

Beating the Algorithm:
Consumer Manipulation, Personalized Pricing, and Big-Data Management

Xi Li

Faculty of Business and Economics, The University of Hong Kong, Hong Kong. Email: xili@hku.hk

Krista J. Li

Kelly School of Business, Indiana University, Bloomington, IN. Email: kjli@indiana.edu

Abstract

1. **Problem definition:** Firms heavily invest in big-data technologies to collect consumer data and infer consumer preferences for price discrimination. However, consumers can use technological devices to manipulate their data and fool firms to obtain better deals. We examine how a firm invests in collecting consumer data and makes pricing decisions and whether it should disclose its scope of data collection to consumers who can manipulate their data.
2. **Methodology/results:** We develop a game-theoretic model to consider a market in which a firm caters to consumers with heterogeneous preferences for a product. The firm collects consumer data to identify their types and issue an individualized price, whereas consumers can incur a cost to manipulate data and mimic the other type. We find that when the firm does not disclose its scope of data collection to consumers, it collects more consumer data. When the firm discloses its scope of data collection, it reduces data collection even when collecting more data is costless. The optimal scope of data collection increases when it is more costly for consumers to manipulate data but decreases when consumer demand becomes more heterogeneous. Moreover, a lower cost for consumers to manipulate data can be detrimental to both the firm and consumers. Lastly, disclosure of data collection scope increases firm profit, consumer surplus, and social welfare.
3. **Managerial implications:** Our findings suggest that a firm should adjust its scope of data collection and prices based on whether the firm discloses the data collection scope, consumers' manipulation cost, and demand heterogeneity. Public policies should require firms to disclose their data collection scope to increase consumer surplus and social welfare. Even without such a mandatory disclosure policy, firms should voluntarily disclose their data collection scope to increase profit. Moreover, public educational programs that train consumers to manipulate their data or raise their awareness of manipulation tools can ultimately hurt consumers and firms.

Keywords: big-data technology, data manipulation, disclosure, pricing, game theory.

1 Introduction

Firms heavily invest in big-data technologies to collect consumers' data and infer their preferences for price discrimination. A survey of 65 executives at Fortune 1000 companies shows that, in 2019, 55% of those companies invested more than \$50 million in big-data projects, 21% invested more than half a billion, and 92% increased their big-data investments year over year (Osborne 2019). In 2019, Target launched its data-driven loyalty program, Target Circle, to track consumers' data and deliver personalized coupons (Perez 2019). Including Target, more than 90% of companies invest in some form of loyalty program (Experian 2014; Morgan 2020), which enables them to collect consumers' data for price discrimination (Smith 2009). Office Depot uses consumers' browsing history and location data to vary prices (Mahdawi 2016). One hundred and seventy brands, including Walmart, use consumers' real-time location data to distribute mobile coupons to them (Chen et al. 2017; Heine 2015).

Many firms disclose their scope of data collection to consumers. Target (2022) for example, announces that it collects consumers' purchase and return histories, geo- and in-store location information. Office Depot (2020) discloses that it collects various data components directly from consumers and uses browser cookies and web beacons. As a wide range of firms collect consumer data for price discrimination and are transparent about it, consumers have become more aware of firms' data collection and data-driven price discrimination.

While firms price discriminates against consumers using their data, sophisticated consumers can also strategically manipulate their data to obtain better deals. Knowing that firms often email coupons to consumers who abandon a purchase, consumers defer purchases by adding products to their shopping cart without checking out to mimic the behavior of buyers with a low willingness to pay (Townley et al. 2017). Orbitz collects consumers' device information and offers Mac higher prices than their PC-using counterparts (Mattioli, 2012). On the other hand, consumers use browser plugins such as the User-Agent Switcher to manipulate their device information. Uber riders have tried to beat the firm's system by requesting and then rejecting quotes for rides and scheduling a longer waiting time for pick up to simulate greater price sensitivity to get lower prices. A female rider goes viral on TikTok by showing how these manipulative tricks get her cheaper rides, which has more than 600,000 views (Fenton 2021).

An interesting tension exists between firms' investment in data collection to identify more consumers and consumers' data manipulation behavior. When consumers manipulate their data, the data that firms collect and use for targeted pricing become contaminated. In this case, a firm may not always benefit from acquiring more data to infer the preferences of more consumers. In this paper, we examine how a firm's optimal scope of data collection, personalized pricing, profit, consumers surplus, and social welfare change with consumers' data manipulation.

A critical factor that needs to be considered when answering these questions is the transparency of data collection. Some firms collect and use consumers' data without notifying consumers, while others are transparent about the scope of their data collection. Knowing that a company collects more data, consumers anticipate a higher chance that the company identifies their preference and can be more willing to manipulate their data. In turn, this transparency may induce the firm to change its pricing and scope of data collection, which ultimately affects its profit as well as consumer surplus and social welfare. It is unclear how a firm's disclosure of data collection affects its profit. As such, managers need to assess how transparency affects their data collection and profit and then decide whether or not to disclose their data collection to consumers.

As more firms use consumer data for price discrimination, there have also been growing public concerns about consumer privacy (Goldfarb and Tucker 2010). Consumer advocates and public policymakers strive to regulate the market while protecting consumers and improving social welfare. Policymakers can create public policy programs that educate consumers on ways to protect themselves through effective data manipulation. Such educational programs can reduce consumers' technological barriers or cost to manipulate their data and enable them to get better deals from firms. In addition, a stricter privacy law could also reduce consumers' cost of data manipulation by limiting firms' ability to collect consumer data from multiple sources (Valletti and Wu 2020). For example, in 2018, the European Union launched the General Data Protection Regulation (GDPR) to regulate companies' collection and handling of consumer data (GDPR 2020). In the United States, all 50 states have passed legislation on data privacy, although the degree of data protection varies (Blog 2018). As a result, consumers' cost to manipulate data can vary by state. Formal analyses need to investigate how these different approaches affect consumers and society as a whole to recommend a viable approach to policymakers.

In this paper, we develop a game-theoretic model to explore managerial and public policy

issues surrounding data collection and consumers' data manipulation. We consider a market in which a firm caters to consumers with heterogeneous preferences (i.e., high-type and low-type) for a product. The firm determines the scope of its data collection: if the firm expands the scope of its data collection, there is a higher likelihood that a particular consumer is tracked. As a result, the firm can identify more consumers' type and then charge a different price to each type of consumer. Meanwhile, consumers can incur a cost to manipulate their data. Specifically, high-type consumers have incentives to manipulate their data to mimic low-type consumers to enjoy the low price that the firm charges to low-type consumers. Using this setting, we examine how consumers' data manipulation affects the firm's pricing and data collection decisions. We also evaluate how changes in consumers' cost of manipulation (enabled by technological devices or educational programs) affect the firm's profit, consumer surplus, and social welfare. Also, we compare equilibrium outcomes when the firm discloses and does not disclose the scope of its data collection; this offers insight into whether or not the firm should disclose its data collection to consumers or if public policies should mandate the firm to do so.

Our analysis reveals several noteworthy findings. First, we find that when consumers have the dual ability to manipulate data and observe the firm's data collection, the firm should cut back on its data collection scope and forgo identifying types and preferences of all consumers, even if doing so is costless. The reason is that collecting more consumers' data benefits the firm by allowing it to infer more consumers' preferences and offer them customized prices. However, collecting more data hurts the firm by inducing more high-type consumers to manipulate their data, which contaminates the accuracy of the data, and weakens the effectiveness of price discrimination. The latter cost of data collection can dominate the former benefit, which makes the firm better off not collecting more data even if it is costless. The optimal scope of data collection increases when it is more costly for consumers to manipulate data but decreases when consumer demand becomes more heterogeneous. When consumers cannot manipulate data or cannot observe the firm's data collection, the firm should increase its data collection scope, collect as much data as possible to infer the preferences of more consumers, and charge them customized prices.

Second, we find that consumer manipulation reduces the firm's profit and social welfare. The intuition is related to the rent-seeking behavior in economics; that is, high-type consumers manipulate data to pay the price targeted to low-type consumers and obtain a rent, which reduces the

firm's ability to perform price discrimination and social efficiency. Interestingly, consumer surplus also decreases with consumer manipulation. This is because as more consumers incur a cost to manipulate data, the firm raises prices for low-type consumers. The pricing and cost effects reduce consumer surplus. Consequently, the emergence of new devices, training programs, and technology that help consumers manipulate data cost-efficiently can end up being detrimental to consumers, the firm, and society.

Third, we find that the firm's disclosure of data collection to consumers improves consumer surplus, the firm's profit, and social welfare. When a firm discloses its data collection scope so that consumers observe it, the firm reduces the amount of consumer data that it collects, which then discourages consumers from manipulating their data. When a firm does not disclose its data collection scope, however, it will collect as much data as possible, which encourages consumers' data manipulation. Thus, disclosure of data collection increases the accuracy of the firm's targeted pricing and leads to reduced prices for low-type consumers. These effects increase consumer surplus and firm profit, which then improves social welfare. Therefore, public policies should mandate firms to disclose their data collection to consumers. Even without mandatory-disclosure policies, the firm should voluntarily disclose its data collection to improve profit.

