omega) \geq r(\omega)$  and  $q^{*}(\omega) < \bar{q}(\omega)$  implies  $p^{*}(\omega) \leq r(\omega)$  since  $x^{*}$  perfectly reveals all records in  $q^{*}$ . By Lemma 1, we know  $p^{*} \in \Phi_{q^{*}}$ . Since  $\gamma_{\pi} > \gamma_{u}$ ,  $\Psi_{q^{*}} = \Phi_{q^{*}}$  by Lemma A.1. Taking  $\psi_{q^{*}} = p^{*}$  we conclude that  $(q^{*}, x^{*})$  is constrained efficient by Proposition 1.

**Proof of Corollary 1.** Let  $(p^*, z^*, q^*, x^*)$  be an equilibrium. If  $q^* = 0$ , then there is no trade, which is inefficient by assumption. Therefore, assume  $q^* \neq 0$ . Next, we argue  $q^*(\underline{\omega}) = 0$ . By Lemma 1 and Lemma A.1,  $p^*(\underline{\omega}) \leq \gamma_u \underline{\omega}$ . Since a type- $\underline{\omega}$  consumer will earn a zero trading surplus by selling her record and by assumption  $\gamma_u \underline{\omega} < r(\underline{\omega})$ , she strictly prefer not to do so. However, this is inefficient since, given  $q^* \neq 0$  and Lemma A.1,  $\psi_{q^*}(\underline{\omega}) \geq (1 + \gamma_u)\underline{\omega} > r(\underline{\omega})$ . Invoking Proposition 1 we conclude the proof.

Next, we introduce another set of sufficient conditions under which all equilibria of the competitive economy are inefficient. Let  $\omega_1$  be the highest type in  $\Omega$  and  $\omega_2$  be the second-highest one. In Corollary A.1, we assume that, if all type- $\omega_1$  consumers sell their records, the merchant strictly prefers to set a fee  $a=\omega_1$ . This is guaranteed by the requirement that  $\bar{q}(\omega_1)\omega_1>\omega\sum_{\omega'\geq\omega}\bar{q}(\omega')$  for all  $\omega<\omega_1$ . Additionally, we assume that if the uninformed

merchant fee is  $a_q = \omega_1$ , the platform wants to pool type- $\omega_2$  and type- $\omega_1$  consumers, to induce a lower fee from the merchant. This is guaranteed by the requirement that  $r(\omega_2) < \gamma_u \omega_2$ . Under these assumptions, we show that for a range of  $r(\omega_1)$ , type- $\omega_1$  consumers strictly prefer to sell their records while it is not efficient. This corresponds to Case 2 in the example of Section 3.1.

**Corollary A.1.** Let  $\gamma_{\pi} \leq \gamma_{u}$ . Assume  $\bar{q}(\omega_{1})\omega_{1} > \omega \sum_{\omega' \geq \omega} \bar{q}(\omega')$  for all  $\omega < \omega_{1}$  and  $r(\omega_{2}) < \gamma_{u}\omega_{2}$ . Then no equilibrium is constrained efficient if  $\gamma_{\pi}\omega_{1} < r(\omega_{1}) < \gamma_{\pi}\omega_{1} + \omega_{1} - \omega_{2}$ .

*Proof.* Let  $(p^*, z^*, q^*, x^*)$  be a competitive equilibrium. We want to show it is not constrained efficient. Note that since  $\gamma_{\pi} \leq \gamma_{u}$ , solving  $\mathcal{SB}$  is equivalent to solving  $\mathcal{FB}$  (since the sequentially rationality constraint of the platform can be relaxed). There are four cases to discuss.

• Case 1: Suppose  $q^*(\omega_2) < \bar{q}(\omega_2), q^*(\omega_1) < \bar{q}(\omega_1)$ . To the contrary, suppose this is constrained efficient. Then by Proposition 1, there exists  $\psi \in \Psi_{q^*}$  such that  $\psi(\omega_1) \leq r(\omega_1)$  and  $\psi(\omega_2) \leq r(\omega_2)$ . By Lemma A.1, we know that  $\psi(\omega_1) \geq \sum_{a \in \Omega} \lambda_a((1 + \gamma_u)\omega_1 + (\gamma_\pi - 1 - \gamma_u)a)$ , where  $(\lambda_a)_{a \in \Omega}$  are weights satisfying  $\lambda_a \geq 0, \sum_{a \in \Omega} = 1$ . We have:

$$\lambda_{\omega_1} \gamma_{\pi} \omega_1 + (1 - \lambda_{\omega_1})((1 + \gamma_u)\omega_1 + (\gamma_{\pi} - 1 - \gamma_u)\omega_2)$$

$$\leq \sum_{a \in \Omega} \lambda_a ((1 + \gamma_u)\omega_1 + (\gamma_{\pi} - 1 - \gamma_u)a)$$

$$\leq \psi(\omega_1) \leq r(\omega_1) < \gamma_{\pi} \omega_1 + \omega_1 - \omega_2,$$

where the first inequality follows because  $(1 + \gamma_u)\omega_1 + (\gamma_\pi - 1 - \gamma_u)a$  is decreasing in a and the last inequality follows by assumption. Comparing the left-hand side and the right-hand side, we conclude  $\lambda_{\omega_1} > \frac{\gamma_u - \gamma_\pi}{1 + \gamma_u - \gamma_\pi}$ . However, this then implies:

$$\psi(\omega_2) \ge \lambda_{\omega_1} (1 + \gamma_u) \omega_2 + (1 - \lambda_{\omega_1}) \gamma_{\pi} \omega_2$$

$$> \frac{\gamma_u - \gamma_{\pi}}{1 + \gamma_u - \gamma_{\pi}} (1 + \gamma_u) \omega_2 + \frac{1}{1 + \gamma_u - \gamma_{\pi}} \gamma_{\pi} \omega_2$$

$$= \gamma_u \omega_2 > r(\omega_2),$$

which is a contradiction to  $\psi(\omega_2) \leq r(\omega_2)$ . Therefore,  $(q^*, x^*)$  is not constrained efficient.

