

# Share: Stackelberg-Nash based Data Markets

Anonymous Author(s)

## ABSTRACT

With the prevalence of data-driven research, data markets with data products in various forms are gaining considerable concerns due to the facilitation of data transaction. In this paper, we present **Stackelberg-Nash based Data Markets (Share)**, which is the first work to introduce both Stackelberg game and Nash game into data markets to realize absolute pricing for data. We propose a three-stage Stackelberg-Nash game to model trading dynamics which not only balances the profits of all participants in buyer-leading markets but also solves the seller selection problem based on sellers' inner competition. We define Stackelberg-Nash Equilibrium and apply two approaches to derive inner Nash equilibrium, conventional direct derivation and novel mean-field based method for complicated cases along with provable approximation guarantees. Experiments on real datasets justify the effectiveness and efficiency of *Share*.

## ACM Reference Format:

Anonymous Author(s). 2018. Share: Stackelberg-Nash based Data Markets. In *Woodstock '18: ACM Symposium on Neural Gaze Detection, June 03–05, 2018, Woodstock, NY*. ACM, New York, NY, USA, 16 pages. <https://doi.org/10.1145/1122445.1122456>

## 1 INTRODUCTION

With great success in a variety of data-driven tasks, data products (e.g., machine learning models, aggregate statistics, and query services) have paved the way for being more widely applied across multifarious industries. High-performance data products hinge on a large amount of high-quality data. Benefiting from the advantage of sensor intelligence, it is more convenient to excavate potential data from individuals, thus greatly enriching data resources. However, these data exploited from various sources are highly dispersed, which brings huge challenges to data aggregation. Besides, there is a gap between data supply and demand, and data demanders usually lack the necessary resources and techniques to turn data into data products, not to mention sufficiently digging out the reserved wealth from data. Thus, despite the remarkable enrichment of data, the wealth of data is far from being fully exploited. Recent studies and practices have approached the commoditization of data with data markets to solve the aforementioned issues. However, these efforts have not satisfied the following properties widely desired.

A typical data market consists of three parties, buyers, brokers, and sellers. Buyers propose demands for data products and pay for them. Brokers help facilitate the transactions between buyers and

sellers as well as take charge of manufacturing data products out of data. Sellers offer data with different quality due to privacy consideration and sell data to brokers in exchange for compensations.

*Medical Data Market Example.* Consider medical data trading. A drug company as a buyer demands for a data model (a kind of data products) based on real medical data to study the effectiveness of the new drug so as to make further business decisions. A broker buys data from hospitals, use the data to train the required model and sells the model to the company. Hospitals own numerous data generated from medical services and sell data to the broker with distinctive quality by privacy preserving mechanism based on the data price offered by the broker and the consent made with patients.

**Desired Properties.** All three parties have their own *revenue* and *cost*. It's imperative to maximize the profit for each of them in order to motivate their participation and boost market vitality. Therefore, it's necessary to design data markets which consider the profits of buyers, brokers, and sellers simultaneously ( $\mathcal{P}_1$ ).

In the above example, hospitals can sell data but never regard data selling as the main business. Rather, drug companies that heavily rely on real medical data to make business decisions on drug supply ask hospitals initiatively for a deal. It thus can be inferred that data transactions are often initiated by data demanders, i.e., buyers. Therefore, it's necessary to construct data markets with the property of buyer-leading (demand-leading) ( $\mathcal{P}_2$ ) which take full account of buyers' priority against others and can help trade data.

Pricing rules play a significant role in the market mechanism. In order to mirror real markets, prices should be absolute, instead of being pegged to a certain benchmark, to reflect actual values of goods. Besides, incentive pricing mechanisms should involve all three parties in the pricing process rather than force them to passively receive certain prices which may discourage their participation. Thus, it's necessary to formulate data markets in which absolute prices can be directly decided by market participants based on their mutual interactions ( $\mathcal{P}_3$ ).

Many researchers have designed data markets dealing machine learning models [4][10][30][38][39], data aggregation statistics [5][24][31][45], or query services [32][33][35]. These works vary in design goals, such as model quality optimization [24], revenue maximization [10], fair revenue allocation [30], social welfare maximization [31], or market protection from strategic participants [22]. While various data market models have been designed, little progress is made on buyer-leading data markets where each party can obtain profit optimization. In terms of pricing, most of the existing works [4][10][32][45][58] employed relative prices for data or data products, which fail to give a real picture of markets. It's therefore tempting to ask: *how to build a well-functioning data market with an absolute pricing mechanism where both the profit needs for all participants and the leading position for buyers are considered.*

**Challenges.** Although many researchers [35][38][39] have proposed three-side data markets which aim to simultaneously satisfy the needs of all three entities, they assume the brokers as neutral.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).

Woodstock '18, June 03–05, 2018, Woodstock, NY

© 2018 Association for Computing Machinery.

ACM ISBN 978-1-4503-XXXX-X/18/06...\$15.00

<https://doi.org/10.1145/1122445.1122456>

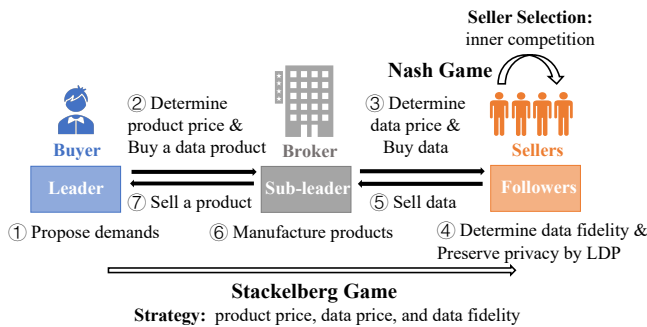
In practice, however, all participants, including the brokers, are *selfish*, i.e., aiming to maximize their own profit. Moreover, how they act in the market affects each other. Consider medical data market example again. If the drug company sets a low price for the desired model for profit, the broker may pay little to buy the training data to guarantee her profit since she bears certain costs (e.g., for training or operating), and therefore the hospitals offer poor-quality data, which causes low-performance models and in turn harms the profits of the company and the broker.

Existing works can't work well in real data markets with three selfish but interdependent parties that all want to realize profit optimization ( $\mathcal{P}_1$ ). A reasonable pricing mechanism is most important to balance the profits of all participants. Absolutely pricing data ( $\mathcal{P}_3$ ), however, is far from trivial. Due to the special properties of data (e.g., freely replicable and inherently combinatorial) summarized in [4], pricing mechanisms for physical goods cannot be directly applied to data. Therefore, the challenge is ( $C_1$ ): *How to design an incentive mechanism including absolute pricing rules for data markets that can balance the profits of all participants.*

As we mentioned before, we focus on the buyer-leading markets where transactions are initiated by buyers. Though efforts have been made to satisfy buyers' needs in [4][38][39], data markets where buyers have a leadership position ( $\mathcal{P}_2$ ) are still understudied. Therefore, the second challenge comes from the buyer's perspective ( $C_2$ ): *How to embody the advantages of buyers against other parties in the demand-leading data markets.*

To give priority to buyers ( $\mathcal{P}_2$ ), it's important to ensure them to gain the *best* data (with highest data quality). Typically there are huge numbers of data sellers with numerous data (e.g., there are large numbers of hospitals with abundant data), so it's critical to select the *best* data from a typical group of sellers, which we call *seller selection problem*. Many of existing works made the brokers responsible for seller selection [5][38][39]. In fact, there should be inner competition among sellers, which can make the winners as the selected sellers without the assistance of brokers or others. Hence, it's worth paying attention to the third challenge from the seller's side ( $C_3$ ): *How to model the inner competition among data sellers to select the best sellers for data transaction.*

**Contributions.** In this paper, we address the identified three challenges by proposing Stackelberg-Nash based Data Markets (*Share*), as shown in Figure 1. The detailed workflow is presented in Section 4.2.



**Figure 1: A data market model with Stackelberg-Nash game**

We adopt game theory to solve the problem of profit maximization for all three parties in data markets ( $C_1$ ). Game theory is a tool for designing multi-objective incentive mechanisms and can provide optimal strategies for involved participants. In our context, we introduce both Stackelberg game and Nash game into data markets for the interactions of the three entities and the inter-seller competition respectively. Specifically, the interactions among profit-interdependent market entities are modeled as a game process, in which each participant can achieve her profit-maximization goal by making her optimal strategy. Moreover, absolute prices of data are modeled as strategies of participants and are directly determined in this game process.

Considering the buyer priority, we formalize the interactions among three parties based on Stackelberg game which can trickily deal with the different status of buyers, brokers, and sellers. Specifically, we formulate the incentive market mechanism as a three-stage Stackelberg-Nash game by regarding buyers as the leaders, brokers as the sub-leaders, and sellers as the followers. Buyers can first announce what data products they demand for and determine the product price based on the profit-maximization goal. Brokers then try to buy data from sellers and each seller chooses what data quality to provide. Buyers are thus endowed with a dominant position in two aspects: the priority of initiating transactions and the intensive influence on prices ( $C_2$ ).

For the problem of seller selection, we elaborate the inner competition among sellers as a Nash game due to its advantage in modeling sellers' equal positions ( $C_3$ ), which is the first work that applies Nash game for the seller selection (data selection) problem. In our Nash model, the buyer needs data which can be from multiple sellers, and the quantity of data that each seller can sell is determined through the competition on data quality (data fidelity) which is correlated to privacy preserving level manipulated by sellers through local differential privacy mechanism. Nash equilibrium among sellers is preferred, which can guarantee that no seller can change her strategy individually and the seller selection result is stable. However, it's challenging to find the equilibrium point especially when the number of sellers is very large since the analytic solutions determining how sellers make strategy are desired which may be complicated and hard to derive. In *Share*, we apply the direct derivation and propose a mean-field based approximation to derive Nash equilibrium.

Our goal in this paper is not to cover all critical issues in buyer-leading data markets, but rather to propose a reasonable and feasible mechanism to satisfy several essential desiderata. Overall, the major contributions are summarized as follows.

- We present *Share*, a data market framework based on a three-stage Stackelberg-Nash game, which not only satisfies the properties of buyer-leading, for-all profit maximization, and absolute pricing, but also considers the inner competition among sellers.
- We first apply Nash game for the seller selection problem, which formulates sellers' inner competition and incorporates seller selection into the game process among three parties.
- We define Stackelberg-Nash Equilibrium in data markets and derive it by backward induction. To solve inner Nash game, we apply the direct derivation as well as design a novel mean-field method for complex cases, for which error analysis is conducted.

- We conduct experiments on real datasets to verify the performance of *Share*.