We generalize our base model to consider consumers with continuous types of preferences. Interestingly, we find that, under this generalization, a segment of high-type consumers mimics a segment of low-type consumers according to a distribution, while low- and medium-type consumers do not manipulate their data. Moreover, the firm adopts a hybrid pricing strategy by 1) charging the same segment-specific price to low-type consumers and high-type consumers who manipulate their data and 2) charging an individual-specific price to medium-type consumers. Lastly, we analyze a situation in which consumers have a heterogeneous ability to manipulate their data and find that enabling more consumers to manipulate their data can be detrimental to consumers. We also discuss how our main results could continue to hold in a competitive setting and when consumers hide their identity to avoid being recognized by the firm.

1.1 Related Literature

Our research builds upon the growing literature investigating how firms use consumer data to identify and target consumers. Thisse and Vives (1988) first consider a case in which firms can offer consumers personalized prices based on their locations and show that personalized pricing intensifies price competition between firms. Nonetheless, these firms cannot help adopting personalized pricing, thereby finding themselves in the form of the prisoner's dilemma. Chen et al. (2001) consider a scenario in which competing firms implement imperfect targeting technologies to classify their consumers. They find that, when targeting is imperfect, a firm improving its targeting technologies benefits not only itself but also its competitor, thereby leading to a "win-win" situation. Koh et al. (2017) investigate voluntary profiling, a concept that has been adopted by many firms, under which consumers can voluntarily opt-in firms' data collection process. They show that, contrary to common wisdom, voluntary profiling often works to the detriment of social welfare and consumer surplus. Chen et al. (2017) study the practice of firms targeting consumers based on their real-time location, which, reversely, gives consumers the incentive to travel across different locations for better offers. They find that a firm's profit can be higher under mobile geo-targeting than under uniform or traditional targeted pricing. Taylor (2004) considers a dynamic market in which one firm sells its consumers' purchase history information to another firm so that it may identify individual consumers and charge them personalized prices. He shows that welfare analyses depend on whether consumers are naive (i.e., do not anticipate sales of their data) or sophisticated (i.e., anticipate sales of their data). Argenziano and Bonatti (2021) study a related scenario in which a consumer interacts with two firms sequentially but firms set not only price but also quality levels. They examine how data linkage between the two firms affects consumer surplus under various privacy regulations. They suggest that such data linkage benefits a naive consumer or a consumer who sufficiently values quality. This stream of research reveals that in a two-period dynamic setting, firm's data collection generates a ratchet effect (Freixas et al. 1985; Laffont and Tirole 1993); that is, first-period high-valuation consumers anticipate the firm to use their revealed preference for quality and price discrimination in the second period and become reluctant to reveal preferences through making first-period purchase, resulting in pooling behaviors. To induce these consumers to make purchase, firms reduce first-period prices, which harms firm

profits. Unlike these studies, we consider a static model in which a firm decides how much data it collects to infer consumer preferences while high-type consumers decide whether to manipulate data to mimic low-type consumers. We show that a higher cost for consumers to manipulate data can improve consumer surplus. Moreover, a firm does not want to collect all consumers' data even if it is costless, because doing so would lead more consumers to manipulate data, reducing data quality and the firm's profit from practicing price discrimination. The mechanism is different from the ratchet effect that shows that consumers' data manipulation harms firm profits in the first period before price discrimination takes place.

Our research also informs analytical studies on consumer privacy. There is a growing literature examining firms' price discrimination strategy when consumers can remain anonymous from the firm, possibly at a privacy cost. Conitzer et al. (2012) explore a setting in which a firm recognizes repeat customers and uses their past purchases to price discriminate, but consumers can, at a cost, hide their identities and avoid being recognized by the firm. They find that, when consumers can freely maintain their anonymity, they individually choose to do so, which results in the highest profit for the firm. Belleflamme and Vergote (2016) investigate a scenario in which a firm uses a technology to identify consumers' willingness to pay, while consumers can adopt a hiding technology to avoid being recognized. They show that the hiding technology can leave consumers worse off. Montes et al. (2019) consider a duopoly in which one or more firms price discriminate against consumers using their information, yet consumers can pay a privacy cost to prevent the firms from doing so. They find that the firms do not always benefit from a higher privacy cost. Valletti and Wu (2020) examine a firm's investment in the accuracy of its consumer profiling when consumers take costly actions to conceal or hide their identity from being recognized by firms. We move one step forward by allowing consumers to not only *hide* their data, but also *manipulate* their data to fool the firm. Consumer hiding and consumer manipulation have different implications for the firm's big data management. When consumers hide their data from being recognized by the firm, this hiding behavior does not change the accuracy or quality of the firm's data; that is, whenever the firm identifies a consumer's type from data, the identification is always correct. By contrast, when a high-type consumer manipulates her data to mimic a low-type, the manipulation contaminates the firm's data by reducing its accuracy or informativeness. The firm may misidentify a high-type consumer to be low-type when the consumer manipulates data. This negative

impact induces the firm to cut back on its data collection. Our result shows that when consumers observe the firm's data collection and manipulate data, the firm should not collect data to identify all consumers' type even if data collection is costless. This result no longer holds when consumers can only hide their data.

Lastly, given that we examine the firm's disclosure of data collection scope, our research is related to research on disclosure of private information. Existing research has examined firms' costless and costly disclosure of quality information (Grossman 1981; Jovanovic 1982; Guo and Zhao 2009), horizontal attribute (Gu and Xie 2013), comparative information (Anderson and Renault 2009). In this paper, we consider the firm's disclosure of its data collection scope, which affects consumers' incentives of data manipulation.

2 The Model

Consider a firm that sells a product to a unit mass of consumers. The marginal cost to produce the product is constant and we standardize it to zero.

Consumer Type. Consumers have heterogeneous preferences for the firm's product. We model the heterogeneity by considering two types of consumers. A fraction α of consumers are high-type consumers with a linear demand $H - p$, where H is the demand capacity that reflects the preference of the consumer and p is the price of the product. The rest $1 - \alpha$ of them are low-type consumers whose demand function is $L - p$, where $H > L > 0$ and $L \geq \frac{H}{2}$. This distribution of consumer preferences, which might be assessed through market research, is assumed to be common knowledge between the firm and the consumers (Wathieu and Bertini, 2007). The assumption $L \geq \frac{H}{2}$ guarantees non-negative demand for all consumers. $t_i \in \{H, L\}$ denotes consumer i 's type.

Consumer Manipulation. Consumers can incur a fixed cost $c \geq 0$ to manipulate their data to mimic the other type of consumers. This cost reflects the efforts and time a consumer undertakes to successfully mimic the other type. If a firm has consumer i 's data that classifies the consumer as type $\tilde{t}_i \in \{H, L\}$, the firm thinks that the type of this consumer is \tilde{t}_i , which equals the real type t_i if the consumer does not manipulate data, or the opposite type if the consumer does. Given that consumers manipulate data individually and firms cannot track such manipulative behaviors, the

firm does not know whether or not a specific consumer manipulates data or not.¹

Data Collection Scope. If the firm increases its scope of data collection, it collects data to identify the types of more consumers; there is a higher chance that a consumer will be tracked and identified by the firm and the consumer coverage of the firm's big data increases. Specifically, for each consumer i , the firm analyzes its collected data to generate a signal $s_i \in \{H, L, \emptyset\}$ about the consumer's type, with

$$s_i = \begin{cases} \tilde{t}_i, & \text{with probability } \rho, \\ \emptyset, & \text{with probability } 1 - \rho. \end{cases} \quad (1)$$

If the firm has data to identify consumer i 's type, the database sends a signal of $s_i \in \{H, L\}$ that indicates the type of the consumer. If the firm does not have the consumer's data to identify his or her type, the database sends a signal of $s_i = \emptyset$ that does not reveal any additional information about consumer i 's type beyond the prior distribution. $\rho \in [0, 1]$ represents the probability that the firm has data about consumer i based on which to infer his or her type. Thus, ρ reflects the scope of the firm's data collection and the consumer coverage of the firm's big data. This is a long-term big-data investment decision that takes place before firms actually go out to collect consumer data. For example, a firm needs to acquire database infrastructures and develop profiling algorithms to collect, manage, and analyze data to issue personalized prices. A higher value of ρ indicates that the firm's database covers data of more consumers, allowing the firm to identify the types of more consumers. When $\rho = 0$, the big data has no consumer coverage, so it does not provide the firm with any useful information about any consumer's type. Reversely, when $\rho = 1$, the big data has complete consumer coverage in the sense that it has data to identify the type of all consumers. The firm can increase ρ by acquiring more consumer data to identify the types of more consumers (e.g., purchasing more consumer data from third parties).² Such a "true-or-noise" information

¹Consumers may take action to delete their data to hide their identity and avoid being recognized by the firm. In Section 6.2, we discuss how the results may change with consumer hiding instead of consumer manipulation. In addition, our model applies to an alternative setting when consumers are heterogeneous in their costs to serve. The firm can use the big-data technology to recognize consumers' cost types, and consumers can manipulate their cost type at a manipulation cost. All our results still hold under this alternative interpretation.