- Case 2: Suppose  $q^*(\omega_2) < \bar{q}(\omega_2), q^*(\omega_1) = \bar{q}(\omega_1)$ . In this case, since by assumption we have  $\bar{q}(\omega_1)\omega_1 > \omega \sum_{\omega' \geq \omega} \bar{q}(\omega')$  for all  $\omega < \omega_1$ , the unique uninformed merchant fee under  $q^*$  is  $a_{q^*} = \omega_1$ . Therefore, by Lemma A.1,  $\psi_{q^*}(\omega_2) \geq (\gamma_u + 1)\omega_2 > r(\omega_2)$ . By Proposition 1, we conclude that  $(q^*, x^*)$  is not constrained efficient.
- Case 3:  $q^*(\omega_2) = \bar{q}(\omega_2), q^*(\omega_1) < \bar{q}(\omega_1)$ . We argue this is also inefficient. First, if  $q^*(\omega_1) = 0$ , then  $\psi_{q^*}(\omega_1) \geq \gamma_\pi \omega_2 + (\gamma_u + 1)(\omega_1 \omega_2) = \gamma_\pi \omega_1 + (\gamma_u \gamma_\pi + 1)(\omega_1 \omega_2) > r(\omega_1)$ . Therefore, by Proposition 1,  $(q^*, x^*)$  is not constrained efficient. Second, suppose  $q^*(\omega_1) > 0$  and  $x^*(\omega_1|\omega_1) > 0$ . Let  $\chi^*(a,\omega) \triangleq x^*(a|\omega)q^*(\omega)$  be the joint distribution over a and  $\omega$  induced by  $(q^*, x^*)$ . Note that by obedience we must have  $\chi^*(a,\omega) > 0$  only if  $a \in \Omega$ . Next, we construct a new allocation (q,x) which is feasible for  $\mathcal{FB}$  but have  $\mathcal{W}(q,x) > \mathcal{W}(q^*,x^*)$ . We construct such (q,x) by directly constructing its induced joint distribution  $\chi(a,\omega)$ . Let  $\chi(a,\omega) = \chi^*(a,\omega)$  for  $a \neq \omega_1$  and  $\chi(\omega_1,\omega) = 0$  for all  $\omega \in \Omega$ . Let (q,x) denote an allocation inducing  $\chi$ . Note that (q,x) is still feasible for  $\mathcal{FB}$  because  $q(\omega) = \sum_{a \in \Omega} \chi(a,\omega) \leq q^*(\omega) \leq \bar{q}(\omega)$  and  $\chi(a,\omega) = \chi^*(a,\omega)$  for all  $a \in \Omega \setminus \{\omega_1\}$ ; for  $a = \omega_1$ , since  $\chi(a,\omega) = 0$  for all  $\omega$ , obedience is satisfied vacuously. We have:

$$\mathcal{W}(q,x) - \mathcal{W}(q^*,x^*) = \sum_{\omega \in \Omega} r(\omega) \chi^*(\omega_1,\omega) - \chi^*(\omega_1,\omega_1) \gamma_\pi \omega_1 > 0,$$

where the inequality follows from the assumption that  $r(\omega_1) > \gamma_\pi \omega_1$ . Therefore, we conclude  $(q^*, x^*)$  is not constrained efficient. Third, suppose  $x^*(\omega_1|\omega_1) = 0$ . Then,  $\sum_{a \in \Omega} x^*(a|\omega_1) u(a,\omega_1) \geq \omega_1 - \omega_2$ , and by Lemma 1 and Lemma A.1 we have  $p^*(\omega_1) \geq \gamma_\pi \omega_1$ . Together, since by assumption  $\omega_1 - \omega_2 + \gamma_\pi \omega_1 > r(\omega_1)$ , it must be that  $q^*(\omega_1) = \bar{q}(\omega_1)$ , a contradiction. Summarizing the three scenarios, we conclude that  $(q^*, x^*)$  cannot be constrained efficient.

• Case 4: Suppose  $q^*(\omega_2) = \bar{q}(\omega_2)$  and  $q^*(\omega_1) = \bar{q}(\omega_1)$ . In this case, since by assumption we have  $\bar{q}(\omega_1)\omega_1 > \omega \sum_{\omega' \geq \omega} \bar{q}(\omega')$  for all  $\omega < \omega_1$ , the unique uninformed merchant fees is  $a_{q^*} = \omega_1$ . Therefore, we have  $\psi_{q^*}(\omega_1) = \gamma_\pi \omega_1 < r(\omega_1)$ . By Proposition 1,  $(q^*, x^*)$  is inefficient.

### A.2 Proofs for Section 4

Next, we provide proofs for results in Section 4. Note that none of the results in that section require the assumption that the platform's payoff is linear in u and  $\pi$ .

**Proof of Proposition 3.** ("If" Direction). Let  $(q^{\circ}, x^{\circ})$  be constrained efficient. First, we argue that  $(q^{\circ}, x^{\circ})$  is a solution of a relaxed version of the data union's problem, which we obtain by discarding the consumers' participation constraints. Additionally, note that the constraint  $\sum_{\omega} \bar{q}(\omega) p(\omega) \leq V(q)$  must hold with equality and, thus, can be substituted into the data union's objective. By doing so, prices p do not appear in the relaxed problem, which can be written as

$$\max_{(q,x)} \sum_{a,\omega} \left( v(a,\omega) + u(a,\omega) \right) x(a|\omega) q(\omega) + \sum_{\omega} (\bar{q}(\omega) - q(\omega)) r(\omega)$$
such that  $q \leq \bar{q}$ ,
and  $x$  solves  $\mathcal{P}_q$ .