**Organization.** The rest of the paper is organized as follows. Section 2 presents the related work. Section 3 provides the preliminaries. We provide a data market model with Stackelberg-Nash game and design the strategy mechanism for market parties in Section 4. Approaches to deriving the market equilibrium and the detailed algorithm for market dynamics are presented in Section 5. We report the experimental results and findings in Section 6. Section 7 draws a conclusion and discusses future work.

## 2 RELATED WORK

In this section, we discuss related work on the data market and game theory.

### 2.1 Data Market

Data markets trade data either in direct or indirect forms (derived data products). Twitter [1], Bloomberg [2], SafeGraph [3], Dawex [15], Iota [28] and et al. implemented data markets where buyers directly purchase data. Koutiris et al. [32][33] proposed query-based data markets which allow buyers to obtain information through querying the database and pay for the query. Recently, model-based data markets [4][10][30][38][39] have been proposed. Agarwal et al. [4] applied algorithmic game theory for two-sided data-driven machine learning model markets. Dealer [38][39], an end-to-end model marketplace with differential privacy, more comprehensively considered the needs of sellers and buyers. It also regulated the broker's role as model pricing and model training.

In Dealer, however, it is assumed that the broker to be neutral without her own profit consideration and determines model prices only for single goal optimization, i.e., revenue maximization for sellers. An et al. [5], who studied transactions for crowdsensing data, were devoted to the multiple goal optimization. Nevertheless, they oversimplified the markets in terms of transaction objects, market participants, and profit function formulations. In this paper, we combine multiple game mechanisms to model for-all profit-maximization data markets with unrestricted data and data products, privacy consideration for data sellers, as well as exquisitely designed profit functions based on economics theories.

As for data pricing, there are several surveys [12][46][47] claiming fundamental principles and reviewing the development and evolution of pricing models. Ghorbani and Zou [25] introduced Shapley value to data valuation. Federated learning was combined with data pricing was developed by Wang et al. [57] while reinforcement learning was adopted by Yoon et al. [61] to price data. In terms of absolute pricing mechanism, Dealer [38][39] provided absolute prices for data models, which, however, highly rely on the survey results and can't be adjusted dynamically. Agarwal et al. [4] applied Myerson's payment function rule to determine absolute model payment but allocated relative compensations to sellers in proportion to their contributions. Therefore, we propose an absolute pricing mechanism to bridge the gap.

Many works looked at the seller selection problem. In Dealer [38][39], brokers chose datasets to achieve Shapley coverage maximum of the trained model. In [5], combinatorial multi-armed bandit

mechanism is used for seller selection by brokers. However, the selection results directly affect the profits of sellers, and therefore the seller selection problem is closely correlated to the profit maximization problem for sellers and should not be considered separately. In fact, the selection problem can be seen as the spontaneous process of the inner competition among sellers, not conducted by the broker or others. In our work, seller selection problem is formalized as the inner Nash game among sellers, which is a part of the incentive mechanism for profit optimization of all participants.

### 2.2 Game Theory

Game theory provides a decision-making and analysis tool for individual behaviors with conflicting objectives. Early game theory originated from the study of competition in a duopoly in microeconomics. Stackelberg game was originally proposed by von Stackelberg [54] to model the asymmetric competition among oligopoly firms. Morgenstern et al. [41] first defined the basic concepts of game, marking the preliminary formation of game theory. Nash accurately described Nash equilibrium [44] and proved its existence in  $N$ -player finite non-cooperative game with mixed strategies [43]. Later, Selten et al. [49] considered the dynamic game and proposed Subgame Perfect Nash Equilibrium. Harsanyi et al. [26] introduced incomplete information into game theory and proposed Bayesian Nash equilibrium. In terms of cooperative game, Lloyd Shapely [48] laid the foundation [23] and provided the formula of Shapley Value [50] which regulated how to allocate revenue among collaborators.

Game theory has been widely used in various situations. Many researchers used Nash game as a powerful tool to formulate and solve problems where there is competition [42][29]. Stackelberg game was first used to formulate the determining process for oligopoly firms producing homogeneous products [54], and has been further applied to many other practical situations with hierarchical organizations [51][56][62]. Besides the original Stackelberg game composed of one leader and one follower, many variants [6][37][53] were proposed and investigated. For example, Bansal et al. [6] used a two-stage Stackelberg game to determine prices for Unmanned Aerial Vehicles. Some studies [37][53] combined Stackelberg game and Nash game together to deal with the problems where both hierarchy and simultaneity exist, but their issues including participant roles, major actions, and optimizing goals in traditional scenarios are quite different from those in data markets.

Since Nash proposed his theory, many researchers have sought algorithms for finding Nash equilibrium. Conitzer et al. [13] showed complexity results of deriving Nash equilibrium and Daskalakis et al. [14] further studied the complexity of computing a mixed Nash equilibrium. In terms of solving Stackelberg game, backward induction approach, an iterative technique to derive dynamic game equilibrium, is often used [5][55][60]. In fact, establishing Stackelberg equilibrium can be formulated as a bilevel optimization problem [52], which has been studied extensively. Some studies also combined other techniques into the equilibrium solving problem [6][36][53].

In this paper, we design the buyer-leading data markets based on Stackelberg game because Stackelberg game concerns interactions of participants with asymmetric status and can thus realize our

desired buyer-leading property while remaining the profit maximization for all parties. Moreover, we first adopt Nash game for the seller selection problem since Nash game models the competition among equals with conflicting profits and can be used for the inner competition among data sellers, which can select sellers based on their optimal strategies.

### 3 PRELIMINARIES

In this section, we describe local differential privacy, Shapley value and game theory used in *Share*. For reference, Table 1 summarizes the frequently used notations.

**Table 1: The summary of frequently used notations.**

	Notation	Definition
Buyer $\mathcal{B}$	$N$	data quantity for production
	$p^M$	unit price of data product
	$v$	product performance
	$\theta_1, \theta_2$	parameters of concern on each attribute
	$\rho_1, \rho_2$	parameters of sensitivity to each attribute
	$U_1(\cdot), U_2(\cdot)$	utility function of each attribute
Broker $\mathcal{A}$	$U(\cdot)$	utility function of the product
	$\Phi(\cdot)$	profit function of the buyer
	$p^D$	unit price of data
	$\sigma_k, k = 0, 1, \dots, 5$	parameters related to cost
	$C(\cdot)$	cost function of production
Seller $S_i$	$\Omega(\cdot)$	profit function of the broker
	$i$	index of seller
	$m$	total number of sellers
	$\tau_i$	data fidelity
	$\epsilon_i$	parameter in local differential privacy
	$\chi_i$	sold data quantity
	$\lambda_i$	parameter of privacy sensitivity
	$L_i(\cdot)$	privacy loss function
	$\Psi_i(\cdot)$	profit function
Data	$D_i$	seller $S_i$ 's raw dataset
	$D_i^t$	seller $S_i$ 's provided dataset
	$D^t$	whole dataset for manufacturing
	$q_i^D$	dataset quality provided by seller $S_i$
	$q^D$	total quality of dataset for manufacturing
	$q^M$	data product quality
	$\omega_i$	weight of seller $S_i$ 's dataset

#### 3.1 Local Differential Privacy

Differential Privacy (DP) [18][19][20] is a framework for privacy protection by providing perturbation against the discovery of presence or absence in a dataset. As for data markets, DP is adopted in Dealer [38][39] to protect sellers' privacy, where the broker utilizes a perturbation algorithm on collected data and trains a set of models with differential privacy guarantees. Since each seller has distinctive preference to privacy preservation, the protection for privacy should be conducted by each seller locally with personal privacy budget. Local Differential Privacy (LDP) [17] is an extended version of DP, which is suited to the above setting.

**Definition 3.1 (Local Differential Privacy).** A randomized algorithm  $\mathcal{A} : \mathcal{Y} \rightarrow \mathcal{Z}$  satisfies  $\epsilon$ -local differential privacy if and only if for any pairs of input tuples  $y, y' \in \mathcal{Y}$ , and for any  $z \in \mathcal{Z}$ , it always holds

$$\mathbb{P}[\mathcal{A}(y) = z] \leq e^\epsilon \cdot \mathbb{P}[\mathcal{A}(y') = z] \quad (1)$$

where  $\epsilon \geq 0$  denotes the privacy budget. The notation  $\mathbb{P}[\cdot]$  denotes probability.

Widely used (local) differential privacy mechanisms include Laplace mechanism [18], Index mechanism [40], and Gaussian mechanism [20]. The smaller  $\epsilon$  is, the less the privacy loss is, and the worse the dataset quality is. In our context, each seller adopts a privacy scheme satisfying local differential privacy to ensure her own privacy protection and form her for-sale dataset with a distinctive quality.

#### 3.2 Shapley Value

Shapley value [50] is an approach to fairly evaluate data importance. Proposed by Shapley [48], Shapley value satisfies the four fundamental requirements of fairness in markets, i.e., balance, symmetry, zero element, and additivity. In our data markets, Shapley value can be used to measure the contribution of each seller's provided data to the data product by measuring the marginal utility improvement, e.g., the accuracy increase for a classification model or the explained variance score raise for a regression model.

**Definition 3.2 (Shapley Value).** Consider a set of  $m$  data sellers such that each seller  $S_i$  owns a dataset  $D_i$  ( $i = 1, 2, \dots, m$ ). A coalition  $\mathbb{D}$  is a subset of  $\{D_1, D_2, \dots, D_m\}$ . Denote by  $U(\mathbb{D})$  a utility function that represents the performance of a data product manufactured by coalition  $\mathbb{D}$  towards a task, e.g., accuracy of a machine learning model. Shapley value of seller  $S_i$  is defined as follows.

$$SV_i = \frac{1}{m} \sum_{\mathbb{D} \subseteq \{D_1, D_2, \dots, D_m\} \setminus D_i} \frac{U(\mathbb{D} \cup \{D_i\}) - U(\mathbb{D})}{\binom{m-1}{|\mathbb{D}|}}. \quad (2)$$

In this paper, Shapley value is adopted to measure the weight of the dataset provided by each seller. Such weights infer the past performance of each seller's dataset based on its contribution to the product performance improvement.

#### 3.3 Game theory

Game theory is the study of mathematical models of strategic interactions among rational agents. A game must specify the players, the strategies, and the payoffs for each outcome. An equilibrium to the game is a stable state in which no player would change her strategy. Two non-cooperative games, the games where players cannot form alliances, are applied in this work, Nash game and Stackelberg game.