²We standardize the firm's cost to improve the consumer coverage of its big data to zero to focus on how demand-side strategic factors affect the firm's data-collection decisions. We intend to show that, even in the absence of cost considerations, the firm may prefer not to pursue full consumer coverage. Alternatively, we could assume a quadratic cost function for the firm to improve the consumer coverage ρ of its big data, which would reduce the firm's equilibrium consumer coverage. This will only strengthen our results.

structure has been used in the literature before (Koh et al. 2017; Lau 2008).

It is worth mentioning that our model admits an alternative interpretation: Even if the firm collects data from all consumers, it may not always be able to identify a consumer's type successfully, and the probability that the firm successfully identifies a consumer, ρ , increases with the amount of data that the firm collects. All our results go through under this alternative interpretation.³

Timing. The timing of the game is as follows. In the first stage, the firm makes its long-term decision on the scope of its data collection by setting the consumer coverage (i.e., ρ) of its big data. Consumers observe ρ when the firm discloses its data collection scope and do not observe it when the firm withholds this information. In the second stage, consumer i privately observes his or her type t_i and decides whether or not to manipulate his or her data at a cost c . In the third stage, the firm uses its data to identify consumer types and receives the signal s_i as described above. In the fourth stage, the firm offers a price p_{s_i} contingent on its signal s_i to each consumer, and the consumer decides how much to purchase from the firm.

3 A No Manipulation Benchmark

We first consider a benchmark model in which consumers cannot manipulate their data, either because they do not have the technical ability to do so or because it is too time- or cost-inefficient to do so.

Proposition 1 *Without data manipulation, the firm's profit increases with the consumer coverage ρ of its big data.*

When consumers cannot manipulate their data, their data reveals their true type. As a result, the firm sets prices based on the consumer type that its big-data technology identifies. For a high-type consumer, the firm charges $\frac{H}{2}$; for a low-type consumer, the firm charges $\frac{L}{2}$. For a consumer that a firm has no data to identify, the firm charges $\frac{\alpha H + (1-\alpha)L}{2}$ based on the expected type $\alpha H + (1-\alpha)L$. The firm's profit is $\Pi = \rho[\alpha H^2 + (1-\alpha)L^2]/4 + (1-\rho)[\alpha H + (1-\alpha)L]^2/4$, increasing with ρ .

The firm's profit increases with the consumer coverage of its big data because identifying a consumer's type allows it to charge the optimally customized price. When the firm cannot identify a consumer's type, it charges a price based on the consumer's expected type; then, the price

³We thank an anonymous referee who suggested this interpretation.

deviates from the optimal one that the firm should charge according to the consumer's true type. This result suggests that, when consumers cannot manipulate their data, the firm always benefits by increasing the consumer coverage of its big data to identify the types of more consumers. Consequently, the firm should invest more in data collection to pursue a full consumer coverage (i.e., $\rho = 1$).

4 Exogenous Consumer Coverage

In this section, we analyze the subgame given the value of ρ , which refers to situations in which the consumer coverage of the firm's big data ρ , is exogenously determined by technology or availability of consumer data.

When consumers can manipulate their data, a low-type consumer has no incentives to manipulate data to mimic a high-type consumer because that would only result in a higher price. Consequently, we focus on high-type consumers' incentives to manipulate data.

If a high-type consumer does not manipulate data, the expected utility will be

$$CS_H = \rho \int_{p_H}^H (v - p_H) dv + (1 - \rho) \int_{p_\emptyset}^H (v - p_\emptyset) dv, \quad (2)$$

where the first term on the right-hand side is the expected utility when the firm identifies the consumer's true type H and charges the high price p_H . The second term is the expected utility when the firm's database does not cover the consumer so that it cannot identify the consumer's type; in this case, the firm charges the consumer the price p_\emptyset .

If a high-type consumer manipulates his or her data to mimic a low-type consumer, the expected utility will be

$$CS_{HL} = \rho \int_{p_L}^H (v - p_L) dv + (1 - \rho) \int_{p_\emptyset}^H (v - p_\emptyset) dv - c, \quad (3)$$

where the first term on the right-hand side is the consumer's expected utility when the firm mistakenly identifies the high-type consumer as a low-type consumer, and the second term is the expected utility when the firm's database does not cover this consumer to identify his or her type.

A high-type consumer is better off manipulating data if and only if $CS_{HL} > CS_H$, which trans-

lates into

$$\frac{\rho}{2}(2H - p_H - p_L)(p_H - p_L) > c, \quad (4)$$

where the left-hand side is the benefit of manipulation that depends on the prices that the firm charges, and the right-hand side is the cost of manipulation.

Given consumers' manipulative behavior, the firm updates its posterior belief about the consumer's type upon receiving a signal from the consumer data and tailors its price offering to the consumer. We delegate the complete analysis to the Appendix and present the equilibrium outcome in Proposition 2. When multiple equilibria exist, we use divinity criterion D1 as the refinement criterion to pin down the unique perfect Bayesian equilibrium:

Proposition 2 *Consumers' data manipulation and the firm's prices and profits vary with the cost of manipulation (c) as follows:*

- a. *If the cost is high (i.e., $c \geq \bar{c} = \rho(H - L)(3H - L)/8$), no consumers manipulate their data. The firm's prices are $p_H = \frac{H}{2}$, $p_L = \frac{L}{2}$, and $p_\emptyset = \frac{\alpha H + (1 - \alpha)L}{2}$ and profit is $\Pi = \rho[\alpha H^2 + (1 - \alpha)L^2]/4 + (1 - \rho)[\alpha H + (1 - \alpha)L]^2/4$.*
- b. *If the cost is low (i.e., $c \leq \underline{c} = \rho(1 - \alpha)(H - L)((3 - \alpha)H - (1 - \alpha)L)/8$), all high-type consumers manipulate their data. Prices are $p_H = \frac{H}{2}$ and $p_L = p_\emptyset = \frac{\alpha H + (1 - \alpha)L}{2}$, and profit is $\Pi = (\alpha H + (1 - \alpha)L)^2/4$.*
- c. *If the cost is medium (i.e., $\underline{c} < c < \bar{c}$), high-type consumers manipulate their data with a probability $\phi = \frac{(1 - \alpha)((H - L)(H\rho + \sqrt{8c\rho + H^2\rho^2}) - 8c)}{8\alpha c}$. Prices are $p_H = \frac{H}{2}$, $p_L = \frac{\alpha\phi H + (1 - \alpha)L}{2\alpha\phi + 2(1 - \alpha)}$, and $p_\emptyset = \frac{\alpha H + (1 - \alpha)L}{2}$, and profit is $\Pi = \rho \left(\frac{\alpha(1 - \phi)H^2}{4} + \frac{(\alpha\phi H + (1 - \alpha)L)^2}{4(\alpha\phi + 1 - \alpha)} \right) + (1 - \rho) \frac{(\alpha H + (1 - \alpha)L)^2}{4}$.*

When the cost to manipulate data is prohibitive (i.e., $c \geq \bar{c}$), no consumers manipulate data, and the results (i.e., part a) are the same as those of the no-manipulation benchmark in Proposition 1. As firms collect more data on consumers and become more sophisticated in data analysis, they will be able to infer consumers' preferences from many other sources, such as their purchase and browsing history, which makes data manipulation much more difficult, if not impossible. In this case, successful data manipulation becomes sufficiently costly and consumers give up manipulating data.

When the cost is minimal (i.e., $c \leq \underline{c}$), all high-type consumers manipulate their data to mimic that of low-type consumers. Essentially, then, all consumers that the firm's data identifies appear to be the low-type, which renders the data just as uninformative as the consumers remaining unidentified. So, the firm charges these consumers the same price based on the expected consumer type.

When the cost of manipulation is medium (i.e., $\underline{c} < c < \bar{c}$), high-type consumers manipulate their data with a probability ϕ :

$$\phi = \frac{(1 - \alpha) \left((H - L)(H\rho + \sqrt{8c\rho + H^2\rho^2}) - 8c \right)}{8\alpha c}. \quad (5)$$

We could interpret ϕ as the fraction of high-type consumers who manipulate their data. The value of ϕ also reflects high-type consumers' manipulation incentive.

Intuition suggests that a lower cost to manipulate data and avoid being identified by the firm should benefit consumers. This intuition further implies that consumer surplus decreases as the cost of manipulation increases. However, Proposition 3 suggests that this intuition may not be correct.

Proposition 3 *A higher cost for consumers to manipulate their data (i.e., an increase in c) weakly increases the firm's profit; it also weakly increases consumer surplus as long as the cost is not too low (i.e., $c > \underline{c}$).*

When the cost of manipulation is so low (i.e., $c \leq \underline{c}$) that all high-type consumers manipulate their data, these consumers all incur the cost of manipulation. A lower cost of manipulation, therefore, increases consumer surplus. On the other hand, when the cost of manipulation is high (i.e., $c \geq \bar{c}$) and no consumers manipulate their data, the cost of manipulation is irrelevant and does not affect consumer surplus.

Now, an interesting result manifests when the cost is medium (i.e., $c \in (\underline{c}, \bar{c})$): high-type consumers randomize in their manipulation decisions and manipulate with a probability ϕ . Proposition 3 reveals that in this case, a lower cost for consumers to manipulate their data decreases consumer surplus; in other words, consumer surplus strictly increases with c (see Figure 1). We can understand this result by analyzing how a high-type and a low-type consumer's surplus changes with c separately.