Note that this relaxed problem coincides with the planner's problem  $\mathcal{SB}$ . Thus, any constrained-efficient allocation must yield a value that is weakly higher than the value of the data union's problem. To complete the proof, we find prices  $p^{\circ}$  such that, given  $(q^{\circ}, x^{\circ})$ , the consumers' participation constraints are satisfied and  $\sum_{\omega} \bar{q}(\omega)p(\omega) = V(q)$ .

To this end, let  $p^{\circ}(\omega) = \tilde{p}(\omega) + t(\omega)$  with  $\tilde{p}(\omega) = \frac{q^{\circ}(\omega)}{\bar{q}(\omega)} \Big( r(\omega) - \sum_a x^{\circ}(a|\omega) u(a,\omega) \Big)$ , where  $t(\omega)$  will be pinned down later. If  $t(\omega) = 0$ , all type- $\omega$  consumers would be indifferent between joining the union or not, and in particular,  $z(\omega) = 1$  is optimal. In this case, the union's budget is:

$$G(q^{\circ}, x^{\circ}) = V(q^{\circ}) - \sum_{\omega} \bar{q}(\omega) \tilde{p}(\omega)$$
$$= \sum_{a, \omega} \Big( v(a, \omega) + u(a, \omega) \Big) x^{\circ}(a|\omega) q^{\circ}(\omega) - \sum_{\omega} q^{\circ}(\omega) r(\omega).$$

Since  $(q^{\circ}, x^{\circ})$  is constrained efficient,  $G(q^{\circ}, x^{\circ}) \geq 0$ . To see this, we add  $\sum_{\omega} \bar{q}(\omega) r(\omega)$  on both sides of this inequality. On the left-hand side, we obtain  $\mathcal{W}(q^{\circ}, x^{\circ}) = W^{\circ} \geq R = \sum_{\omega} \bar{q}(\omega) r(\omega)$ , which is the right-hand side.

Finally, to guarantee  $\sum_{\omega} \bar{q}(\omega) p^{\circ}(\omega) = V(q^{\circ})$ , we can uniformly distribute  $G(q^{\circ}, x^{\circ})$  to the consumers. Specifically, let  $t(\omega) = G(q^{\circ}, x^{\circ})$  (recall that  $\sum_{\omega} \bar{q}(\omega) = 1$ ). Therefore, since  $z(\omega) = 1$  was optimal under  $\tilde{p}(\omega)$ , it is still optimal under  $p^{\circ}(\omega) \geq \tilde{p}(\omega)$ .

Thus, we constructed a profile  $(p^{\circ}, q^{\circ}, x^{\circ})$  that is feasible for the data union. Moreover, since  $(q^{\circ}, x^{\circ})$  solves the relaxed problem, it must be that  $(p^{\circ}, q^{\circ}, x^{\circ})$  solves the data union's problem.

("Only If" Direction). Let  $(p^*, q^*, x^*)$  be a solution to the data union's problem but suppose it is not constrained efficient. Let  $(q^\circ, x^\circ)$  be a constrained-efficient allocation. We have that:

$$\sum_{a,\omega} \left( v(a,\omega) + u(a,\omega) \right) x^*(a|\omega) q^*(\omega) - \sum_{\omega} q^*(\omega) r(\omega)$$

$$< \sum_{a,\omega} \left( v(a,\omega) + u(a,\omega) \right) x^{\circ}(a|\omega) q^{\circ}(\omega) - \sum_{\omega} q^{\circ}(\omega) r(\omega) = W^{\circ}$$
(A.4)

By the "only if" direction, we know there exist  $p^{\circ}$  such that  $(p^{\circ}, q^{\circ}, x^{\circ})$  is feasible for the data union and achieves  $W^{\circ}$ . This contradicts  $(p^*, q^*, x^*)$  being a solution to the data union's problem.

**Proof of Proposition 4.** As explained in the main text, given a constrained efficient allocation  $(q^{\circ}, x^{\circ})$ , we construct  $\tau^*, p^*, z^*$  as follows. Take any  $p^* \in \Phi_{q^{\circ}}$ , which by Lemma A.1 is nonempty. Let  $\tau^*$  be defined as in Equation (8) and define  $z^*(\omega) \triangleq q^{\circ}(\omega)/\bar{q}(\omega)$  for all  $\omega$ . Next, we check that  $(p^*, z^*, q^{\circ}, x^{\circ})$  is an equilibrium of the economy with taxation  $\tau^*$ . First, since  $p^* \in \Phi_{q^{\circ}}$ ,  $q^{\circ}$  solves the platform's problem in the first period by Lemma A.2. Second, since  $(q^{\circ}, x^{\circ})$  is constrained efficient,  $x^{\circ}$  must solve  $\mathcal{P}_{q^{\circ}}$ , i.e., the platform's problem in the second period. Third, for any  $z(\omega) \in [0,1]$ , the payoff of an  $\omega$ -consumer is  $z(\omega)(p^*(\omega) + \sum_{q} x^{\circ}(a|\omega)u(a,\omega) - \tau^*(\omega)) + (1-z(\omega))r(\omega) = r(\omega)$ , so  $z^*(\omega)$  solves the consumer's problem. Finally, by definition we have  $z^*(\omega)\bar{q}(\omega) = q^{\circ}(\omega)$  for all  $\omega$ , so data markets clear. Therefore,  $(p^*, z^*, q^{\circ}, x^{\circ})$  is an equilibrium of the economy with taxation  $\tau^*$ .