**Definition 3.3 (Nash game, Nash equilibrium).** Nash game is a game where two or more players take strategy simultaneously to maximize their own expected payoff. Nash equilibrium is the stable state where no player has anything to gain by changing only one's own strategy. Formally, let  $S_i$  be the set of all possible strategies for player  $i$ , where  $i = 1, 2, \dots, n$ . Let  $s^* = (s_i^*, s_{-i}^*)$  be a strategy profile, a set consisting of one strategy for each player, where  $s_{-i}^* = (s_1^*, \dots, s_{i-1}^*, s_{i+1}^*, \dots, s_n^*)$  denotes the  $n-1$  strategies of all the players except  $i$ . Let  $\pi_i(s_i, s_{-i})$  be player  $i$ 's payoff as a function of the strategies. The strategy profile  $s^*$  is a Nash equilibrium if

$$\pi_i(s_i^*, s_{-i}^*) \geq \pi_i(s_i, s_{-i}^*), \forall s_i \in S_i, i = 1, 2, \dots, n. \quad (3)$$

**Definition 3.4 (Stackelberg game, Stackelberg equilibrium).** Stackelberg game is a game where one player (the *leader*  $L$ ) moves first, and the other player (the *follower*  $F$ ) can observe the *leader*'s behavior and move sequentially. Stackelberg equilibrium is a refinement

of Nash equilibrium used in dynamic games. Similarly, the strategy profile  $(s_L^*, s_F^*)$  is a Stackelberg equilibrium if

$$\begin{aligned}\pi_L(s_L^*, s_F^*) &\geq \pi_L(s_L, s_F^*), \\ \pi_F(s_L^*, s_F^*) &\geq \pi_F(s_L^*, s_F).\end{aligned}\quad (4)$$

#### 4 MARKET MODEL: PARTICIPANTS, MECHANISM, AND EQUILIBRIUM

In this section, we first describe the market model from the perspectives of buyers, brokers, and sellers respectively in Section 4.1. Then, we formulate the market mechanism as a three-stage Stackelberg-Nash game in Section 4.2 and define the market equilibrium, Stackelberg-Nash Equilibrium in Section 4.3.

##### 4.1 Market Participants

Considering the data markets composed of three types of participants, i.e., buyers, brokers, and sellers, we define the role of each party as follows.

- **Buyer.** Buyer  $\mathcal{B}$  wants to fulfill her data-driven task through data trading. Buyer  $\mathcal{B}$  needs to purchase a data product from broker  $\mathcal{A}$  with her demands claimed, i.e., the required product performance indicated as  $\nu$  and the size of the dataset for manufacturing denoted by  $N$ .
- **Broker.** Broker (Arbiter)  $\mathcal{A}$  wants to make profits by bridging the transactions between buyers and sellers. Broker  $\mathcal{A}$  needs to buy  $N$  data pieces from  $m$  sellers and make the product from the data to satisfy the product performance demand  $\nu$  which takes cost because manufacturing data products consumes resources (e.g., computing resources). Then, broker  $\mathcal{A}$  sells the product to buyer  $\mathcal{B}$  in exchange of rewards.
- **Seller.** Each seller  $S_i, i = 1, 2, \dots, m$  owns a dataset  $D_i$  and wants to sell it for compensation. Seller  $S_i$  needs to preprocess her dataset for privacy protection and sell the protected  $\chi_i$  data pieces to broker  $\mathcal{A}$ . Note that  $\sum_{i=1}^m \chi_i = N$ , the total number of data pieces that buyer  $\mathcal{B}$  requires. Seller  $S_i$  will get the compensation while suffering from the privacy loss which comes from the potential privacy leaks from the data she offers to broker  $\mathcal{A}$ .

Since there are various types of markets with distinctive characteristics which cannot be all covered, we choose one type of markets with the following reasonable assumptions.

- The market is buyer-leading, and one trade starts up when a buyer raises the demand. Buyers orientate the market in turn (coming one at a time) as in work [4] so that in each round of transactions only one buyer  $\mathcal{B}$  is considered.
- There exists one broker  $\mathcal{A}$  in our market, e.g., the large-scale data trading center in the real world, which infers that the competition among brokers for the same data product is beyond our consideration.
- There are a large number of sellers  $\{S_i | i = 1, 2, \dots, m\}$ , and each seller  $S_i$  has a dataset  $D_i$  to participate in the trading which is big enough, i.e., for any required number  $\chi_i \in \mathbb{N}^+$  of data pieces,  $|D_i| \geq \chi_i$ .
- The market is highly transparent, which means each participant can see the behaviors of other participants, and

cooperation is excluded, inferring that no two entities would collude to get a better outcome.

Each of the participants has the *revenue* (the received compensation or the gained utility) and the *cost* (the payment, the manufacturing cost, or the privacy loss). We assume that all the participants in the market are profit-driven and want to maximize their own profit, i.e., the difference between the *revenue* and *cost*, and the profit functions are known public. The detailed definitions of the profit functions of buyers, brokers, and sellers are given, respectively.

##### 4.1.1 Profit Function of Buyer.

When buyer  $\mathcal{B}$  comes to the market and asks for a data product with the required performance, she cares about her *revenue*, the utility she can get from the product and her *cost*, the payment she should give to the broker.

*Revenue.* The revenue of buyer  $\mathcal{B}$  is the product utility which is quantified as a utility function  $U(\cdot)$ . As for data products, what buyer  $\mathcal{B}$  concerns includes both the product performance and the quality of the dataset used in the manufacture. Apparently, product performance matters to buyer  $\mathcal{B}$ . However, it only infers how the product performs under a certain testing environment. For example, the performance of the data model is confined to a specific validation dataset, thus causing an out-of-sample error on a different dataset. Dataset quality measures how *good* the raw materials are, adding a new dimension to the judgment of product utility. Combining product performance and dataset quality together to measure product utility can alleviate the above problem and make the quantification of product utility more comprehensive. Dataset quality is defined as  $q_i^D = g(\chi_i, \tau_i)$  which is positively correlated with both data fidelity  $\tau_i$  and data quantity  $\chi_i$  provided by seller  $S_i$ . Specifically, data fidelity  $\tau_i$  is closely correlated with the privacy protection by local differential privacy mechanism, and data quantity  $\chi_i$  can be determined based on sellers' inner competition. The detailed formulations for  $\tau_i$  and  $\chi_i$  will be elaborated later.

According to the utility theory in economics [59], the utility of a data product is defined as the weighted sum of the utility of the dataset quality and the product performance,  $U_1(q^D)$  and  $U_2(\nu)$ . We further formulate them as the logarithmic functions as below which follow the principle of diminishing marginal utility in economics.

$$\begin{aligned}U_1(q^D) &= \ln(1 + \rho_1 q^D), \\ U_2(\nu) &= \ln(1 + \rho_2 \nu).\end{aligned}\quad (5)$$

The data product we consider composes of two attributes, the quality  $q^D$  of dataset used in the production process, and the performance  $\nu$  of the product.  $q^D = \sum_{i=1}^m q_i^D = \sum_{i=1}^m g(\chi_i, \tau_i)$  refers to the total quality of the dataset for manufacturing composed of the datasets from sellers.  $\nu > 0$  is positively correlated to  $\mathcal{B}$ 's required product performance measured by one or several specific indicators. For data models,  $\nu$  can be the explained variance, accuracy, or other indicators measuring model performance.  $U_1(\cdot)$  and  $U_2(\cdot)$  define the utility of dataset quality and product performance.  $\rho_1 > 0$  and  $\rho_2 > 0$  refer to buyer  $\mathcal{B}$ 's sensitivity to these two attributes respectively. More sensitive, more utility added when the attribute gets better. For example, if higher dataset quality can bring buyer  $\mathcal{B}$  much more utility, its  $\rho_1$  would be big, meaning that buyer  $\mathcal{B}$  is highly sensitive to the quality of production materials. Moreover, it

can be intuitively seen that better dataset quality or better product performance brings about better product utility.

Then, we can make the following definition for the total utility of a data product

$$U(\chi, \tau, v) = \theta_1 U_1(q^D) + \theta_2 U_2(v), \quad (6)$$

where  $\chi = (\chi_1, \chi_2, \dots, \chi_m)$  and  $\tau = (\tau_1, \tau_2, \dots, \tau_m)$ .  $\theta_1$  and  $\theta_2$  satisfy  $\theta_1, \theta_2 \in (0, 1)$ ,  $\theta_1 + \theta_2 = 1$ , which measure buyer  $\mathcal{B}$ 's concern on the dataset quality and the product performance, indicating the relative significance between these two attributes. In our example, if dataset quality plays a greater role in the decision-making of the drug company, the company may set  $\theta_1 = 0.7$  and  $\theta_2 = 0.3$ .

*Cost.* The cost of buyer  $\mathcal{B}$  is the payment to broker  $\mathcal{A}$ .  $q^M = h(q^D, v)$  is defined as the quality of the data product, which infers that the quality of the data product depends on both dataset quality  $q^D$  in production and product performance  $v$ . Also,  $p^M$  is defined as the unit price of  $q^M$  (or, simplified as the price of data product). Therefore, the payment for the product can be formulated as the product of the price times the product quality, i.e.,  $p^M q^M$ , which corresponds to our common sense that the payment for goods is equal to the unit price multiplied by the quantity.

*Profit.* The profit  $\Phi(\cdot)$  of buyer  $\mathcal{B}$  is the difference between the utility of the product and the payment to the broker as follows.

$$\Phi(p^M, \tau) = U(\chi, \tau, v) - p^M q^M. \quad (7)$$

#### 4.1.2 Profit Function of Broker.

When broker  $\mathcal{A}$  receives the data product requirements from buyer  $\mathcal{B}$ , she cares about her *revenue*, the payment from buyer  $\mathcal{B}$  and her *cost* which consists of the compensations to sellers to buy the data and the manufacturing cost in the process of producing the data product.

*Revenue.* The revenue of broker  $\mathcal{A}$  is the payment from buyer  $\mathcal{B}$ , i.e.,  $p^M q^M$  (the cost of buyer  $\mathcal{B}$ ).

*Cost.* The cost of broker  $\mathcal{A}$  is the sum of the compensations to sellers and the manufacturing cost. For the former, broker  $\mathcal{A}$  needs to manipulate the compensations to sellers according to their provided dataset quality.  $p^D q^D$  is defined as the compensations to sellers where  $q^D$  means the total quality of the manufacturing dataset and  $p^D$  is the unit price of  $q^D$  (or, simplified as the price of data). The value of compensations equals the product of the unit data price times the total dataset quality, similar with the above discussion of payment for the data product.