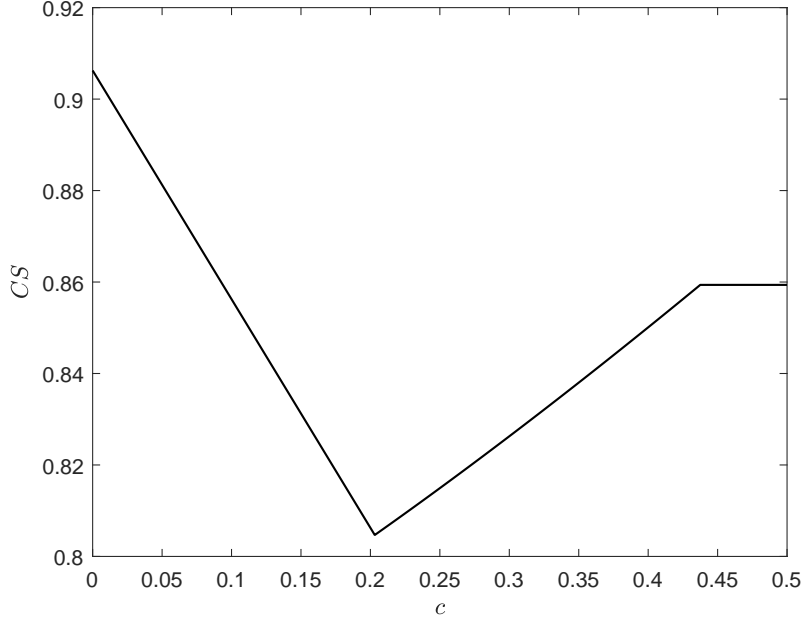


Figure 1: Consumer surplus changes with manipulation cost c ($H = 3, L = 2, \alpha = 0.5, \rho = 0.5$)

Consider first a high-type consumer's surplus when he or she does not manipulate data:

$$CS_H = \rho \int_{p_H}^H (v - p_H) dv + (1 - \rho) \int_{p_\emptyset}^H (v - p_\emptyset) dv = \frac{\rho H^2 + (1 - \rho)((2 - \alpha)H - (1 - \alpha)L)^2}{8}, \quad (6)$$

which does not depend on c . If a high-type consumer does manipulate, the expected surplus will be

$$CS_{HL} = \rho \int_{p_L}^H (v - p_L) dv + (1 - \rho) \int_{p_\emptyset}^H (v - p_\emptyset) dv - c, \quad (7)$$

where p_L is given in Proposition 2. Calculations suggest that $CS_{HL} = CS_H = (\rho H^2 + (1 - \rho)((2 - \alpha)H - (1 - \alpha)L)^2)/8$, which again does not depend on c . Note that $CS_{HL} = CS_H$ is not a coincidence because the high-type consumer's willingness to mix implies that the consumer must be indifferent about manipulating and not; therefore, in this case, an increase in c has no effect on a high-type consumer's surplus.

As for low-type consumers, they never manipulate so the changes in c have no direct effects on them. Instead, c exerts an indirect effect of changing the price that these low-type consumers pay. Specifically, an increase in c discourages high-type consumers from manipulating their data

to mimic those of low-type consumers. As a result, the low-type consumers that the firm identifies become more informative, accurately revealing true low-type consumers. Consequently, the firm reduces the price that it charges to consumers identified as low-type (i.e., $\partial p_L / \partial c = \partial p_L / \partial \phi \cdot \partial \phi / \partial c < 0$), which benefits them.

In sum, when c is in a medium range so that only a fraction of high-type consumers manipulates their data, a higher cost of manipulation has no consequences or effects on high-type consumers' surplus but it does increase low-type consumers' surplus. As a result, consumer surplus, as a whole, increases with the cost of manipulation. Over the entire region of c , consumer surplus first decreases and then weakly increases with c (see Figure 1).

Finally, a higher cost for consumers to manipulate data is weakly beneficial to firms. The intuition is related to the rent-seeking behavior in economics; that is, high-type consumers manipulate data to pay the price targeted to low-type consumers and obtain a rent, which reduces the firm's ability to perform price discrimination. Specifically, as c increases, consumers manipulate their data less often, i.e., $\partial \phi / \partial c < 0$, and the firm's collected data become more informative. Then, the firm has a greater ability to customize prices to target different types of consumers accurately. This result suggests that public policies or technological tools that enable consumers to manipulate their data more cost-efficiently can end up hurting both consumers and firms.

5 Endogenous Consumer Coverage

A firm can collect more consumer data to expand the consumer coverage of its big data (i.e., ρ). In this section, we endogenize the consumer coverage of the firm's big data to understand how the firm should make this strategic decision. In this case, the value of ρ is not necessarily public information; instead, the firm can decide whether or not to voluntarily disclose how much data and what types of data it collects so consumers can assess the firm's capability to identify consumer types and infer the consumer coverage of the firm's big data. Alternatively, public policies may make these decisions for the firm by mandating firms to disclose their collection and usage of consumer data to the public. To assess how firms make voluntary disclosure decisions or how public policymakers should regulate firms' disclosure of data collection, we analyze the respective scenario when consumers can or cannot observe the firm's data collection.

5.1 Consumers Observe the Firm's Data Collection

When consumers observe the firm's data collection to infer the consumer coverage of its big data ρ , as Proposition 2 shows, the firm's equilibrium profit is

$$\Pi = \rho \left(\frac{\alpha(1-\phi)H^2}{4} + \frac{(\alpha\phi H + (1-\alpha)L)^2}{4(\alpha\phi + 1-\alpha)} \right) + (1-\rho) \cdot \frac{(\alpha H + (1-\alpha)L)^2}{4}, \quad (8)$$

where ϕ is the probability that a high-type consumer manipulates his or her data. From Proposition 2, we can rewrite the consumers' manipulative behavior as a function of ρ , where $\phi = 1$ ($\phi = 0$) if the high-type consumers always (never) manipulate data. More specifically, we have

$$\phi = \begin{cases} 0 & \text{if } \rho \leq \frac{8c}{(H-L)(3H-L)}, \\ \frac{(1-\alpha)((H-L)(H\rho + \sqrt{8c\rho + H^2\rho^2}) - 8c)}{8\alpha c} & \text{if } \frac{8c}{(H-L)(3H-L)} < \rho < \frac{8c}{(1-\alpha)(H-L)((3-\alpha)H - (1-\alpha)L)}, \\ 1 & \text{if } \rho \geq \frac{8c}{(1-\alpha)(H-L)((3-\alpha)H - (1-\alpha)L)}. \end{cases} \quad (9)$$

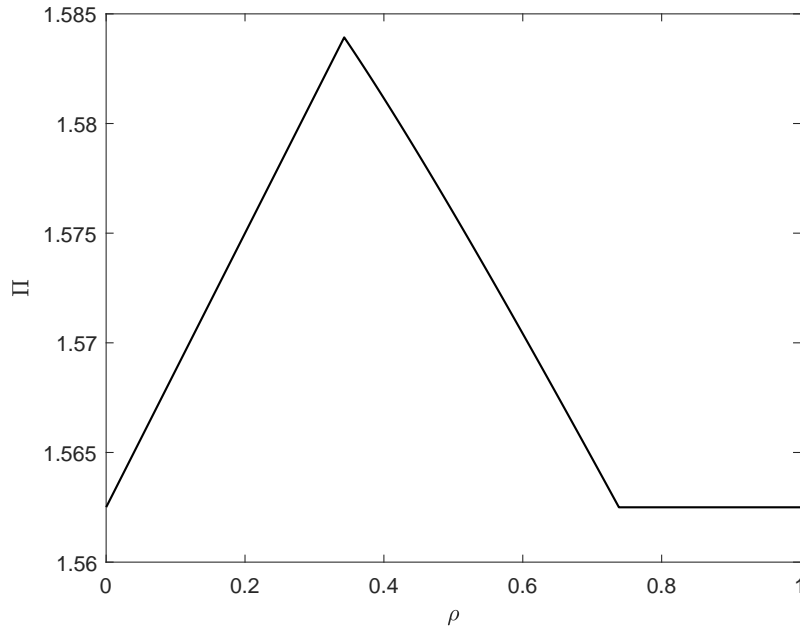


Figure 2: The firm's profit changes with investment in data collection ρ ($H = 3, L = 2, \alpha = 0.5, c = 0.3$)

Figure 2 depicts how the firm's profit changes with the consumer coverage ρ of its big data. The firm chooses a ρ that maximizes its profit. We summarize the result in Proposition 4:

Proposition 4 *If consumers observe the firm's data collection to infer ρ , $c < \frac{(H-L)(3H-L)}{8}$, and $\alpha \leq \frac{3H-L}{4H-2L}$:*

- a. *The firm chooses a partial consumer coverage; that is, $\rho^* = \frac{8c}{(H-L)(3H-L)} < 1$.*
- b. *The optimal consumer coverage ρ^* increases with c and decreases with consumers' demand heterogeneity (i.e., $H - L$).*

Proposition 1 shows that, when consumers cannot manipulate their data, the firm always benefits from increasing the consumer coverage of its big data and identifying more consumers' type. In contrast, Proposition 4 suggests that when consumers can manipulate their data, increasing the consumer coverage of a firm's big data can actually hurt the firm and decrease its profit. Therefore, a firm may deliberately forfeit collecting data to identify the types of more consumers, even if it is costless to do so.