In this equilibrium, the government's budget balance is:

$$\begin{split} \sum_{\omega} \tau^*(\omega) q^{\circ}(\omega) &= \sum_{\omega} (p^*(\omega) + \sum_{a} x^{\circ}(a|\omega) u(a,\omega)) q^{\circ}(\omega) - \sum_{\omega} r(\omega) q^{\circ}(\omega) \\ &= \sum_{\omega} \sum_{a} x^{\circ}(a|\omega) (v(a,\omega) + u(a,\omega)) q^{\circ}(\omega) - \sum_{\omega} r(\omega) q^{\circ}(\omega) \\ &= \mathcal{W}(q^{\circ}, x^{\circ}) - R = W^{\circ} - R \ge 0, \end{split}$$

where the second equality follows because in equilibrium the platform has a zero payoff and the last equality follows because  $(q^{\circ}, x^{\circ})$  is constrained efficient. Therefore, the government does not run a budget deficit.

**Proof of Proposition 5**. ("Only If" Direction). Let  $(p^*, z^*, q^*, x^*)$  be a Lindahl equilibrium. We first prove that  $(q^*, x^*)$  must solve  $\mathcal{FB}$ . Since  $(q^*, x^*)$  solves the platform's problem (9), we have that

$$\sum_{a,\omega} v(a,\omega) x^*(a|\omega) q^*(\omega) - \sum_{a,\omega} v(a,\omega) x(a|\omega) q(\omega)$$

$$\geq \sum_{a,\omega} p^*(a,\omega) x^*(a|\omega) q^*(\omega) - \sum_{a,\omega} p^*(a,\omega) x(a|\omega) q(\omega)$$
(A.5)

for all (q, x) that satisfies obedience. Similarly, by the maximization problem of type- $\omega$  consumers, we get

$$\sum_{a} u(a,\omega)z^{*}(a,\omega) + r(\omega)\left(1 - \sum_{a} z^{*}(a,\omega)\right) - \sum_{a} u(a,\omega)z(a,\omega) - r(\omega)\left(1 - \sum_{a,} z(a,\omega)\right)$$

$$\geq -\sum_{a} p^{*}(a,\omega)z^{*}(a,\omega) + \sum_{a} p^{*}(a,\omega)z(a,\omega)$$

for all  $z(a, \omega) \in \mathbb{R}_+^A$  such that  $\sum_a z(a, \omega) \leq 1$ . Summing over consumers of the same type and across types, we get that for all (q, x) such that  $q \leq \bar{q}$ :

$$\sum_{a,\omega} u(a,\omega)x^{*}(a|\omega)q^{*}(\omega) - \sum_{\omega} r(\omega)q^{*}(\omega) - \sum_{a,\omega} u(a,\omega)x(a|\omega)q(\omega) + \sum_{\omega} r(\omega)q(\omega)$$

$$\geq -\sum_{a,\omega} p^{*}(a,\omega)x^{*}(a|\omega)q^{*}(\omega) + \sum_{a,\omega} p^{*}(a,\omega)x(a|\omega)q(\omega).$$
(A.6)

Equations (A.5) and (A.6) jointly imply that for all (q, x) satisfying feasibility and obedience:

$$\sum_{a,\omega} (v(a,\omega) + u(a,\omega)) x^*(a|\omega) q^*(\omega) - \sum_{\omega} r(\omega) q^*(\omega)$$

$$\geq \sum_{a,\omega} (v(a,\omega) + u(a,\omega)) x(a|\omega) q(\omega) - \sum_{\omega} r(\omega) q(\omega).$$

Therefore,  $(q^*, x^*)$  solves  $\mathcal{FB}$ .

("If" Direction). We now prove that for any allocation  $(q^{\dagger}, x^{\dagger})$  that is unconstrained efficient, there is a  $(p^*, z^*)$  such that  $(p^*, z^*, q^{\dagger}, x^{\dagger})$  is a Lindahl equilibrium. First of all, notice that  $\mathcal{FB}$  admits an optimal solution. Second, we can define  $p^*(a, \omega) = r(\omega) - u(a, \omega)$  for all  $a, \omega$ , so that each  $\omega$  consumer is indifferent across all possible  $z(\cdot, \omega)$  and we can therefore choose  $z^*$  such that  $z^*(\cdot, \omega)\bar{q}(\omega) = x^{\dagger}(\cdot|\omega)q^{\dagger}(\omega)$ .

We can equivalently rewrite  $\mathcal{FB}$  in terms of  $\chi$ :

$$\max_{\chi \in \mathbb{R}_{+}^{A \times \Omega}} \sum_{a, \omega} \Big( v(a, \omega) + u(a, \omega) \Big) \chi(a, \omega) + \sum_{\omega} \Big( \bar{q}(\omega) - \sum_{a} \chi(a, \omega) \Big) r(\omega)$$

such that 
$$\sum_{a}\chi(a,\omega)\leq \bar{q}(\omega)$$
,  $\forall \omega\in\Omega$  and  $\sum_{\omega}\left(\pi(a,\omega)-\pi(\hat{a},\omega)\right)\chi(a,\omega)\geq 0$   $\forall \ a,\hat{a}\in A$ 

Since  $(q^{\dagger}, x^{\dagger})$  is unconstrained efficient, we know  $\chi^{\dagger}(a, \omega) \triangleq x^{\dagger}(a|\omega)q^{\dagger}(\omega)$  solves the above problem. Define  $z^{*}(a, \omega) = \chi^{\dagger}(a, \omega)/\bar{q}(\omega)$ . Since  $\chi^{\dagger}$  is an optimal solution, by strong duality, we know its dual, which is defined by  $\mathcal{P}'_{\bar{q}}(v+u)$ , admits an optimal solution  $(\mu^{*}(\omega), \lambda^{*}(\hat{a}|a))$ . Define  $p^{*}(a, \omega) = \mu^{*}(\omega) + r(\omega) - u(a, \omega)$ .