The second part of broker  $\mathcal{A}$ 's cost comes from the product production. Broker  $\mathcal{A}$  needs to consume some resources to make the product, e.g., computation resources for training models, which is measured by cost function  $C(\cdot)$  related to data size  $N$  and product performance  $v$ . It's obvious that how many data put into production and what performance level needed to achieve significantly affect the manufacturing cost. According to the work [11], a kind of widely used transcendental logarithmic cost function form is adopted due to its adaptability to varied economies of scale and manufacturing strategy (e.g., how to allocate computing resources). The cost function is defined as follows which can be manipulated through various parameters by broker  $\mathcal{A}$  according to the practical manufacturing situation.

$$C(N, v) = \exp \left( \sigma_0 + \sigma_1 \ln(N) + \sigma_2 \ln(v) + \frac{1}{2} \sigma_3 \ln^2(N) + \frac{1}{2} \sigma_4 \ln^2(v) + \sigma_5 \ln(N) \cdot \ln(v) \right). \quad (8)$$

$\sigma_k, k \in \{0, 1, 2, 3, 4, 5\}$  are the parameters of the translog cost function which can be fitted by broker  $\mathcal{A}$  based on the actual manufacturing procedure.

*Profit.* The profit  $\Omega(\cdot)$  of broker  $\mathcal{A}$  is defined as the received payment from buyer  $\mathcal{B}$  minus the manufacturing cost and the compensations to sellers as follows.

$$\Omega(p^M, p^D, \tau) = p^M q^M - C(N, v) - p^D q^D. \quad (9)$$

#### 4.1.3 Profit Function of Seller.

When seller  $\mathcal{S}_i$  gets the purchase request for data from broker  $\mathcal{A}$ , she cares about her *revenue*, the compensation from broker  $\mathcal{A}$  and her *cost* coming from her privacy loss.

*Revenue.* The revenue of seller  $\mathcal{S}_i$  is the compensation from broker  $\mathcal{A}$ , i.e.,  $p^D q_i^D$  (the first part of the cost of broker  $\mathcal{A}$ ).

*Cost.* The cost of seller  $\mathcal{S}_i$  is the privacy loss she suffers based on data fidelity  $\tau_i$  she provides. As we discussed before, the data fidelity is closely correlated to the data protection process.  $\tau_i = f(\epsilon_i)$  is defined where  $\epsilon_i$  is the parameter in local differential privacy as in Section 3.1. Each seller would adopt LDP before she sells the data to broker  $\mathcal{A}$ . In fact, what broker  $\mathcal{A}$  receives is the *processed data* protected by each seller to some extent. More  $\epsilon_i$ , less noise added, thus causing better fidelity  $\tau_i$ . There are many alternative function forms for  $f(\cdot)$ , as long as following the diminishing trend of marginal effect. In fact,  $f(\cdot)$  should satisfy Inada Conditions [27] in economics which can stipulate the marginal change trend and encourage the market equilibrium. We conclude the following characteristics that  $f(\cdot)$  should satisfy.

1. The data piece has poor fidelity  $\tau_i = 0$  when  $\epsilon_i = 0$  which means random noise has been added to the data, causing severe data distortion.
2. Bigger  $\epsilon_i$ , bigger  $\tau_i$  due to less noise to data.
3.  $\tau_i$  increases slower as  $\epsilon_i$  becomes bigger because when  $\epsilon_i$  is big enough, it can hardly make a significant difference to further increasing the data fidelity because the fidelity is high enough as very little noise is added. On the other hand, when  $\epsilon_i$  is very small, reducing noise, i.e., increasing  $\epsilon_i$ , can remarkably increase the data fidelity. Besides,  $\tau_i$  cannot increase perpetually and should be upper bounded.

Based on Inada Conditions, we choose an inverse trigonometric function form as  $f(\cdot)$  and give the following definition of  $\tau_i$ .

$$\tau_i = \frac{2}{\pi} \arccos(\epsilon_i + 1), \quad \epsilon_i \in [0, \infty), \quad (10)$$

which infers that  $\tau_i \in [0, 1)$ . Additionally,  $\tau_i = 1$  when no noise is added. Therefore,  $\tau_i \in [0, 1]$ .

Bigger  $\tau_i$ , the better fidelity of data, but more privacy loss for  $\mathcal{S}_i$ . We define such loss of  $\mathcal{S}_i$  through the function  $L_i(\cdot)$  which is positively related to  $\tau_i$ . Various function forms can be applied to formulate sellers' privacy loss. For the sake of simplicity in the solving process, a widely used quadratic function with variable  $\tau_i$  is

chosen, which reflects the trend of loss with  $\tau_i$ . It's reasonable that the loss function is not only increasing but also increasing faster for bigger  $\tau_i$ , which corresponds with the principle of increasing marginal cost in economics. Moreover, the loss seller  $S_i$  suffers is positively related to the data quantity  $\chi_i$  she provides. Thus, The privacy loss function is defined as follows.

$$L_i(\tau_i) = \lambda_i(\chi_i \tau_i)^2, \quad (11)$$

where  $\lambda_i > 0$  is  $S_i$ 's privacy sensitivity. In our example, the privacy loss of the hospital infers that the hospital needs to make amends for leaking the patients' personal information, and  $\lambda_i$  is associated with their beforehand agreement on the usage of personal medical data. In fact,  $\lambda_i$  tunes  $S_i$ 's price elasticity of privacy loss.

**Profit.** The profit of seller  $S_i$  is the difference between the compensation from broker  $\mathcal{A}$  and the privacy loss as follows.

$$\Psi_i(p^D, \tau_i) = p^D q_i^D - L_i(\tau_i). \quad (12)$$

## 4.2 Market Mechanism

As we mentioned before, it's assumed that participants of three types come to the market and take strategies in order. We present the market workflow first. Then we specify the strategies of buyer  $\mathcal{B}$ , broker  $\mathcal{A}$ , and each seller  $S_i$ , respectively. Based on the strategies, the market mechanism is proposed.

**Market Workflow.** The market workflow is consistent with Fig. 1. ① Buyer  $\mathcal{B}$  comes and puts forward the demands for the product including the required product performance and the corresponding indicator(s) as well as the size of dataset for manufacturing. ② Buyer  $\mathcal{B}$  determines the unit product price to buy the data product from broker  $\mathcal{A}$ . ③ Broker  $\mathcal{A}$ , acting as the bridge for the trade between buyer  $\mathcal{B}$  and  $m$  sellers, determines the unit data price to buy the data from sellers. ④ Each seller chooses what data, strictly speaking, what data fidelity to sell, and conducts corresponding privacy protection locally. ⑤ Sellers sell the protected datasets to broker  $\mathcal{A}$  in exchange of the compensations. ⑥ Using the dataset bought from sellers, broker  $\mathcal{A}$  manufactures the product. ⑦ Broker  $\mathcal{A}$  sells the product to buyer  $\mathcal{B}$ . After  $\mathcal{B}$  receives the product and gives payment to  $\mathcal{A}$ , the trade is finished.

**Buyer's Strategy.** Buyer  $\mathcal{B}$  first makes her strategy. The strategy of buyer  $\mathcal{B}$  is the product price  $p^M$ . Buyer  $\mathcal{B}$  determines  $p^M$  to maximize her profit by considering the utility of the product and predicting the responses of the broker and sellers to the  $p^M$ .

**Broker's Strategy.** Broker  $\mathcal{A}$  takes her strategy secondly. The strategy of broker  $\mathcal{A}$  is data price  $p^D$ . Broker  $\mathcal{A}$  determines  $p^D$  to maximize her profit by considering the given  $p^M$  from buyer  $\mathcal{B}$ , predicting what data fidelity sellers would provide according to such  $p^D$  as well as thinking about the corresponding manufacturing cost.

**Seller's Strategy.** Sellers make their strategies thirdly. The strategy of each seller  $S_i$  is data fidelity  $\tau_i$ . Seller  $S_i$  determines  $\tau_i$  to maximize her own profit by balancing the revenue of selling data and the cost of the privacy loss from the sold data.

Meanwhile, inner competition among  $m$  sellers should be considered. To be specific, given the unit data price  $p^D$ , data with bigger  $\tau_i$  provided, more quantity of  $S_i$ 's data,  $\chi_i$ , would be sold, while if

other sellers provide data with better fidelity, less data quantity of  $S_i$  could be chosen. Therefore,  $\chi_i$  can be calculated according to all sellers'  $\tau$  as below.

$$\chi_i = N \frac{\omega_i \tau_i}{\sum_{j=1}^m \omega_j \tau_j}, \quad (13)$$

where  $\omega_1, \omega_2, \dots, \omega_m$  refer to the weights of sellers' data, which are maintained by the broker. Such weights reflect the historical performance of each seller's data in past deals and can therefore mirror previous buyers' influence on the current buyer to some extent. The broker would update these weights after each round of transaction. For example, new weights can be generated based on both old weights and sellers' contributions to the data product measured by Shapley value [50] in the last transaction.

We define such inner competition in sellers as a Nash game, that is, each seller  $S_i$  determines their strategy  $\tau_i$  simultaneously to maximize her own profit which is meantime affected by other sellers' strategies. Nash equilibrium would be achieved where no seller can increase her profit by unilaterally changing her strategy with all other sellers' strategies fixed.

**Three-stage Stackelberg-Nash Game.** Strategies of buyer  $\mathcal{B}$ , broker  $\mathcal{A}$ , and sellers  $S_i$  ( $i = 1, 2, \dots, m$ ) constitute the strategy profile  $\langle p^M, p^D, \tau \rangle$  of data markets. Such a profile determines market trading rules including what data (data fidelity) to sell ( $\tau$ ), selling at what price for both data ( $p^D$ ) and data product ( $p^M$ ) as well as how to select sellers (the calculated  $\chi$  based on  $\tau$ ). The market mechanism is formulated as a three-stage Stackelberg-Nash game, where buyer  $\mathcal{B}$  is the leader, broker  $\mathcal{A}$  is the sub-leader, and  $m$  sellers act as the followers. Each of them tries to maximize her own profit by determining her optimal strategy variable. The three-stage Stackelberg-Nash game is defined as follows.

*Definition 4.1 (three-stage Stackelberg-Nash game).* The game consists of three stages for buyer  $\mathcal{B}$ , broker  $\mathcal{A}$ , and sellers  $S_1, S_2, \dots, S_m$ , respectively.

*Stage 1 Buyer  $\mathcal{B}$ :*  $p^{M*} = \arg \max_{p^M} \Phi(p^M, \tau)$ .

*Stage 2 Broker  $\mathcal{A}$ :*  $p^{D*} = \arg \max_{p^D} \Omega(p^M, p^D, \tau)$ .

*Stage 3 Seller  $S_i$ :*  $\tau_i^* = \arg \max_{\tau_i} \Psi_i(p^D, \tau)$ ,  $i = 1, 2, \dots, m$ .

The above three-stage Stackelberg-Nash game involves both hierarchy and simultaneity. Hierarchy infers that a certain participant, buyer  $\mathcal{B}$  in our context, has some sort of advantage over others that enables her to act first, broker  $\mathcal{A}$  takes her strategy secondly, and sellers make their strategies last, while simultaneity indicates the equal positions of  $m$  sellers, who take strategy simultaneously in their inner Nash game.