To understand this result, note that an increase in ρ has two effects on the firm's profit: First, a higher ρ suggests that the firm can identify the types of more consumers, which allows it to offer customized prices to more consumers, which benefits the firm. Second, a higher ρ also implies that high-type consumers are more likely to be identified by the firm and charged a high price. If high-type consumers can observe the value of ρ , they know how an increase in ρ also increases the probability of the firm identifying them and charging them a high price, which gives them a stronger incentive to manipulate their data to mimic low-type consumers, i.e., $\partial\phi/\partial\rho > 0$ (see Equation 5), which hurts the firm. Then, in the specified parameter range, the second effect dominates the first effect, and the firm is better off acquiring fewer consumers' data in its database that does not identify all consumers' type. This result is consistent with observations of business data collection practice.⁴

Note that ρ^* , or the optimal consumer coverage of the firm's big data, increases with c . This is because, when c increases, high-type consumers are less likely to manipulate their data. Then, consumer data become more informative in guiding the firm's price discrimination decisions, and the firm is incentivized to invest more in data collection.

Interestingly, ρ^* decreases as the heterogeneity (i.e., $H - L$) increases amongst consumers. In other words, when high-type and low-type consumers exhibit increasingly different demands,

⁴Outside of this parameter range, when the cost of manipulation is so high (i.e., $c \geq \frac{(H-L)(3H-L)}{8}$), no consumers manipulate their data regardless of what the value of ρ is. When the manipulation cost is low (i.e., $c < \frac{(H-L)(3H-L)}{8}$) but a sufficient number of consumers are high type (i.e., $\alpha > \frac{3H-L}{4H-2L}$), there are no tractable solutions to ρ^* .

the firm has fewer incentives to invest in data collection. While this result seems counterintuitive, the rationale is as follows: when $H - L$ increases, the prices charged to different consumer types become more dispersed, giving high-type consumers a greater incentive to manipulate their data to mimic low-type consumers. The consumer data then become less informative and valuable to the firm. Anticipating this, the firm has fewer incentives to collect more data and, as a result, chooses a lower level of consumer coverage ρ .

5.2 Consumers Do Not Observe the Firm's Data Collection

Now, consider the scenario in which consumers do not observe what types of data and how much data the firm collects to infer the consumer coverage of its big data. Because consumers do not directly observe ρ , they rely on $\tilde{\rho}$, their belief of ρ when making their manipulation decisions. As such, consumers do not respond to ρ , and this is the main difference between the case in which consumers observe the firm's data collection or not. However, consumers' beliefs must be correct along the equilibrium path, i.e., $\tilde{\rho} = \rho$ in equilibrium. Thus, consumers should be able to rationally anticipate the firm's choice and act accordingly.

In this case, a high-type consumer's expected surplus from no manipulation and manipulation is

$$\widetilde{CS}_H = \tilde{\rho} \int_{p_H}^H (v - p_H) dv + (1 - \tilde{\rho}) \int_{p_\emptyset}^H (v - p_\emptyset) dv, \quad (10)$$

and

$$\widetilde{CS}_{HL} = \tilde{\rho} \int_{p_L}^H (v - p_L) dv + (1 - \tilde{\rho}) \int_{p_\emptyset}^H (v - p_\emptyset) dv - c, \quad (11)$$

respectively. A high-type consumer prefers to manipulate his or her data if and only if $\widetilde{CS}_H < \widetilde{CS}_{HL}$. Following the discussion above, a high-type consumer's probability of manipulation, ϕ , does not depend on ρ but depends on $\tilde{\rho}$ instead, i.e., $\partial\phi/\partial\rho = 0$. The firm's profit is

$$\Pi = \rho \left(\frac{\alpha(1 - \phi)H^2}{4} + \frac{(\alpha\phi H + (1 - \alpha)L)^2}{4(\alpha\phi + 1 - \alpha)} \right) + (1 - \rho) \frac{(\alpha H + (1 - \alpha)L)^2}{4}, \quad (12)$$

and the firm chooses ρ to maximize its profit. Note that consumers hold rational expectations of ρ ,

that is, $\tilde{\rho} = \rho$ in equilibrium (but not off the equilibrium). By solving the firm's profit maximization problem, we obtain the following proposition:

Proposition 5 *If consumers do not observe the firm's data collection, the firm chooses the full consumer coverage (i.e., $\rho^* = 1$).*

Proposition 1 implies that, when consumers cannot manipulate their data, the firm is better off collecting more consumer data to pursue the full consumer coverage. Proposition 4 indicates that the firm should cut back on its data collection if consumers can manipulate their data and observe the firm's data collection to infer ρ . Now, Proposition 5 shows that the firm should revert to pursuing the full consumer coverage if consumers can manipulate their data but cannot observe the firm's data collection.

The intuition behind Proposition 5 is as follows: As shown above, consumer coverage of the firm's database has a direct effect and an indirect effect on the firm's profit. The direct effect is that it allows the firm to identify more consumers' type and charge more consumers customized prices; this price-discrimination effect increases the firm's profit. The indirect effect is that a higher capability to identify consumers' type induces more high-type consumers to manipulate their data, which makes the firm's data less accurate and reduces the effectiveness of its targeting and price discrimination strategies. When consumers can observe the firm's data collection to infer ρ , the firm carefully evaluates the trade-off between the two effects when choosing the consumer coverage of its big data.

However, when consumers cannot observe the firm's data collection, the indirect effect is no longer the firm's concern, because consumers do not observe ρ and, thus, do not respond to it. As a result, the firm can increase the consumer coverage of its big-data technology without encouraging consumers to manipulate their data. As a result, the firm increases its data collection to pursue the full consumer coverage.

5.3 Disclosure of Data Collection

We considered the firm's data collection decisions when consumers observe or do not observe such decisions. We showed that, while the firm may deliberately choose not to collect all consumers' data to identify their type (i.e., pursue a partial consumer coverage) when consumers can observe

its data collection, the firm always chooses to collect all available data to identify all consumers' type when consumers cannot observe its data collection. Rational consumers take this into account when making their manipulation decisions.

In this section, we compare the equilibrium profit, consumer surplus, and social welfare under observability and unobservability of data collection to see their effects on the firm, consumers, and society. This analysis helps managers and public policymakers assess the implications of disclosing a firm's data collection. Without laws or regulations, the firm's disclosure is cheap talk and cannot be verified, and therefore consumers do not take disclosure into account. With regulations, however, firms must truthfully disclose their collection of consumer data. If the firm lies about its data collection, it could damage the brand's image, or the firm may be subject to other penalties. For example, the GDPR sets a maximum fine of €20 million or 4% of annual global turnover, whichever is greater, if a company does not comply with its data disclosure requirements (Itgovernance 2020). Thus, we focus on truthful disclosures that are either regulated or verifiable.

Proposition 6 *Disclosing the firm's data collection (weakly) increases the firm's profit, consumer surplus, and social welfare.*

While it may be intuitive that a firm cannot be worse off disclosing its data collection to consumers, it is *a priori* unclear whether such disclosure benefits consumers and society as a whole. For example, Li et al. (2020) find that firms' disclosure of data collection often works to the detriment of consumers and society. Proposition 6, by contrast, shows that disclosing data collection can also benefit consumers. This result suggests that public policymakers should consider mandating firms to disclose their data collection to consumers. Moreover, we find out that such disclosure can strictly benefit the firm under certain circumstances. While we assume that information disclosure is costless, this result implies that the firm may find it profitable to disclose such information even under a positive disclosure cost (Guo and Zhao 2009).

When consumers cannot observe the firm's data collection, the firm cannot help collect more data to pursue the full consumer coverage (see Proposition 5). Then, in equilibrium, consumers correctly anticipate that the firm's database has the full consumer coverage. As a result, high-type consumers have strong incentives to manipulate their data. And as consumers manipulate data with an increased probability, the firm's data become less informative, which offsets the firm's

increased ability to customize prices according to consumer types. However, the firm still takes the opportunity to choose a more powerful technology, even though doing so ultimately hurts its own profit. Therefore, the firm is actually better off disclosing its data collection to consumers, which helps it commit to choosing a partial consumer coverage that dissuades consumers from manipulating their data.

Proposition 6 also shows that observing the firm's data collection leaves consumers better off. The rationale is that, without observability, the firm cannot help collecting more data to identify the types of all consumers. As a result, high-type consumers are more likely to be identified and price-discriminated against, so they manipulate their data more often to counteract the increasing likelihood of being accurately identified by the firm. In doing so, these consumers incur a dead-weight loss of manipulation since they incur the cost to manipulate their data, which decreases consumer surplus. Moreover, low-type consumers also suffer as more high-type consumers engage in data manipulation, because manipulation raises the prices that low-type consumers pay. By contrast, when the firm discloses its data collection, it reduces data collection, which gives high-type consumers fewer incentives to incur the cost of manipulation. Consumer surplus increases as a result. And, because ρ is lower, the firm price discriminates less, which also benefits consumers.