We first argue that given  $p^*$ ,  $z^*(\omega)$  is optimal for type- $\omega$  consumers. When  $\mu^*(\omega)=0$ , we have  $p^*(a,\omega)=r(\omega)-u(a,\omega)$ . Thus, type- $\omega$  consumers are indifferent between keeping the data and selling it in market  $(a,\omega)$  for all  $a\in A$ . Therefore,  $z^*(\cdot,\omega)$  is optimal. When  $\mu^*(\omega)>0$ , by complementary slackness, we have that  $\sum_a z^*(a,\omega)=1$ . Therefore, no type- $\omega$  consumer keeps the data. For any  $z(a,\omega)$  such that  $\sum_a z(a,\omega)=1$ , by construction we have  $\sum_a z(a,\omega)(u(a,\omega)+p^*(a,\omega))=\mu^*(\omega)+r(\omega)$ . Therefore, the consumers are indifferent among these strategies, and thus  $z^*(\omega)$  is optimal.

Next, we argue that  $\chi^{\dagger}$  solves the platform's problem given  $p^*$ . We first show that the platform's payoff is non-positive under  $p^*$ . To show this, we only need to show the dual problem of the platform's problem, defined by (A.3), is feasible under  $p^*$ . The dual feasible set is given by:

$$\sum_{\hat{a}} (\pi(\hat{a}, \omega) - \pi(a, \omega)) \lambda(\hat{a}|a) \ge v(a, \omega) - p^*(a, \omega)$$
$$= v(a, \omega) + u(a, \omega) - \mu^*(\omega) - r(\omega)$$

for all  $a, \omega$ , with  $\lambda \geq 0$ . But we know this is feasible because  $\lambda^*$  satisfies these constraints. Given dual feasibility, weak duality implies:

$$\sum_{a,\omega} (v(a,\omega) - p^*(a,\omega))\chi(a,\omega) \le 0$$

for all  $\chi$  that is feasible to the platform.

Finally, by strong duality we have:

$$\sum_{a,\omega} \Big( v(a,\omega) + u(a,\omega) \Big) \chi^{\dagger}(a,\omega) - \sum_{a,\omega} \chi^{\dagger}(a,\omega) r(\omega) = \sum_{\omega} \mu^{*}(\omega) \bar{q}(\omega).$$

This implies:

$$\sum_{a,\omega} (v(a,\omega) - p^*(a,\omega) + \mu^*(\omega)) \chi^{\dagger}(a,\omega) = \sum_{\omega} \mu^*(\omega) \bar{q}(\omega).$$

By complementary slackness we know  $\sum_{a,\omega} \mu^*(\omega) \chi^{\dagger}(a,\omega) = \sum_{\omega} \mu^*(\omega) \bar{q}(\omega)$ , which implies:

$$\sum_{a,\omega} (v(a,\omega) - p^*(a,\omega)) \chi^{\dagger}(a,\omega) = 0.$$

Therefore, we conclude  $\chi^{\dagger}$  solves the platform's problem given  $p^*$ .

# **Online Appendix (For Online Publication Only)**

# **B** Equilibrium Existence

In this section, we prove the existence of an equilibrium of the competitive economy, allowing for arbitrary specification of v. We start by showing that the solution correspondence of  $\mathcal{P}_q$  has nice properties.

### Lemma B.1.

- 1. The solution correspondence  $x^*(q)$  of  $\mathcal{P}_q$  is nonempty-valued, compact-valued, and upper-hemicontinuous.
- 2. V(q) is continuous in q.

*Proof.* Fix q. Note that  $\mathcal{P}_q$  can be reformulated as:

$$\begin{split} \max_{\chi \geq 0} \quad & \sum_{a,\omega} v(a,\omega) \chi(a,\omega) \\ \text{such that} \quad & \sum_{\omega} \left( \pi(a,\omega) - \pi(a',\omega) \right) \chi(a,\omega) \geq 0 \qquad \forall \ a,a' \in A. \\ \text{and} \quad & \sum_{a} \chi(a,\omega) = q(\omega) \qquad \forall \ \omega \in \Omega \end{split} \tag{B.1}$$

In this problem, the objective is continuous in  $\chi$  and the feasible set is nonempty (because  $\chi(\omega,\omega)=q(\omega)$  is always feasible) and compact. Therefore, the solution correspondence is nonempty- and compact-valued. The continuity of  $\chi^*(q)$  follows from Theorem 2 in Böhm (1975). Since  $\chi^*$  is continuous in q and the objective function  $\sum_{a,\omega} v(a,\omega)\chi(a,\omega)$  is continuous in  $\chi$ , the value function V(q) is continuous.

Note that x is a solution to  $\mathcal{P}_q$  if and only if  $\chi(a,\omega):=x(a|\omega)q(\omega)$  is a solution to (B.1), we claim that  $x^*(q)$  is upper-hemicontinuous. It is clear that  $x^*$  is closed-valued, so we only need to show it has a closed graph. Take any  $(q_n,x_n)\to (q,x)$  such that  $x_n\in x^*(q_n)$ , we want to show  $x\in x^*(q)$ . Note that  $\chi_n(a,\omega)\to \chi(a,\omega):=x(a|\omega)q(\omega)$ . By continuity of  $\chi^*$  we know  $\chi\in\chi^*(q)$  and thus  $x\in\chi^*(q)$ .

In light of Lemma B.1, we are ready to prove the existence of an equilibrium.

### **Proposition B.1.** An equilibrium of the competitive economy exists.

*Proof.* We start by introducing a correspondence whose fixed points characterize the set of competitive equilibria. Let  $P = [-M, M]^{|\Omega|}$  be the space of possible equilibrium prices, where M is chosen to be large so that any possible equilibrium prices are within that range. Let  $Q \times X$  be the space of feasible data allocations. Taken together,  $P \times Q \times X$  is a nonempty, compact, and convex set. Define a correspondence  $F: P \times Q \times X \Rightarrow P \times Q \times X$  such that  $(p', q', x') \in F(p, q, x)$  if:

- 1. x' solves problem  $\mathcal{P}_q$ .
- 2. q' solves the consumers' problem given (p, x).<sup>21</sup>
- 3. p' is such that q solves the platform's problem (1) in the first period.