## 4.3 Market Equilibrium

In the above game, our objective is to find an optimal strategy profile  $\langle p^{M*}, p^{D*}, \tau^* \rangle$ , by which each participant can maximize her own profit. Meanwhile, the optimal solution must satisfy some equilibrium so that no one is willing to adopt other strategies, which indicates the market stability and sustainability, making our design reasonable. We novelly define a Stackelberg-Nash Equilibrium (SNE) in data markets as follows.



**Definition 4.2 (Stackelberg-Nash Equilibrium).** An optimal strategy profile  $\langle p^{M*}, p^{D*}, \tau^* \rangle$  constitutes a Stackelberg-Nash Equilibrium (SNE) if and only if the following set of inequalities is satisfied.

$$\Phi(p^{M*}, \tau^*) \geq \Phi(p^M, \tau^*), \quad (14)$$

$$\Omega(p^{M*}, p^{D*}, \tau^*) \geq \Omega(p^{M*}, p^D, \tau^*), \quad (15)$$

$$\Psi_i(p^{D*}, \tau^*) \geq \Psi_i(p^{D*}, \tau_{-i}^*, \tau_i), i = 1, 2, \dots, m, \quad (16)$$

where  $\tau_{-i}$  means other sellers' strategies except  $\mathcal{S}_i$ 's, i.e.,  $\tau_j, j \neq i, j = 1, 2, \dots, m$ .

SNE indicates that each participant takes her optimal strategy which maximizes her own profit in a buyer-leading sequence. No one can add her own profit by unilaterally changing her strategy with all other participants' strategies fixed.

## 5 MARKET CONSTRUCTION: EQUILIBRIUM SOLVING APPROACH AND DATA TRADING ALGORITHM

In this section, we first derive the market equilibrium by backward induction in Section 5.1. We solve Stage 3, i.e., sellers' Nash equilibrium by two approaches, direct derivation and mean-field approximation for complex cases. Error analysis for the mean-field approximation and equilibrium analysis for SNE are also presented. Then, we present an algorithm to describe the market dynamics in Section 5.2.

### 5.1 Equilibrium Solving: Backward Induction Approach

To determine the optimal strategy profile  $\langle p^{M*}, p^{D*}, \tau^* \rangle$ , we adopt the backward induction approach [7]. We first investigate Stage 3 to solve Nash equilibrium among sellers and derive the expression of each seller's optimal strategy  $\tau_i^*, i = 1, 2, \dots, m$  (Eq. 20) for any given unit data price  $p^D$  in Section 5.1.1. To find the solution to Nash equilibrium, we explore two methods, i.e., the direct derivation and an approximate method using the mean-field state which can deal with complicated cases. Next, we consider Stage 2 to determine the expression of the optimal strategy  $p^{D*}$  (Eq. 25) of broker  $\mathcal{A}$  for any given unit product price  $p^M$  in Section 5.1.2. In this process, the expression of  $\tau_i^*, i = 1, 2, \dots, m$  solved from Nash game can be used as sellers' optimal reactions to  $p^D$ . Then, we back to Stage 1 to find the value (rather than the expression) of buyer  $\mathcal{B}$ 's optimal strategy  $p^{M*}$  (Eq. 27) based on the optimal reactions of the broker as well as sellers in Section 5.1.3. After that, we can get the value of the optimal strategy  $p^{D*}$  by substituting  $p^{M*}$  into the result (Eq. 25) in Stage 2. Finally, we can compute the value of each seller's optimal strategy  $\tau_i^*$  by substituting  $p^{D*}$  into the result (Eq. 20) in Stage 3. Till now, the whole optimal incentive strategy profile  $\langle p^{M*}, p^{D*}, \tau^* \rangle$  has been determined. The detailed deduction is presented as follows.

#### 5.1.1 Expression of $\tau^*$ in Stage 3.

We present two approaches to derive the expression of  $\tau_i^*$  for sellers, i.e., the direct derivation and a mean-field based approximation method.

**Direct Derivation.** By substituting Eqs. 11,13 into Eq. 12 and instantiating  $q_i^D = g(\chi_i, \tau_i)$  as  $\chi_i \tau_i$  since  $q_i^D$  is positively correlated with  $\chi_i$  and  $\tau_i$ , we get each seller's profit

$$\begin{aligned} \Psi_i(p^D, \tau_i) &= p^D \chi_i \tau_i - \lambda_i (\chi_i \tau_i)^2 \\ &= p^D \cdot N \frac{\omega_i \tau_i^2}{\sum_{j=1}^m \omega_j \tau_j} - \lambda_i \left( N \frac{\omega_i \tau_i^2}{\sum_{j=1}^m \omega_j \tau_j} \right)^2, i = 1, 2, \dots, m. \end{aligned}$$

$\Psi_i$  is correlated to not only seller  $\mathcal{S}_i$ 's strategy  $\tau_i$  but also other sellers' strategies  $\tau_j, j \neq i$  due to the inner competition formulated as Nash game among sellers. As we discussed before, each seller aims to maximize her own profit. Therefore, we derive each of the first-order derivatives for  $m$  sellers' profit functions and let each of them equal to zero, thus getting  $m$  equations. The equation for seller  $\mathcal{S}_i$  is

$$p^D \frac{\partial \left( N \frac{\omega_i \tau_i^2}{\sum_{j=1}^m \omega_j \tau_j} \right)}{\partial \tau_i} - 2\lambda_i \cdot N \frac{\omega_i \tau_i^2}{\sum_{j=1}^m \omega_j \tau_j} \cdot \frac{\partial \left( N \frac{\omega_i \tau_i^2}{\sum_{j=1}^m \omega_j \tau_j} \right)}{\partial \tau_i} = 0. \quad (17)$$

If  $\frac{\partial \left( N \frac{\omega_i \tau_i^2}{\sum_{j=1}^m \omega_j \tau_j} \right)}{\partial \tau_i} = 0$ , it is an all-zero solution, which does not meet

our problem situation, so we can directly eliminate  $\frac{\partial \left( N \frac{\omega_i \tau_i^2}{\sum_{j=1}^m \omega_j \tau_j} \right)}{\partial \tau_i}$ , and then get

$$\begin{cases} p^D \sum_{i=1}^m \omega_i \tau_i - 2N\lambda_1 \omega_1 \tau_1^2 = 0 \\ p^D \sum_{i=1}^m \omega_i \tau_i - 2N\lambda_2 \omega_2 \tau_2^2 = 0 \\ \vdots \\ p^D \sum_{i=1}^m \omega_i \tau_i - 2N\lambda_m \omega_m \tau_m^2 = 0, \end{cases} \quad (18)$$

where each  $\mathcal{S}_i$ 's equation not only relates to her own strategy  $\tau_i$  but also contains other sellers' strategies, requiring us to solve  $m$  simultaneous equations together. Finding that

$$2N\lambda_1 \omega_1 \tau_1^2 = 2N\lambda_2 \omega_2 \tau_2^2 = \dots = 2N\lambda_m \omega_m \tau_m^2 = p^D \sum_{i=1}^m \omega_i \tau_i. \quad (19)$$

By adding all  $m$  equations in Eq. 18, we can get

$$mp^D \sum_{i=1}^m \omega_i \tau_i - 2N \sum_{i=1}^m \lambda_i \omega_i \tau_i^2 = 0,$$

and using  $\tau_1$  to indicate other  $\tau_i$  ( $i = 2, 3, \dots, m$ ) from Eq. 19, we get

$$mp^D \tau_1 \sum_{i=1}^m \sqrt{\frac{\lambda_1 \omega_1 \omega_i}{\lambda_i}} - 2N m \lambda_1 \omega_1 \tau_1^2 = 0.$$

Therefore,

$$\tau_1^* = \frac{p^D}{2N\sqrt{\omega_1 \lambda_1}} \sum_{i=1}^m \sqrt{\frac{\omega_i}{\lambda_i}},$$

and using Eq. 19 again, we can get all sellers' optimal strategies  $\tau_i^*$  as

$$\tau_i^* = \frac{p^D}{2N\sqrt{\omega_i \lambda_i}} \sum_{j=1}^m \sqrt{\frac{\omega_j}{\lambda_j}}, i = 1, 2, \dots, m. \quad (20)$$



Note that we justify that the second-order derivative  $\frac{\partial^2 \Psi_i(p^D, \tau_i)}{\partial \tau_i^2} < 0$ , so these solutions can maximize each seller's profit.

**Mean-field based Approximate Method.** It's theoretically feasible that the optimal  $\tau$  can be derived by directly using the derivation method for each seller's profit function and then solving  $m$  simultaneous equations as above. However, for complicated function forms (e.g., more complicated loss function rather than the used quadratic one), since the number of sellers  $m$  is quite large in practice, it may be hard to derive analytical expressions by solving so many simultaneous equations each with complex forms. Specifically, the  $m$  equations are highly coupled, i.e., each with all  $\tau_i, i = 1, 2, \dots, m$ , and eliminating the similar terms to simplify the equations as we did in Eq. 17 is not always feasible. Note that even computers cannot help in this solving process because we desire the analytical solutions rather than the numerical solutions. Therefore, we propose an approximate method which makes each equation with a single  $\tau_i$  and independent from others. We take another privacy loss function form as an example where the direct derivation is not practically feasible in order to illustrate the mean-field method. Specifically, we replace Eq. 11 with  $L_i(\tau_i) = \lambda_i \chi_i \tau_i^2$ .

The approximation is based on the mean-field theory [34]. The mean-field theory deals with situations that involve a great number of agents, i.e., sellers in our context. When there are great numbers of sellers in Nash game, it's reasonable that a single seller has a *tiny* (infinitesimal) influence on the equilibrium and is affected by other sellers through a mean-field state, which we formulate as  $\bar{\tau}$ , the weighted mean of all sellers' strategies.

$$\bar{\tau} = \frac{\sum_{i=1}^m \omega_i \tau_i}{m}. \quad (21)$$

The mean-field state  $\bar{\tau}$  infers the overall data fidelity provided by sellers at equilibrium and is not intensively affected by the data fidelity from one specific seller.