Finally, because disclosing data collection benefits both the firm and consumers, this action increases social welfare accordingly. Therefore, this finding suggests that firms should voluntarily disclose their data collection, thereby increasing their profit and consumer surplus. Public policies should mandate firms to disclose their data collection, which is mutually beneficial to firms and consumers.

6 Extensions

In this section, we relax several assumptions made in the base model to show the robustness of our key findings and provide new insights.

6.1 Continuous Consumer Types

In the base model, we model consumer heterogeneity by assuming that a consumer's type is drawn from a two-point distribution. In this section, we generalize the model to allow for continuous heterogeneity in consumer type. In the base model with two consumer types, the consumer's manipulation decision is binary. In contrast, with consumer types are continuously distributed, high-type consumers need to carefully determine the behavior of which type of consumers they would like to mimic. Let a consumer's type be drawn from a continuous distribution. To obtain tractable results, assume that consumer i 's demand function is $D_i = t_i - p$, where t_i is drawn from a uniform distribution over $[L, H]$, with $H > L > 0$ and $L \geq \frac{H}{2}$. If consumer i manipulates his or her data, he or she can choose to mimic any other consumer type $\tilde{t}_i \in [L, H]$ at the same cost, c . Otherwise, the consumer's data reveal his or her true type, i.e., $\tilde{t}_i = t_i$.

In this context, the firm's big-data technology operates similarly as it does in the base model. For each consumer i , the technology generates a signal $s_i \in [H, L] \cup \{\emptyset\}$, with

$$s_i = \begin{cases} \tilde{t}_i, & \text{with probability } \rho, \\ \emptyset, & \text{with probability } 1 - \rho. \end{cases} \quad (13)$$

That fraction ρ indicates the consumer coverage of the firm's database. The game proceeds in the same sequence as in the base model.

We solve for the equilibrium outcome and summarize it in the following proposition:

Proposition 7 *With continuous consumer types, consumers' manipulative behavior and the firm's prices vary with c as follows:*

- a. *If the cost is high (i.e., $c \geq \bar{c} = \rho(3H - L)(H - L)/8$), no consumers manipulate data.*
- b. *Otherwise, there exist $L \leq x_L \leq x_H \leq H$ where x_L and x_H are defined implicitly by $\frac{H^2 - L^2 - x_H^2 + x_L^2}{2(H - L - x_H + x_L)} = x_L$ and $c = \frac{\rho(x_H - x_L)(3x_H - x_L)}{8}$:*
 - (1) *Only high-type consumers (i.e., $t_i \in [x_H, H]$) manipulate.*
 - (2) *The type of consumers that high-type consumers mimic \tilde{t}_i follows the probability density func-*

tion that

$$f(x) = \begin{cases} \frac{2(x-x_L)}{(H-x_H)(2x_L-H-x_H)}, & \text{if } L \leq x \leq x_L, \\ 0, & \text{otherwise.} \end{cases} \quad (14)$$

The firm's pricing strategy is:

$$p(s_i) = \begin{cases} x_L/2, & \text{if } L \leq s_i \leq x_L, \\ s_i/2, & \text{if } x_L < s_i \leq H, \\ (H+L)/4, & \text{if } s_i = \emptyset. \end{cases} \quad (15)$$

The implicit function $c = \frac{\rho(x_H-x_L)(3x_H-x_L)}{8}$ is derived from

$$-c + \rho \int_{\frac{x_L}{2}}^{x_H} \left(v - \frac{x_L}{2}\right) dv + (1-\rho) \int_{p_\emptyset}^{x_H} (v - p_\emptyset) dv = \rho \int_{\frac{x_H}{2}}^{x_H} \left(v - \frac{x_H}{2}\right) dv + (1-\rho) \int_{p_\emptyset}^{x_H} (v - p_\emptyset) dv, \quad (16)$$

which implies that a consumer of type x_H is indifferent about whether or not to manipulate. Next, upon observing a consumer's type $x \in [L, x_L]$, the firm believes that the consumer's expected type is $E[t_i|x] = x_L$. Mathematically, this presents as

$$\frac{E[t_i|t_i \geq x_H](H-x_H)f(x) + x}{(H-x_H)f(x) + 1} = x_L, \quad (17)$$

From Equation (17), we derive function $f(\cdot)$. Finally, for $f(\cdot)$, we must have

$$\int_L^{x_L} f(x) dx = 1. \quad (18)$$

This gives the implicit function $\frac{H^2-L^2-x_H^2+x_L^2}{2(H-L-x_H+x_L)} = x_L$.

If the cost of manipulation is not prohibitive, a segment of high-type consumers (i.e., $t_i \in [x_H, H]$) manipulates their data to pool together with a segment of low-type consumers (i.e., $t_i \in [L, x_L]$). Note that, with continuous consumer types, high-type consumers randomize in the type of consumers that they mimic. One may wonder why these consumers do not always mimic the lowest-type consumer whose $\tilde{t}_i = L$. The reason is, if they all mimicked the lowest-type consumer whose $\tilde{t}_i = L$, the firm, upon observing $s_i = L$, will rationally expect that this identified

consumer type is likely to be fake and that some, if not most, of those consumers are actually high-type consumers, i.e., $E(t_i|s_i = L)$ is high. As a result, the firm charges all these consumers a high price. On the other hand, because no consumers would choose to manipulate their data to mimic the type $\tilde{t}_i = L + \epsilon$ for some small $\epsilon > 0$, the firm believes that consumers who are identified with a slightly higher type $L + \epsilon$ are truly this type (i.e., $E[t_i|s_i = L + \epsilon] = L + \epsilon$). Therefore, the firm charges this type of consumers a lower price than what it charges consumers identified to be the lowest type (i.e., $E[t_i|s_i = L + \epsilon] < E[t_i|s_i = L]$ and $p_L > p_{L+\epsilon}$), which induces high-type consumers to mimic consumers whose type is $L + \epsilon$. Therefore, fabricating data to the lowest type $\tilde{t}_i = L$ is not always in the best interest of high-type consumers. In equilibrium, high-type consumers mimic the segment of low-type consumers according to the distribution that Equation (14) specifies. In addition, consumers with moderate preferences, i.e., $t_i \in (x_L, x_H)$, are not pooled with any other type. If they manipulate their data, they can potentially get a better deal; however, the potential gain from manipulation is not sufficient to cover its cost. As such, they do not manipulate in equilibrium. As for low-type consumers with $t_i \leq x_L$, manipulating never renders them a lower price and they do not manipulate in equilibrium.

Interestingly, the firm's pricing strategy with continuous consumer types is a hybrid of segment-specific pricing and individual-specific pricing. The firm charges the segment of low-type consumers (i.e., $s_i \leq x_L$) the same price and, therefore, high-type consumers are willing to mix. For medium-type consumers (i.e., $s_i > x_L$), the firm charges each consumer a personalized price. This hybrid pricing strategy is different from a segment-specific pricing strategy used in behavior-based pricing by which a firm charges two segments of repeat and new customers different prices (Fudenberg and Tirole 2000). It is also different from individual-specific personalized pricing (Chen and Iyer 2002).

The rationale behind this hybrid pricing equilibrium is as follows. First, if a consumer with type $t_i = x$ manipulates data, any consumer with type $t_i > x$ will also manipulate. Therefore, there must exist a threshold x_H above which all consumers manipulate. Second, all consumers who manipulate their data will fabricate them to get the lowest possible price, p_{\min} ; in other words, they end up getting the same price. Third, as discussed above, it is irrational for all high-type consumers to mimic consumers with the lowest type $\tilde{t}_i = L$, which only boosts up the price p_L they pay. In the equilibrium characterized above, high-type consumers mimic consumers with

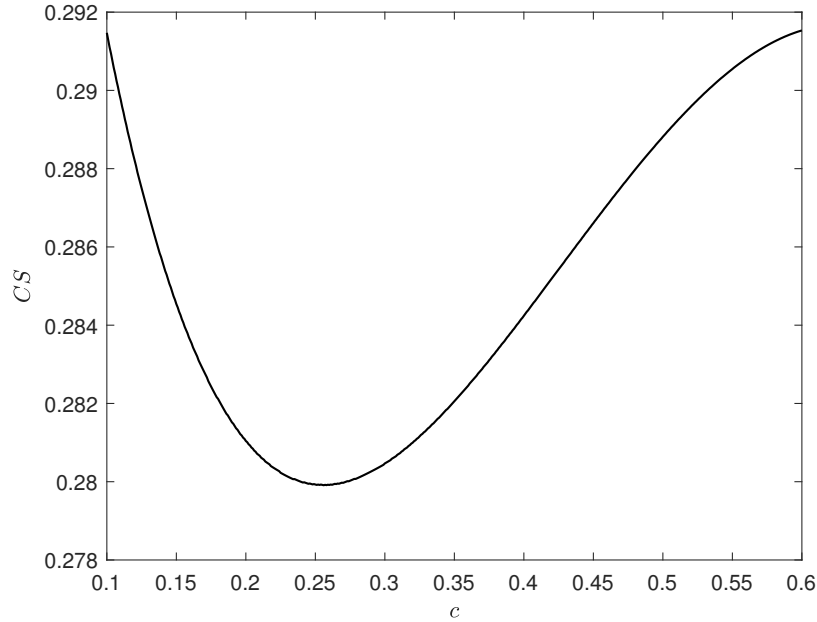


Figure 3: Consumer surplus changes with manipulation cost c ($H = 2, L = 1, \rho = 1$)