Note that (p, q, x) is a competitive equilibrium if and only if it is a fixed point of F. Therefore, a competitive equilibrium exists if F admits a fixed point. Toward this, we first prove the following claim and then apply Kakutani's fixed point theorem.

**Claim.** *F is nonempty-valued, convex-valued, and has a closed graph.* 

*Proof of the Claim.* We first show that F is nonempty-valued. Fix any (p,q,x). By Lemma B.1,  $\mathcal{P}_q$  admits a solution x'; given (p,x), the consumers' problem always has a solution q'; given q, since  $\mathcal{P}_q$  admits an optimal solution, by strong duality  $\mathcal{P}'_q(v)$  also admits an optimal solution. Lemma A.2 then implies that a price p' under which q solves the platform's problem exists. Therefore,  $(p', q', x') \in F(p, q, x)$ .

Next we show F is convex-valued. Note that by definition of F, given (p,q,x), the choice of p', q', and x' are independent with each other. Therefore, it is sufficient to check convexity for each dimension. If x' and x'' both solve  $\mathcal{P}_q$ , clearly any convex combination also solves it; If q' and q'' both solve the consumers' problem, then any convex combination also solves the consumers' problem. To see this, if under (p,x) consumer  $\omega$  has a strict preference, then  $q'(\omega) = q''(\omega)$ . if under (p,x) consumer  $\omega$  is indifferent, then any  $q(\omega)$  is optimal. If under both p' and p'', q solves the platform's first-stage problem, then by Lemma A.2 we know both

<sup>&</sup>lt;sup>21</sup>Formally, we should impose market clearing saying that  $z'=q'/\bar{q}$  solves the consumers' problem. We skip this step to abbreviate notation.

p' and p'' are solutions to  $\mathcal{P}'_q(v)$ . Therefore, any convex combination of them is still a solution to  $\mathcal{P}'_q(v)$ . Again by Lemma A.2, q solves the platform's first-stage problem under that convex combination.

Finally, we argue F has a closed graph. Suppose  $(p_n, q_n, x_n) \to (p, q, x), (p'_n, q'_n, x'_n) \to (p', q', x')$ , and  $(p'_n, q'_n, x'_n) \in F(p_n, q_n, x_n)$ . We want to show  $(p', q', x') \in F(p, q, x)$ . By Lemma B.1, we know the solution correspondence of  $\mathcal{P}_q$  is upper-hemicontinuous, so x' is a solution to  $\mathcal{P}_q$ ; To see q' solves the consumers' problem, note that for all  $\omega$  and  $z \in [0, \bar{q}(\omega)]$ :

$$q'_n(\omega)(p_n(\omega) + \sum_a u(a,\omega)x_n(a|\omega)) + (\bar{q}(\omega) - q'_n(\omega))r(\omega)$$

$$\geq z(p_n(\omega) + \sum_a u(a,\omega)x_n(a|\omega)) + (\bar{q}(\omega) - z)r(\omega)$$

By continuity we get:

$$q'(\omega)(p(\omega) + \sum_{a} u(a,\omega)x(a|\omega)) + (\bar{q}(\omega) - q'(\omega))r(\omega)$$

$$\geq z(p(\omega) + \sum_{a} u(a,\omega)x(a|\omega)) + (\bar{q}(\omega) - z)r(\omega)$$

Therefore, q' is optimal for the consumers given (p, x); To see under p', q solves the platform's problem, note that for all  $\tilde{q} \geq 0$ :

$$V(q_n) - \sum_{\omega} p'_n(\omega) q_n(\omega) \ge V(\tilde{q}) - \sum_{\omega} p'_n(\omega) \tilde{q}(\omega)$$

Since V is continuous by Lemma B.1, taking limit we get:

$$V(q) - \sum_{\omega} p'(\omega)q(\omega) \ge V(\tilde{q}) - \sum_{\omega} p'(\omega)\tilde{q}(\omega).$$

This completes the proof that  $(p', q', x') \in F(p, q, x)$ .

Using the Claim, we can apply Kakutani's fixed-point theorem to F and conclude that F admits a fixed point. Therefore, a competitive equilibrium exists.

# C Complete Equilibrium Characterization for Section 3.1

In this section, we characterize the entire set of equilibria for our example from Section 3.1. We first note that in order for the platform's problem to admit a solution, we must have  $p^*(1) \ge$ 

 $0, p^*(2) \ge 0, p^*(1) + p^*(2) \ge \gamma_u$ . Moreover, in order for the platform to trade, we must have  $p^*(1) + p^*(2) = \gamma_u$ .

Case 1:  $2\bar{r} - 1 < \gamma_u < \bar{r}$ . The unique equilibrium allocation is no trade, i.e.,  $q^*(\omega) = 0$  for all  $\omega$ , and it is supported by a price vector  $p^*$  satisfying:

$$p^*(1) \in [0, \bar{r}]$$
 and  $p^*(2) \in [\max\{0, \gamma_u - p^*(1)\}, \bar{r}].$  (C.1)

Next we explain why this is the solution. Note that type-1 consumers are willing to sell only if  $p^*(1) \geq \bar{r} > \gamma_u$ . However, in this case we cannot have  $p^*(1) + p^*(2) = \gamma_u$ . Therefore, the unique equilibrium allocation is  $q^*(\omega) = 0$ . It can be supported by  $x^*(\omega|\omega) = 1$ . It can be checked that with the prices in (C.1), it is optimal for the consumers not to sell and for the platform not to buy. Any other price will induce some type of consumers to strictly prefer selling.