Using the new privacy loss function, the profit function of seller  $S_i$  in Eq. 12 is changed into

$$\Psi_i(p^D, \tau_i) = p^D(\chi_i \tau_i) - \lambda_i \chi_i \tau_i^2. \quad (22)$$

Using  $\bar{\tau}$ , now  $\chi_i$  can be simplified as  $N \frac{\omega_i \tau_i}{m \bar{\tau}}$ . Since  $\bar{\tau}$  is not strongly affected by specific  $\tau_i$ , we can easily derive the first-order derivative of each seller's profit function  $\Psi_i(p^D, \tau_i)$  with respect to  $\tau_i$

$$\begin{cases} p^D \cdot N \frac{\omega_1 \tau_1^2}{m \bar{\tau}} - \lambda_1 \cdot N \frac{\omega_1 \tau_1^3}{m \bar{\tau}} = 0 \\ p^D \cdot N \frac{\omega_2 \tau_2^2}{m \bar{\tau}} - \lambda_2 \cdot N \frac{\omega_2 \tau_2^3}{m \bar{\tau}} = 0 \\ \vdots \\ p^D \cdot N \frac{\omega_m \tau_m^2}{m \bar{\tau}} - \lambda_m \cdot N \frac{\omega_m \tau_m^3}{m \bar{\tau}} = 0. \end{cases}$$

We derive  $S_i$ 's optimal strategy

$$\tau_i^* = \frac{2p^D}{3\lambda_i}, i = 1, 2, \dots, m. \quad (23)$$

Note that we justify that the second-order derivative  $\frac{\partial^2 \Psi_i(p^D, \tau_i)}{\partial \tau_i^2} < 0$ , so these solutions can maximize each seller's profit.

**Error Analysis.** We use fixed  $\bar{\tau}$  to replace  $\frac{\sum_{i=1}^m \omega_i \tau_i}{m}$  when deriving the derivatives. Such replacement is an approximation and its error depends on the form of the profit function. In this part, we analyze the error bound of the approximated mean-field approach.

**THEOREM 5.1.** *The exact weighted mean of all sellers' strategies by the direct derivation is defined as  $\bar{\tau}^{DD}$ , and the approximated one by the mean-field method is  $\bar{\tau}^{MF}$ . The error is  $\bar{\tau}^{DD} - \bar{\tau}^{MF}$ . Consider the case that the privacy loss function is  $L_i(\tau_i) = \lambda_i \chi_i \tau_i^2$ . When the number of sellers  $m$  is large and by scaling  $\omega_1, \omega_2, \dots, \omega_m$  such that  $\frac{\omega_i}{\lambda_i} \leq \frac{1}{p^D m^2}$ , we get*

$$-\frac{1}{6m^2} < \bar{\tau}^{DD} - \bar{\tau}^{MF} < \frac{1}{m} - \frac{2}{3m^2}.$$

*Note that what makes sense is the proportional relationship among  $\omega_i, i = 1, 2, \dots, m$ , allowing us to arbitrarily scale them.*

**PROOF.** We first calculate the upper bound of  $\bar{\tau}^{DD} - \bar{\tau}^{MF}$ . By applying direct derivation to Eq. 22, we can get

$$2p^D \sum_{j=1}^m \omega_j \tau_j - p^D \omega_i \tau_i = 3\lambda_i \tau_i \sum_{j=1}^m \omega_j \tau_j - \lambda_i \omega_i \tau_i^2.$$

By splitting  $\sum_{j=1}^m \omega_j \tau_j$  into  $\sum_{j=1, j \neq i}^m \omega_j \tau_j$  and  $\omega_i \tau_i$ , we can easily obtain a quadratic equation about  $\tau_i$  by deforming the above formula, and using root formula for the quadratic equation, we can get

$$\tau_i^* = \frac{p^D \omega_i - 3\lambda_i \sum_{j=1, j \neq i}^m \omega_j \tau_j + \sqrt{(3\lambda_i \sum_{j=1, j \neq i}^m \omega_j \tau_j - p^D \omega_i)^2 + 16p^D \lambda_i \omega_i \sum_{j=1, j \neq i}^m \omega_j \tau_j}}{4\lambda_i \omega_i}, \quad (24)$$

where  $\sum_{j=1, j \neq i}^m \omega_j \tau_j$ . With the constraint  $\frac{\omega_i}{\lambda_i} \leq \frac{1}{p^D m^2}$  we can justify that  $3\lambda_i \sum_{j=1, j \neq i}^m \omega_j \tau_j - p^D \omega_i > 0$  when  $m$  is very large. Thus according to  $\sqrt{x+y} < \sqrt{x} + \sqrt{y}$ , we can scale and deform the above formula to get

$$\omega_i \tau_i^* < \frac{\sqrt{16p^D \lambda_i \omega_i \sum_{j=1, j \neq i}^m \omega_j \tau_j}}{4\lambda_i}.$$

Further simplifying and scaling the above formula, we can get

$$\omega_i \tau_i^* < \sqrt{p^D \frac{\omega_i}{\lambda_i} \sum_{j=1}^m \omega_j \tau_j} \leq \sqrt{p^D \frac{\omega_i}{\lambda_i} \sum_{j=1}^m \omega_j \tau_j^*},$$

which applies to all  $\tau_i^*, i = 1, 2, \dots, m$ . Then, by adding  $m$  inequalities together and simplifying it, we can obtain

$$\sum_{i=1}^m \omega_i \tau_i^* < \left( \sum_{i=1}^m \sqrt{p^D \frac{\omega_i}{\lambda_i}} \right)^2,$$

and thus,

$$\bar{\tau}^{DD} = \frac{1}{m} \sum_{i=1}^m \omega_i \tau_i^* < \frac{1}{m} \left( \sum_{i=1}^m \sqrt{p^D \frac{\omega_i}{\lambda_i}} \right)^2.$$

Additionally, using Eqs. 21 and 23, we can derive  $\bar{\tau}^{MF}$  as below.

$$\bar{\tau}^{MF} = \frac{1}{m} \sum_{i=1}^m \frac{2p^D \omega_i}{3\lambda_i}.$$

Then we use  $\frac{\omega_i}{\lambda_i} \leq \frac{1}{p^D m^2}$  and get

$$\begin{aligned} \bar{\tau}^{DD} - \bar{\tau}^{MF} &< \frac{1}{m} \left( \sum_{i=1}^m \sqrt{p^D \frac{\omega_i}{\lambda_i}} \right)^2 - \frac{1}{m} \sum_{i=1}^m \frac{2p^D \omega_i}{3\lambda_i} \\ &\leq \frac{1}{m} \left( \sum_{i=1}^m \sqrt{\frac{1}{m^2}} \right)^2 - \frac{1}{m} \sum_{i=1}^m \frac{2}{3m^2} \\ &= \frac{1}{m} - \frac{2}{3m^2}. \end{aligned}$$

Next, we calculate the lower bound. Since  $(3\lambda_i \Sigma_{\tau_{-i}} - p^D \omega_i)^2 + 16p^D \lambda_i \omega_i \Sigma_{\tau_{-i}} > (p^D \omega_i + 3\lambda_i \Sigma_{\tau_{-i}})^2$ , using Eq. 24 we can get

$$\begin{aligned} \bar{\tau}^{DD} &= \frac{1}{m} \sum_{i=1}^m \omega_i \tau_i^* \\ &> \frac{1}{m} \sum_{i=1}^m \frac{p^D \omega_i - 3\lambda_i \Sigma_{\tau_{-i}} + \sqrt{(p^D \omega_i + 3\lambda_i \Sigma_{\tau_{-i}})^2}}{4\lambda_i} \\ &= \frac{1}{m} \sum_{i=1}^m \frac{p^D \omega_i}{2\lambda_i}, \end{aligned}$$

and using  $\frac{\omega_i}{\lambda_i} \leq \frac{1}{p^D m^2}$  again, we get

$$\begin{aligned} \bar{\tau}^{DD} - \bar{\tau}^{MF} &= \frac{1}{m} \sum_{i=1}^m \omega_i \tau_i^* - \frac{1}{m} \sum_{i=1}^m \frac{2p^D \omega_i}{3\lambda_i} \\ &> \frac{1}{m} \sum_{i=1}^m \frac{p^D \omega_i}{2\lambda_i} - \frac{1}{m} \sum_{i=1}^m \frac{2p^D \omega_i}{3\lambda_i} \\ &\geq -\frac{1}{6m^2}. \end{aligned}$$

Therefore, Theorem 5.1 holds.  $\square$

Through the above error analysis, we draw the following empirical conclusion: by scaling the value of  $\omega_i$  ( $i = 1, 2, \dots, m$ ) to satisfy certain conditions, the error of the mean-field approximation method will be bounded in an acceptable range and decrease as  $m$  increases when  $m$  is very large, especially when  $m$  approaches infinity, the error is approximately zero. Actually, this result is in line with the mean-field theory. When the number of sellers  $m$  is big enough, our proposed mean-field method appears reasonable in terms of error.

### 5.1.2 Expression of $p^{D*}$ in Stage 2.

We use the direct derivation to derive the expression of  $p^{D*}$  for the broker.

**Direct Derivation.** By substituting Eq. 20 into Eq. 13, we can get each  $\chi_i^*$

$$\chi_i^* = N \frac{\omega_i \tau_i^*}{\sum_{j=1}^m \omega_j \tau_j^*} = N \frac{\sqrt{\frac{\omega_i}{\lambda_i}}}{\sum_{j=1}^m \sqrt{\frac{\omega_j}{\lambda_j}}}.$$

Then we get

$$q^{D*} = \sum_{i=1}^m \chi_i^* \tau_i^* = \sum_{i=1}^m \frac{p^D}{2\lambda_i}.$$

Since  $q^M = h(q^D, v)$  is positively correlated to  $q^D$  and  $v$ , we instantiate  $h(q^D, v)$  as  $q^D v$ . We get

$$q^{M*} = q^{D*} v = \frac{1}{2} \sum_{i=1}^m \frac{1}{\lambda_i} p^D v.$$

By substituting  $q^{D*}$  and  $q^{M*}$  into  $\mathcal{A}$ 's profit function in Eq. 9, we get

$$\Omega(p^M, p^D, \tau) = p^M \cdot \left( \frac{1}{2} \sum_{i=1}^m \frac{1}{\lambda_i} p^D v \right) - C(N, v) - p^D \cdot \left( \frac{1}{2} \sum_{i=1}^m \frac{1}{\lambda_i} p^D \right).$$

We derive the first-order derivative with respect to  $p^D$  and let it equal to 0

$$\frac{\partial \Omega(p^M, p^D, \tau)}{\partial p^D} = \frac{1}{2} \sum_{i=1}^m \frac{1}{\lambda_i} v p^M - \sum_{i=1}^m \frac{1}{\lambda_i} p^D = 0,$$

and can get the expression of  $p^{D*}$

$$p^{D*} = \frac{v p^M}{2}. \quad (25)$$

Note that we justify that the second-order derivative  $\frac{\partial^2 \Omega(p^M, p^D, \tau)}{\partial p^{D^2}} = -\sum_{i=1}^m \frac{1}{\lambda_i} < 0$ , so the solution can maximize the broker's profit.

### 5.1.3 Value of $p^{M*}$ in Stage 1.

We as well use the direct derivation in this stage, and by using the results in Sections 5.1.1 and 5.1.2 we can directly derive the value rather than the expression of  $p^{M*}$  for the buyer.