$t_i \in [L, x_L]$, and, as a result, the firm must offer all consumers with these identified types the same price $p_{t_i} = p_{\min}$. It follows, then, that consumers with $t_i \in [L, x_L] \cup [x_H, H]$ get the same price p_{\min} . Finally, consumers with $t_i \in [x_L, x_H]$ are perfectly recognized: On the one hand, their types are not so low that no high-type consumers are willing to mimic them; on the other hand, their types are not so high that they themselves do not have incentives to mimic others. As a result, their types are truthfully revealed to the firm. Thus, the firm charges them personalized prices based on their individual-specific true type. As discussed above, in equilibrium, consumers with types $t_i \in [L, x_L] \cup [x_H, H]$ end up paying the same price $p_{\min} = \frac{x_L}{2}$, while consumers with types $t_i \in (x_L, x_H)$ end up paying $p = \frac{t_i}{2} > p_{\min}$. This result suggests that the final price charged by the firm is non-monotone in t_i : consumers with a moderate t_i pay the highest price.

Because of our model's complexity, we could not analytically examine comparative statics pertaining to the effects of c and ρ on the equilibrium outcomes. Instead, we illustrate our key findings with numerical examples: First, let $L = 1, H = 2$, and $\rho = 1$. We present consumer surplus in Figure 3. It follows immediately that consumer surplus can increase with c , thereby replicating our first main result. Second, let $L = 1, H = 2$, and $c = 0.1$, Figure 4 illustrates

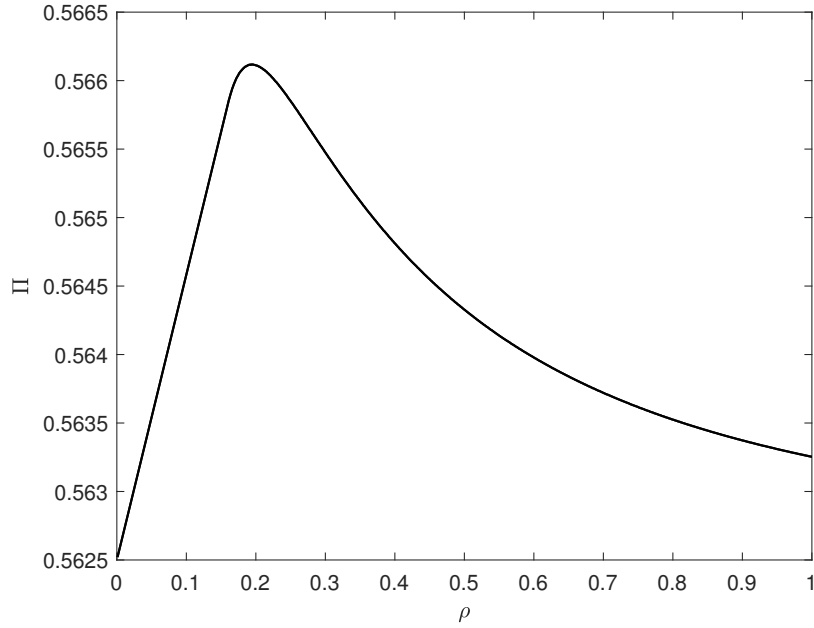


Figure 4: The firm's profit changes with investment in data collection ρ ($H = 2, L = 1, c = 0.1$)

how the firm's profit changes with ρ . We see from the figure that the firm's profit is, once again, nonmonotone in ρ , and is maximized at a moderate ρ . Therefore, we replicate our second main result that the firm can be worse off with a high ρ .

6.2 Consumer Hiding

In the base model, we examine how a firm's targeted pricing and data collection change when high-type consumers deliberately manipulate their data to mimic low-type consumers. However, some consumers may use various technologies to simply maintain anonymity and avoid being identified by firms (Conitzer et al. 2012; Valletti and Wu 2020). In this section, we consider how our results may change with consumer hiding instead of consumer manipulation.

Consider the following model: As before, α of the consumers are high-type consumers while $1 - \alpha$ of them are low-type consumers, and we use $t_i \in \{H, L\}$ to denote the consumer type. We use $\tilde{t}_i \in \{H, L, \emptyset\}$ to denote the consumer data profile: If a consumer does not hide, then $\tilde{t}_i = t_i$; otherwise $\tilde{t}_i = \emptyset$. The firm collects data from ρ of the consumers, and receives a signal s_i about

consumer i 's type:

$$s_i = \begin{cases} \tilde{t}_i, & \text{with probability } \rho, \\ \emptyset, & \text{with probability } 1 - \rho. \end{cases} \quad (19)$$

Note that in the current model, when the firm receives a signal $s_i = H$ ($s_i = L$), the firm knows for sure that consumer i is a high (low) type consumer. However, if the firm receives a signal $s_i = \emptyset$, there are two possibilities: (1) the consumer's type is unknown (e.g., not covered by the firm's data) or (2) the consumer hides her information to mimic the unknown type. The firm is unable to distinguish between these two scenarios.

We obtain the following insights from the model. First, assume that the consumer coverage of the firm's big data is exogenously given: When consumers do not hide, identified high-type consumers receive the price $p_H = H/2$ while unidentified consumers receive the price $p_\emptyset = (\alpha H + (1 - \alpha)L)/2 < p_H$. This price difference gives high-type consumers an incentive to hide their identities and pool together with unidentified consumers. As more high-type consumers decide to take this route, the firm believes that an unidentified consumer is more likely a high-type consumer and raises the price p_\emptyset . Therefore, high-type consumers' hiding behavior exerts negative externalities on unidentified consumers. Thus, as with the base model, the total consumer surplus can increase with the cost of hiding.

Next, assume that the firm chooses the consumer coverage of its big data: Unlike the base model, now the firm chooses $\rho = 1$. To understand this, consider the case of $\rho = 1$. Because low-type consumers already get a low price $p_L = L/2$, they have no incentives to hide and are perfectly recognized by the firm. Given that no low-type consumers hide, the firm classifies any unidentified consumers as high-type consumers and offers them the price $p_H = H/2$. In other words, the firm still perfectly recognizes all consumers and maximizes its profit.⁵

⁵We also analyze the case when consumers have heterogeneous abilities to manipulate their data and discuss how our intuitions can hold in a competitive market (for details, see the Appendix).

7 Conclusions

As big-data and information technologies advance, firms continue to invest heavily in these tools to collect consumer data from multiple sources. They also use machine learning and AI algorithms to analyze big data, identify consumer types, and infer consumer preferences so they may tailor the prices they charge to specific consumer types. Meanwhile, consumers have become increasingly aware that firms collect and use their data for targeted pricing. As a result, they can take action to manipulate their data and avoid being price-discriminated against. In this paper, we examine the tension between a firm's efforts to identify consumers using consumer data and its consumers' efforts to manipulate their data to get a better deal. Furthermore, we examine whether or not firms should disclose their data collection to consumers. Our analysis shows that the firm's optimal data collection scope depends on whether or not consumers can manipulate their data and observe the firm's data collection. A firm's customized prices vary with its consumers' cost of manipulation. A higher cost for consumers to manipulate their data can benefit both the firm and consumers. Disclosure of data collection improves the firm's profit, consumer surplus, and social welfare. Our key insights and managerial implications are as follows.

How should a firm invest in data collection to identify consumers who can manipulate their data? Our analysis shows that the firm's optimal data collection scope depends on whether or not consumers can 1) manipulate their data and 2) observe the firm's data collection. When consumers cannot manipulate their data or observe the firm's data collection, the firm's profit increases with the consumer coverage of its big data, suggesting that the firm should invest more in data collection to identify the types of more consumers. However, when consumers can manipulate their data and observe the firm's data collection, the firm should not identify all consumers, even when collecting more data to identify consumers is costless. These findings warn managers that increasing data collection can sometimes backfire on them.

Should public policies mandate a firm to disclose its data collection to consumers? Should a firm voluntarily disclose its data collection to consumers even without mandatory disclosure policies? Our analysis shows that the firm's profit, consumer surplus, and social welfare are all higher when consumers can observe the firm's data collection. Therefore, public policies that mandate firms to disclose their data collection to consumers help all parties. Even in the absence of such mandatory-

disclosure policies, firms should voluntarily disclose their data collection to consumers to improve profit.

How should a firm customize prices according to consumer data when consumers can manipulate their data? Our findings suggest that a firm's customized prices vary with its consumers' cost of manipulation. As the cost decreases and more high-type consumers manipulate their data to mimic low-type consumers, the firm should raise the price that it charges to consumers who the firm's data identifies to be a low-type.