Case 2:  $\gamma_u > 2\bar{r}$ . There is a unique equilibrium such that:  $p^*(1) = \gamma_u$  and  $p^*(2) = 0$ ;  $z^*(1) = 1$  and  $z^*(2) = \min\{1, \frac{\bar{q}(1)}{\bar{r}\bar{q}(2)}\}$ ;  $q^*(1) = \bar{q}(1)$  and  $q^*(2) = \min\{\bar{q}(2), \frac{\bar{q}(1)}{\bar{r}}\}$ ;  $x^*(1|1) = 1$  and  $x^*(1|2) = \frac{q^*(1)}{q^*(2)}$ . It can be easily checked that this is an equilibrium. To show uniqueness, note that since  $\gamma_u > 2\bar{r}$ , at least one type has a strict incentive to sell because  $p^*(1) + p^*(2) \ge \gamma_u$ . If type-1 has a strict incentive, we have  $q^*(2) \ge \min\{\bar{q}(2), \frac{\bar{q}(1)}{\bar{r}}\}$  since  $p^*(2) \ge 0$ , but this requires  $p^*(1) = \gamma_u$ ,  $p^*(2) = 0$ ; if type-2 has a strict incentive, in order for the platform to be willing to buy, we must have  $p^*(1) = \gamma_u$ ,  $p^*(2) = 0$ .

Case 3:  $\bar{r} < \gamma_u \le 2\bar{r}$ . It can be easily checked that both equilibria of Case 1 and Case 2 continue to be an equilibrium in this case. Next we argue those are all possible equilibria. On one hand, for the equilibria with no trade, the price has to satisfy (C.1), otherwise some type will have a strict incentive to sell. On the other hand, given any equilibrium with trade, we must have  $p^*(1) + p^*(2) = \gamma_u$ . Moreover, we must have  $q^*(1) > 0$ , otherwise the platform is not willing to buy  $\omega = 2$  at a positive price while type-2 consumers are not willing to sell at 0 price. If  $q^*(1) > 0$ , we must have  $q^*(2) \ge \min\{\bar{q}(2), \frac{q^*(1)}{\bar{r}}\}$  because  $p^*(2) \ge 0$  and the platform will choose  $x^*(1|1) = 1, x^*(1|2) = \frac{q^*(1)}{q^*(2)}$ . Since  $q^*(2) > q^*(1)$ , in order for the platform to be willing to buy, it must be the case that  $p^*(1) = \gamma_u, p^*(2) = 0$ . The unique equilibrium with trade then follows.

Case 4:  $\gamma_u = \bar{r}$ . In this case, the equilibria with no trade is the same as Case 1. The equilibria

with trade satisfy  $p^*(1) = \gamma_u = \bar{r}$ ,  $p^*(2) = 0$  with

$$0 < q^*(1) \le \bar{q}(1), \ q^*(2) = \min\{\bar{q}(2), \frac{q^*(1)}{\bar{r}}\}.$$

It is easy to check these are equilibria. To see these capture all equilibria with trade, we can follow the same argument as Case 3 to derive the unique equilibrium price under trade:  $p^*(1) = \gamma_u = \bar{r}$ ,  $p^*(2) = 0$ . With these prices, since type-1 consumers are indifferent, any  $0 \le q^*(1) \le \bar{q}(1)$  is optimal for them.  $q^*(2)$  is then pinned down by the indifference condition of type-2 consumers.

To sum up, the constrained-efficient allocation  $q^{\circ}(1) = q^{\circ}(2) = \bar{q}(1)$  can never be an equilibrium, so all equilibria of this competitive economy are inefficient.

## **D** Social Welfare

In the main text, we focused on a notion of welfare that excludes the merchant's profit (see Equation (2)). In this section, we show that an analogous result to Proposition 2 holds if we allow the welfare function to include the merchant's profit. We refer to this as the "social welfare," defined as

$$SW(q,x) = \sum_{a,\omega} \left( v(a,\omega) + u(a,\omega) + \pi(a,\omega) \right) x(a|\omega) q(\omega) + \sum_{\omega} \left( \bar{q}(\omega) - q(\omega) \right) r(\omega). \tag{D.1}$$

In light of this, the new efficiency benchmark is as follows.

**Definition 5.** An allocation  $(q^{\circ}, x^{\circ})$  is **constrained socially efficient** if it solves

$$\max_{q,x} \quad SW(q,x)$$
such that  $q \leq \bar{q}$ ,
and  $x \text{ solves } \mathcal{P}_q$ .

**Remark D.1.** Define "unconstrained social efficiency" by replacing the second constraint in Definition 5 with obedience constraints, in the spirit of Definition 4. Note that the notion of "constrained social efficiency" and "unconstrained social efficiency" are equivalent. This is because when trading-off  $\pi$  and u, the planner uses weights  $1 + \gamma_{\pi}$  and  $1 + \gamma_{u}$  and the platform uses weights  $\gamma_{\pi}$  and  $\gamma_{u}$ . Therefore, as shown in Figure 1, their optimal choice of x coincides given any database y. In this sense, their incentives are aligned in choosing y.

Indeed, whenever the platform's incentive in choosing x is aligned with the planner's preference, constrained efficiency is equivalent to unconstrained efficiency. For example, when the planner does not take into account the merchant's payoff as in the main text, the platform's incentive is aligned with the planner when either  $\gamma_{\pi} \geq 1 + \gamma_{u}$  or  $\gamma_{u} \geq \gamma_{\pi}$ . In these cases, the positive statements of Proposition 2, 3, and 4 can be equivalently stated in terms of unconstrained efficiency (cf. Definition 4).

Under this new criterion of social efficiency, we have the following result, which extends Proposition 2 to this more demanding efficiency benchmark.

**Proposition D.1.** Let  $(p^*, z^*, q^*, x^*)$  be an equilibrium of the competitive economy. If  $\gamma_{\pi} > \gamma_u$  and, in addition,  $r(\omega) \notin [\gamma_{\pi}\omega, (1+\gamma_{\pi})\omega)$  for all  $\omega$ , the equilibrium allocation  $(q^*, x^*)$  is constrained socially efficient. Otherwise, the equilibrium allocation can be socially inefficient.