**Direct Derivation.** By substituting Eq. 20 and Eq. 25 into  $\mathcal{B}$ 's profit function in Eq. 7, we can obtain the profit of  $\mathcal{B}$

$$\begin{aligned} \Phi(p^M, \tau) &= \theta_1 \ln(1 + \rho_1 q^{D*}) + \theta_2 \ln(1 + \rho_2 v) - p^M q^{M*} \\ &= \theta_1 \ln(1 + c_1 p^M) + \theta_2 \ln(1 + \rho_2 v) - \frac{c_2 \theta_1}{2} p^{M^2}, \end{aligned}$$

where  $c_1 = \frac{\rho_1 v}{4} \sum_{i=1}^m \frac{1}{\lambda_i}$  and  $c_2 = \frac{v^2}{2\theta_1} \sum_{i=1}^m \frac{1}{\lambda_i}$ .

Then, we derive the first-order derivative of  $\Phi(p^M, \tau)$  as follows.

$$\frac{\partial \Phi(p^M, \tau)}{\partial p^M} = \frac{\theta_1 c_1}{1 + c_1 p^M} - c_2 \theta_1 p^M. \quad (26)$$

By letting  $\frac{\partial \Phi(p^M, \tau)}{\partial p^M}$  in Eq. 26 equal to zero, we obtain the following equation

$$c_1 c_2 \cdot p^{M^2} + c_2 \cdot p^M - c_1 = 0.$$

Using characteristic root method, we find buyer  $\mathcal{B}$ 's optimal strategy  $p^{M*}$  (after discarding the negative solution)

$$p^{M*} = \frac{-c_2 + \sqrt{c_2^2 + 4c_1^2 c_2}}{2c_1 c_2}. \quad (27)$$

Note that we justify that the second-order derivative  $\frac{\partial^2 \Phi(p^M, \tau)}{\partial p^{M^2}} < 0$ , so the solution can maximize the buyer's profit.

After that, we can determine the optimal value of  $p^{D*}$  by substituting  $p^{M*}$  into Eq. 25 as well as each seller's optimal value of  $\tau_i^*$

by substituting  $p^{D*}$  into Eq. 20. Till now, the whole optimal incentive strategy profile  $\langle p^{M*}, p^{D*}, \tau^* \rangle$  has been determined, based on which the market trade can be conducted.

#### 5.1.4 Equilibrium Analysis.

Based on the above strategies, the market trade can be conducted. In this part, we prove the existence and uniqueness of SNE in *Share*.

**THEOREM 5.2.** *The whole optimal incentive strategy profile  $\langle p^{M*}, p^{D*}, \tau^* \rangle$  determined by backward induction approach uniquely constitutes SNE.*

**PROOF.**

*Existence.* For the buyer, when the broker and sellers hold the optimal strategy expressions in Eq. 25 and Eq. 20, the buyer's profit  $\Phi(p^M, \tau^*)$  only changes with  $p^M$ . In the process of deriving optimal  $p^{M*}$ , the first-order derivation is set to be 0 and the second-order is strictly less than 0, which means that the maximum profit is obtained at  $p^{M*}$ . Thus, Eq. 14 holds at  $p^{M*}$ . For the broker, when the buyer and sellers hold the optimal strategies  $p^{M*}$  in Eq. 27 and  $\tau^*$  in Eq. 20, the broker's optimal profit can be obtained at  $p^{D*}$  since the profit function is strictly concave and has a single extreme point  $p^{D*}$ . Thus, Eq. 15 holds at  $p^{D*}$ . For the sellers, each seller determines each  $\tau_i^*$  simultaneously in the same way, and  $\tau_i^*, i = 1, 2, \dots, m$  are jointly decided. For each seller  $S_i$ , her optimal strategy  $\tau_i^*$  is determined by letting the first-order derivation equal to 0. Since the profit function is strictly concave, the extreme point  $\tau_i^*$  maximizes seller  $S_i$ 's profit if  $\tau_i^* \leq 1$ . Otherwise, when the extreme point is larger than 1, the optimal value  $\tau_i^* = 1$  can also maximize  $S_i$ 's profit since the profit function is monotonically increasing in the feasible range of  $\tau_i$  and maximized at the right endpoint 1. Thus, Eq. 16 holds at  $\tau_i^*$ . Therefore, it's proved that SNE exists in our mechanism.

*Uniqueness.* For the buyer, since her profit function is strictly concave, the maximum profit is obtained only at  $p^{M*}$ . Any other value of  $p^M \neq p^{M*}$  will yield an inferior profit. Such result can be also explained by Convex Optimization [8], i.e., the strategy space of  $p^M$  is a convex and compact subspace of Euclidean space, and the profit function  $\Phi(\cdot)$  is a convex function of  $p^M$ , leading to the unique optimal  $p^{M*}$  that maximizes  $\Phi(\cdot)$ . Thus, Eq. 14 holds only at  $p^{M*}$ . For the broker, her profit function is also strictly concave and only has a single extreme point  $p^{D*}$ . Any other value of  $p^D \neq p^{D*}$  will lower the broker's profit. Thus, Eq. 15 holds only at  $p^{D*}$ . For sellers, seller  $S_i$  can only have lower profit by deciding  $\forall \tau_i \neq \tau_i^*$ . If  $\tau_i^* \leq 1$ , seller  $S_i$ 's profit function is concave and only maximized at  $\tau_i^*$ . Otherwise, the profit can only be maximized at the right endpoint,  $\tau_i^* = 1$ , since the profit function is monotonically increasing. For each seller, if she chooses other  $\forall \tau_i \neq \tau_i^*$ , she can only have lower profit with  $\tau_{-i}^*$  kept the same. Thus, the unique Nash equilibrium among sellers is achieved and Eq. 16 holds only at  $\tau_i^*$ . Therefore, it's proved that other strategy profile except our solution cannot satisfy SNE, which infers the uniqueness of SNE in our mechanism.  $\square$

## 5.2 Complete Data Trading Dynamics

We summarize the complete dynamics of data markets in Algorithm 1, which integrates the equilibrium solving process in the previous section.

The first phase is *Parameter Collection*. We assume that each party can give their specific input parameters. The buyer sets appropriate parameters  $\theta_1, \theta_2, \rho_1, \rho_2$  for her utility function and proposes demand parameters  $v$  and  $N$  for the product (Line 2). Note that the form of the product is not restricted from simple data aggregation to deep learning models. The broker determines  $\sigma_k, k \in \{0, 1, 2, 3, 4, 5\}$  for her translog cost function and maintains the weights  $\omega_i, i = 1, 2, \dots, m$  of sellers' datasets (Line 3). Specifically, the weights are set the same when the market is established. To decide the real weights before the first transaction, the broker can use dummy buyers to iterate several times, and in each iteration, Shapley values of sellers' datasets are calculated after production and can be used to update weights in the next iteration. Sellers give their privacy sensitivity  $\lambda_i, i = 1, 2, \dots, m$  (Line 4). We assume that participants provide their truthful parameters in line with the practical situation under the supervision of market regulators (e.g., by regular spot-check). Plus,  $m$  is fixed according to the number of sellers in practice (Line 5).

The second phase is *Strategy Decision*. Using the strategy mechanism, buyer  $\mathcal{B}$ , broker  $\mathcal{A}$ , and sellers  $S_i, i = 1, 2, \dots, m$  give strategies of the product price  $p^{M*}$ , data price  $p^{D*}$ , and data fidelity  $\tau_i^*$  in order according to Eqs. 27, 25, 20, respectively (Line 7).

Then *Data Transaction* between the broker and sellers begins. The data quantity chosen from each seller can be calculated according to Eq. 13 (Line 9). Next, each seller picks  $\chi_i^*$ -sized dataset (Line 11) and pre-processes it for privacy protection based on  $\epsilon_i^*$  calculated from Eq. 10 (Lines 12-13). After that, seller  $S_i$  gives her protected dataset  $D_i^t$  to the broker in exchange of the compensation  $p^{D*} q_i^{D*}$  (Line 14).

Next phase is *Product Production*. The broker collects the data as  $D^t$  and uses it to make the product (Line 16). Moreover, the weights of sellers' datasets are updated by the broker based on their corresponding contributions to the data product (Line 17). We give one update formula based on Shapley value as an example:  $\omega'_i = 0.2\omega_i + 0.8S^tV_i, i = 1, 2, \dots, m$  where  $S^tV_i$  is the Shapley value of  $D_i^t$  to the product, and the updated weights  $\omega'_i$  can be used in the subsequent transaction.

The last phase is *Product Transaction* between the broker and the buyer. The broker gives the product to the buyer and the buyer pays  $p^{M*} q^{M*}$  to the broker (Line 19). So far, the current round of data trading among buyer  $\mathcal{B}$ , broker  $\mathcal{A}$ , and sellers  $S_i, i = 1, 2, \dots, m$  has finished. When the next buyer comes, the next round of data trading will start and can use the updated  $\omega'_i, i = 1, 2, \dots, m$ .

**Time Complexity.** As seen from Algorithm 1, the phase of *Parameter Collection* costs  $O(m)$  since  $m$  sellers need to provide  $\lambda_i, i = 1, 2, \dots, m$ . *Strategy Decision* costs  $O(m)$  based on the optimal strategy profile. *Data Transaction* costs  $O(m + N)$  because each seller needs to form and protect her  $\chi_i$ -sized dataset, and totally  $N$  data pieces are preprocessed and sold to the broker. The time cost of *Product Production* depends on the exact product type, production mode, and the way of updating the weights  $\omega_1, \omega_2, \dots, \omega_m$ . The last phase of *Product Transaction* takes constant time. Therefore, the

**Algorithm 1:** Data trading algorithm.

- 
- 1 %% Parameter Collection;
  - 2 From the current buyer  $\mathcal{B}$ , parameters  $N, v, \theta_1, \theta_2, \rho_1, \rho_2$  are provided;
  - 3 From broker  $\mathcal{A}$ ,  $\sigma_k (k \in \{0, 1, 2, 3, 4, 5\}), \omega_i (i = 1, 2, \dots, m)$  are given;
  - 4 From existing  $m$  sellers, each seller  $\mathcal{S}_i$  decides  $\lambda_i$ ;
  - 5  $m$  is given according to the practical situation;
  - 6 %% Strategy Decision;
  - 7 Through three-stage Stackelberg-Nash game, the optimal strategy profile  $\langle p^{M*}, p^{D*}, \tau^* \rangle$  is determined by the buyer, the broker, and sellers, respectively;
  - 8 %% Data Transaction;
  - 9 The quantity of data each seller can sell, i.e.,  $\chi^*$ , is calculated according to Eq. 13;
  - 10 **for** each seller  $\mathcal{S}_i, i = 1, 2, \dots, m$  **do**
    - 11 Randomly pick  $\chi_i^*$  data pieces from her dataset  $D_i$ ;
    - 12 Calculate  $\epsilon_i^*$  from the strategy  $\tau_i^*$  according to Eq. 10;
    - 13 Conduct LDP with  $\epsilon_i^*$  on her  $\chi_i^*$ -sized dataset, and then give the protected  $D_i^t$  to broker  $\mathcal{A}$ ;
  - 14 Broker  $\mathcal{A}$  gets data from sellers to form dataset  $D^t$  for production and pays compensation  $p^{D*} q_i^{D*}$  to each seller;
  - 15 %% Product Production;
  - 16 Broker  $\mathcal{A}$  then uses  $D^t$  to produce the data product;
  - 17 After manufacturing the product, broker  $\mathcal{A}$  updates  $\omega_1, \omega_2, \dots, \omega_m$  (might scale down proportionally as needed) based on the contribution to the product from each seller's  $D_i^t$ ;
  - 18 %% Product Transaction;
  - 19 Broker  $\mathcal{A}$  gives the product to buyer  $\mathcal{B}$  and meantime  $\mathcal{B}$  pays  $p^{M*} q^{M*}$  to  $\mathcal{A}$ .
- 

complexity of data trading algorithm excluding *Product Production* is  $O(m + N)$ .