How does consumers' cost to manipulate their data affect the firm and consumers? We find that a higher cost for consumers to manipulate their data can benefit both the firm and consumers. This finding suggests that the emergence of new devices, training programs, and technology that cost-efficiently help consumers get better deals can actually end up being detrimental to consumers, the firm, and society as a whole. Similarly, public programs that educate or encourage consumers to protect themselves and get better deals by manipulating their data can lead to unintended consequences.

How do consumers' manipulative behavior and a firm's targeted pricing strategies change when consumers have continuous types? Our analysis shows that, when consumers exhibit continuous types of preferences instead of the two discrete types and the cost of manipulation is not prohibitive, a segment of high-type consumers mimics a segment of low-type consumers according to a distribution, while medium-type consumers do not manipulate. The firm uses a hybrid pricing strategy by charging low-type consumers and high-type consumers who manipulate the same segment-specific price while charging medium-type consumers an individual-specific price.

Our study offers several directions for future researchers. First, we analyze how a monopolist makes data collection and disclosure decisions to keep our analyses tractable. Although our conjecture suggests that the main results and intuitions should hold in a competitive setting, future research could formally analyze a model with competition to uncover how competition generates new insights. Second, our model with continuous consumer types provides interesting results on consumers' manipulative behaviors and firms' pricing strategies. Future research could build on this model to examine how firms make other strategic decisions. Lastly, our research examines how firms utilize consumer data for targeted price discrimination that hurts consumers. Firms could also use data to understand consumer preferences and taste to design better products or

more effective advertising, which could benefit consumers (Koh et al. 2017; Ichihashi 2020). Future research could explore firms' big-data strategies and public policy recommendations on data collection that could provide value to consumers.

Acknowledgement

Xi Li acknowledges financial support from the Hong Kong Research Grants Council (Grant 21500920). Krista J. Li acknowledges the financial support from the Weimer Faculty Fellowship.

The authors thank the department editor, the associate editor, two anonymous referees, and seminar participants at the Hong Kong Polytechnic University, University of Hong Kong, University of Science and Technology of China, Zhongnan University of Economics and Law, and the 2022 China Marketing International Conference.

References

- Anderson, S. P. and Renault, R. 2009, Comparative advertising: disclosing horizontal match information. *RAND Journal of Economics*, 40, 558–581.
- Argenziano, R. and Bonatti, A., 2021. Data Linkages and Privacy Regulation. Working Paper, MIT Sloan School of Management.
- Blog JIPEL 2018. Post-GDPR: Will the U.S. implement a comprehensive data privacy law? *JIPEL*. Available at <https://blog.jipel.law.nyu.edu/2018/11/post-gdpr-will-the-u-s-implement-a-comprehensive-data-privacy-law/>.
- Belleflamme, P. and Vergote, W., 2016. Monopoly price discrimination and privacy: The hidden cost of hiding. *Economics Letters*, 149, 141–144.
- Chen, Y. and Iyer, G. 2002. Research note consumer addressability and customized pricing. *Marketing Science*, 21(2), 197–208.
- Chen, Y., Li, X. and Sun, M., 2017. Competitive mobile geo targeting. *Marketing Science*, 36(5), 66–682.
- Chen, Y., Narasimhan, C. and Zhang, Z.J., 2001. Individual marketing with imperfect targetability. *Marketing Science*, 20(1), 23–41.
- Conitzer, V., Taylor, C.R. and Wagman, L., 2012. Hide and seek: costly consumer privacy in a market with repeat purchases. *Marketing Science*, 31(2), 277–292.
- Experian, 2014, Driving customer loyalty: maximize loyalty program data collection to drive insight and revenue. Available at <https://docplayer.net/63899516-Driving-customer-loyalty-maximize-loyalty-program-data-collection-to-drive-insight.html>.

- Fenton, R. 2021, Woman goes viral on TikTok with this Uber hack to get cheaper rides after night out. *Mirror*. Available at <https://www.mirror.co.uk/money/savings-banks/woman-goes-viral-tiktok-uber-23903949>.
- Freixas, X., Gueserie, R., and Tirole, J. 1985, Planning under incomplete information and ratchet effect, *Review of Economic Studies*, 52(2), 173-191.
- Fudenberg, D. and Tirole, J. 2000, Customer poaching and brand switching *The RAND Journal of Economics*, 31(4), 634-657.
- Fudenberg, D. and Villas-Boas, J. M. 2006. Behavior-based price discrimination and customer recognition. *Handbook on Economics and Information Systems*, 377-436, T.J. Hendershott, Ed., Elsevier.
- Grossman SJ (1981) The informational role of warranties and private disclosure about product quality. *J. Law Econom.* 24(3):461–483.
- GDPR EU 2020. A guide to GDPR data privacy requirements. *GDPR.EU*, available at <https://gdpr.eu/data-privacy/>.
- Goldfarb, A. and Tucker, C. 2010, Privacy regulation and online advertising. *Management Science*, 57(1): 57-71.
- Gu, Z. and Xie, Y. 2013, Facilitating fit revelation in the competitive market. *Management Science*, 59, 1196–1212.
- Guo, L. and Zhao, Y. 2009, Voluntary quality disclosure and market interaction. *Marketing Science*, 28(3):488–501.
- Heine, C. 2015. 170 U.S. brands are already using this ad tech that can target people in a specific building, *Adweek*, available at <https://www.adweek.com/digital/170-us-brands-are-already-using-ad-tech-can-target-people-specific-building-163272/>.
- Ichihashi, S. 2020. When online sellers use different prices for different consumers, *London School of Economics and Political Science*. Available at <https://blogs.lse.ac.uk/businessreview/2020/01/20/when-online-sellers-use-different-prices-for-different-consumers/>.
- Itgovernance 2020. GDPR fines: administrative fines and other penalties for non-compliance with the EU General Data Protection Regulation. Available at <https://www.itgovernance.co.uk/dpa-and-gdpr-penalties>.
- Jovanovic B (1982) Truthful disclosure of information. *Bell J. Econom.* 13(1):36–44.
- Koh, B., Raghunathan, S., and Nault, B. R. 2017. Is voluntary profiling welfare enhancing? *MIS Quarterly*, 41(1), 23–41.
- Laffont, J.-J. and Tirole, J., 1993. The dynamics of incentive contracts, *Econometrica*, 56(5), 1153-1175.
- Lau, S., 2008. Information and bargaining in the hold-up problem. *RAND Journal of Economics*, 39(1), 266–282.
- Li, X., Li, K.J. and Wang, X., 2020. Transparency of behavior-based pricing. *Journal of Marketing Research*, 57(1), 78–99.

- Mahdawi, A. 2016. Cookie monsters: why your browsing history could mean rip-off prices. *The Guardian*, available at <https://www.theguardian.com/commentisfree/2016/dec/06/cookie-monsters-why-your-browsing-history-could-mean-rip-off-prices>.
- Mattioli, D. 2012. On Orbitz, Mac users steered to pricier hotels. *The Wall Street Journal*, available at <https://www.wsj.com/articles/SB10001424052702304458604577488822667325882>.
- Montes, R., Sand-Zantman, W. and Valletti, T., 2019. The value of personal information in online markets with endogenous privacy. *Management Science*, 65(3), 1342–1362.
- Morgan, B. 2020. 50 stats that show the importance of good loyalty programs, even during a crisis. *Forbes*, available at <https://www.forbes.com/sites/blakemorgan/2020/05/07/50-stats-that-show-the-importance-of-good-loyalty-programs-even-during-a-crisis/#4c327ce12410>.
- Office Depot 2020. Office Depot's privacy statement - your privacy rights. *NewVantage Partners*, available at <https://www.officedepot.com/cm/help/privacy-statement>.
- Osborne, C. 2019. Fortune 1000 to 'urgently' invest in Big Data, AI in 2019 in fear of digital rivals. *ZD Net*, available at <https://www.zdnet.com/article/fortune-1000-to-urgently-invest-in-big-data-ai-in-2019-in-fear-of-digital-rivals/>.
- Perez, S. 2019. Target's personalized loyalty program launches nationwide next month. *Tech Crunch*. Available at <https://techcrunch.com/2019/09/09/targets-personalized-loyalty-program-launches-nationwide-next-month/>.
- Smith, L. 2009. Price discrimination and loyalty programs – a deadly combo *Retail Marketing Management*, available at <http://bus4411.blogspot.com/2009/03/price-discrimination-and-loyalty.html>.
- Target, 2020. Target privacy policy. Available at <https://www.target.com/c/target-privacy-policy/-/N-4sr7p?Nao=0#Product>.
- Taylor, C. R. 2004. Consumer privacy and the market for customer information. *The RAND Journal of Economics*, 35(4), 631-650.
- Thisse, J.F. and Vives, X., 1988. On the strategic choice of spatial price policy. *American Economic Review*, 122–137.
- Townley, C., Morrison, E., and Yeung, K. 2017, Big data and personalised price discrimination in EU competition law. King's College London Law School Research Paper No. 2017-38.
- Valletti, T. and Wu, J., 2020. Consumer profiling with data requirements: structure and policy implications. *Production and Operations Management*, 29(2), 309–329.
- Wathieu, L. and Bertini, M., 2007. Price as a stimulus to think: The case for willful overpricing. *Marketing Science*, 26(1), 118–129.