*Proof.* We prove sufficiency here. Let  $(p^*, z^*, q^*, x^*)$  be a competitive equilibrium. By Proposition 2, we know the equilibrium allocation  $(q^*, x^*)$  is constrained efficient. Moreover, following the argument in the proof of Proposition 2, we also know that  $x^*(a|\omega) = \hat{x}(a|\omega)$  whenever  $q^*(\omega) > 0$ , where  $\hat{x}(\omega|\omega) = 1$  is the full-disclosure mechanism. Therefore,

$$q^* \in \arg \max_{q \leq \bar{q}} \sum_{a,\omega} \Big( v(a,\omega) + u(a,\omega) \Big) \hat{x}(a|\omega) q(\omega) - \sum_{\omega} r(\omega) q(\omega)$$
$$= \arg \max_{q \leq \bar{q}} \sum_{\omega} \Big( \gamma_{\pi} \omega - r(\omega) \Big) q(\omega).$$

The solution to this problem is  $q^*(\omega) = \bar{q}(\omega)$  if  $\gamma_\pi \omega > r(\omega)$ ,  $q^*(\omega) = 0$  if  $\gamma_\pi \omega < r(\omega)$ , and  $q^*(\omega) \in [0,1]$  if  $\gamma_\pi \omega = r(\omega)$ . The constrained socially efficient allocation  $(q^\circ, x^\circ)$  also features  $x^\circ(a|\omega) = \hat{x}(a|\omega)$  whenever  $q^\circ(\omega) > 0$ . Therefore, the solution of the social welfare problem is  $q^\circ(\omega) = \bar{q}(\omega)$  if  $(1 + \gamma_\pi)\omega > r(\omega)$ ,  $q^\circ(\omega) = 0$  if  $(1 + \gamma_\pi)\omega < r(\omega)$ , and  $q^\circ(\omega) \in [0,1]$  if  $(1 + \gamma_\pi)\omega = r(\omega)$ . When  $r(\omega) \notin [\gamma_\pi \omega, (1 + \gamma_\pi)\omega)$  for all  $\omega$ , the equilibrium allocation  $(q^*, x^*)$  is also a solution to social welfare problem, and thus constrained socially efficient.

Intuitively, if we take into account the merchant's profit, the inefficiency can arise from two sources. The first one is the pooling externality discussed in the main text. When  $\gamma_{\pi} > \gamma_{u}$ , the only optimal mechanism for the platform given any q is full disclosure, so the pooling externality disappears. The second one is a traditional externality. Since the platform does not

take into account the merchant's payoff, it refuses to buy data when the price is high, even when trade is still socially optimal. When the sufficient condition of the proposition is not satisfied, the equilibrium can be inefficient.

Next, we illustrate the two sources of inefficiency using the example of Section 3.1. We will denote the constrained-efficient allocation by  $(q^{\circ}, x^{\circ})$ . (This is the same constrained-efficient allocation as in Section 3.1, defined by Definition 2.) We also denote the equilibrium allocation in Case 1 (inefficiently low trade) by  $(q_L^*, x_L^*)$  and in Case 2 (inefficiently high trade) by  $(q_H^*, x_H^*)$ . These are characterized in Section 3.1.

We first argue that in both cases, the social welfare of the equilibrium,  $SW(q^*, x^*)$ , is strictly lower than  $SW(q^\circ, x^\circ)$ . As before, this is originated from the pooling externality. Using the characterizations in Section 3.1, we can directly compute:

$$\begin{split} SW(q^{\circ}, x^{\circ}) &= \bar{q}(1)(3 + \gamma_{u}) + \bar{r}(\bar{q}(2) - \bar{q}(1)), \\ SW(q_{L}^{*}, x_{L}^{*}) &= \bar{r} < SW(q^{\circ}, x^{\circ}), \\ SW(q_{H}^{*}, x_{H}^{*}) &= \bar{q}(3 + \gamma_{u}) + \bar{r} \max\{0, \bar{q}(2) - \frac{\bar{q}(1)}{\bar{r}}\} < SW(q^{\circ}, x^{\circ}). \end{split}$$

The take is that, even if we measure efficiency using social welfare (Equation (D.1)), the equilibria are still suboptimal compared to the constrained-efficient allocation. One may suspect that in Section 3.1, the inefficiency is an artifact that we did not take into account the merchant's profit, but as we highlight here, that is not the case.

In addition to the pooling externality, there is a new source of inefficiency: since in this case we have  $(1+\gamma_\pi)\omega > \bar{r} > \gamma_\pi\omega = 0$ , even  $(q^\circ, x^\circ)$  is not constrained socially efficient. The social welfare is maximized at  $q^\bullet(\omega) = \bar{q}(\omega)$  and  $x^\bullet(1|1) = 1$ ,  $x^\bullet(1|2) = \frac{\bar{q}(1)}{\bar{q}(2)}$ , which gives a social welfare of

$$SW(q^{\bullet}, x^{\bullet}) = \bar{q}(1)(\gamma_u + 1) + 2(\bar{q}(2) - \bar{q}(1)) > SW(q^{\circ}, x^{\circ}).$$

Therefore, the constrained-efficient allocation is not constrained socially efficient. This additional gap is created by the fact that the profit of the merchant is not taken into account by the platform or the consumers. This is a traditional externality that can arise even without the informational friction discussed in our paper. For instance, consider the case where there is only one type of consumers  $\omega = 1$  with 0 < r(1) < 1. The platform's objective has  $\gamma_u > \gamma_{\pi} = 0$ . Then the constrained socially efficient allocation is  $q^{\bullet}(1) = 1$ , but the only equilibrium is no trade.