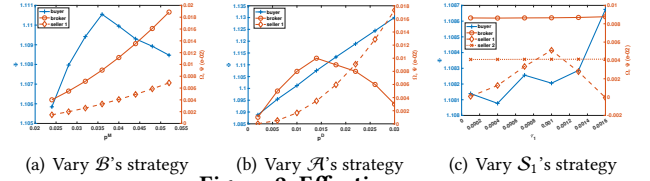
## 6 EXPERIMENTS

In this section, we present experimental studies validating the effectiveness and efficiency of *Share*. We first describe our experiment setup including the used datasets and parameter settings in Section 6.1. Sections 6.2 and 6.3 show the results verifying the effectiveness and efficiency of *Share*, respectively. Section 6.4 shows the effect of parameters used in *Share*.

### 6.1 Experiment Setup

We conduct experiments on a machine with an Intel Core i7-11700KF running Ubuntu with 64GB memory. We choose Linear Regression model as the data product.

**Datasets.** We conduct extensive experiments on a real dataset of a Combined Cycle Power Plant (CCPP) [16]. The dataset contains 9,568 data points, each having four features, i.e., hourly average ambient variables Temperature (AT), Exhaust Vacuum (V), Ambient Pressure (AP), and Relative Humidity (RH). The Linear Regression



**Figure 2: Effectiveness.**

task is to predict the net hourly electrical energy output of the plant. Besides the real dataset, we augment CCPP by first replicating 100 times and then adding Gaussian noise  $\mathcal{N}(0, 0.1^2)$  to generate a synthetic dataset with the size of 1,000,000 to test the efficiency of *Share*.

**Parameter Settings.** Our parameters include the number of sellers  $m$ , model requirements demanded by the buyer, i.e., the needed total data quantity  $N$  and the model performance  $v$  which refers to the model explained variance, and individual parameters of each party, i.e., buyer  $\mathcal{B}$ 's  $\theta_1, \theta_2, \rho_1, \rho_2$  related to model utility, broker  $\mathcal{A}$ 's cost parameters  $\sigma_k, k \in \{0, 1, 2, 3, 4, 5\}$ , and seller  $\mathcal{S}_i$ 's privacy sensitivity  $\lambda_i, i = 1, 2, \dots, m$ . Without loss of generality, we set the buyer's requirement parameters as  $N = 500, v = 0.8$  and utility parameters as  $\theta_1 = 0.5, \theta_2 = 0.5, \rho_1 = 0.5, \rho_2 = 250$  (in order to balance the impacts of the product quality and the dataset quality). The cost parameters of the broker are related to the practical manufacturing situation and are set as default values  $\sigma_0 = 1 \times 10^{-3}, \sigma_1 = -2, \sigma_2 = -3, \sigma_3 = 1 \times 10^{-3}, \sigma_4 = 2 \times 10^{-3}, \sigma_5 = 1 \times 10^{-3}$ . We set  $m = 100$ , and sellers'  $\lambda_i, i = 1, 2, \dots, m$  are picked randomly in  $(0, 1)$ .

We assume each seller owns a dataset of the same size and distribute 9,000 data pieces of the CCPP dataset (the remaining data is used for test) equally to 100 sellers. In the real world, the datasets of sellers vary in quality. To simulate this characteristic, we first sort data by quality measured by Shapley value, which indicates the contribution of each data piece to model training. The Shapley value is calculated based on Monte Carlo Method [9][21] with the permutation number set as 100. Then by distributing data over sellers, we get sellers each of whom owns 90 data pieces but with different quality. Laplace mechanism [18] is used to ensure local differential privacy for each seller.

$\omega_1, \omega_2, \dots, \omega_m$  are generated by using buyer  $\mathcal{B}$  as the dummy buyer to iterate the mechanism which takes five times to stabilize the profits. We consider buyer  $\mathcal{B}$  as a general buyer coming after several transactions have finished. Shapley values of sellers' datasets can be calculated after model training to update weights for the next transaction.

We will show the results of using the direct derivation in the mechanism. Note that the mean-field approach used when the direct derivation fails functions the same to the markets in terms of the effectiveness, efficiency, and parameter influence.

### 6.2 Effectiveness

We implement the mechanism and unilaterally change the strategies of the buyer, the broker, and sellers respectively to verify the profit maximization of all parties as well as the corresponding equilibrium state.

Fig.2(a) shows the results of profits when we change  $p^M$  around the optimal strategy  $p^{M*}$  while maintaining the rest. Seller  $S_1$  acts as a representative of sellers. It's found that the peak of the buyer's profit  $\Phi(\cdot)$  appears when her optimal strategy  $p^{M*} = 0.036$  determined in SNE is adopted. Whatever strategy the buyer chooses except  $p^{M*}$ , she will get lower profit when all other participants' strategies fixed. The change of the profits of the broker and the seller is intuitive. Specifically, with growing  $p^M$ , the broker can gain more profit, which can further add the compensations for sellers and make their profits higher as well.

Fig.2(b) shows the results of profits when we change  $p^D$  around the optimal strategy  $p^{D*} = 0.014$  while maintaining the rest. Similarly, it is found that the broker cannot increase her profit by unilaterally changing her strategy. The change of the profits of the buyer and seller is also intuitive. Specifically, the growing  $p^D$  brings more compensations to sellers, adding their profits. Due to more compensations, the dataset quality from sellers can therefore be improved, which causes the rise of the buyer's profit.

Fig.2(c) shows the results of profits when we change  $\tau_1$  around the optimal strategy  $\tau_1^* = 0.001$  while maintaining the rest. Note that the first two sellers  $S_1$  and  $S_2$  are chosen as representatives. It is the same that the seller who changes her strategy unilaterally gets no more profit. Even if one seller changes her strategy, the broker can nearly keep her profit as before, which indicates the *inner* competition among sellers formulated as Nash game. Specifically, the effect of sellers' bounded rationality is almost limited among sellers and is corrected automatically in Stage 3, which signifies the transparency of Stage 3 to upper stages. The change of the buyer's profit may be due to the effect of data on the model, which is not always predictable, causing the irregular curve of  $\Phi(\cdot)$ . In theory, varying  $\tau_1$  surely makes differences on other sellers. However, since the number of sellers is large, this effect is *diluted* and negligible, making the profit of  $S_2$  almost unchanged.

### 6.3 Efficiency

Fig.3(a) and Fig.3(b) show the runtime of the proposed data trading algorithm with and without Shapley value to update weights, respectively. We use the synthetic dataset with 1,000,000 data tuples and adjust the number of sellers  $m$  from 5 to 10,000 while fixing the other parameters and the average number of data pieces chosen from one seller being 100. Fig.3(a) shows that the runtime grows as  $m$  goes higher but with an acceptable rate. Even when  $m = 10,000$ , it does not take too much time. Note that our mechanism contains an extremely time-consuming part to calculate Shapley values. Fig.3(b) shows that our mechanism without Shapley value calculation can run very fast with a linear time complexity, which corresponds to the complexity analysis of Algorithm 1 in Section 5.2.

### 6.4 Parameter Influence

In this section, we make sensitivity analyses on the parameters in our mechanism. Fig.4, Fig.5, Fig.6, Fig.7, and Fig.8 investigate how the parameters affect the strategies and profits of three parties.

Fig.4(a) and Fig.4(b) present the effect of  $\theta_1$  on strategies and profits, respectively. Note that  $\theta_1$  and  $\theta_2$  refer to the buyer's concern on the manufacturing dataset and the product performance, respectively. Since the sum of  $\theta_1$  and  $\theta_2$  has been set as a constant,

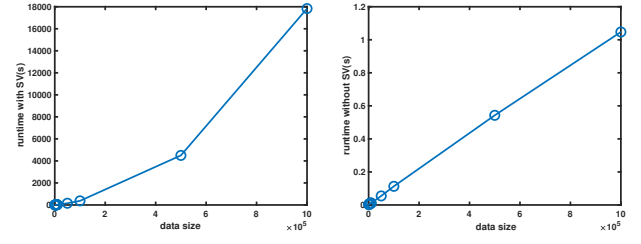


Figure 3: Efficiency.

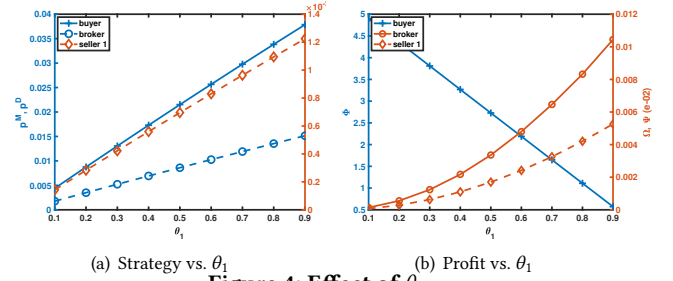


Figure 4: Effect of  $\theta_1$ .

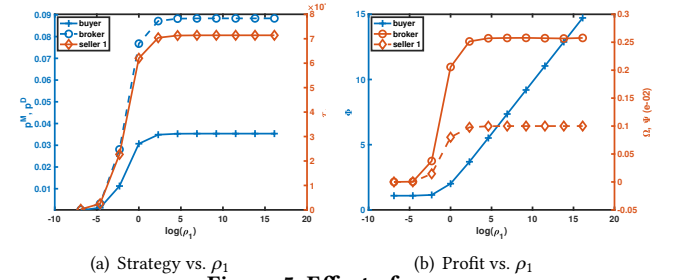
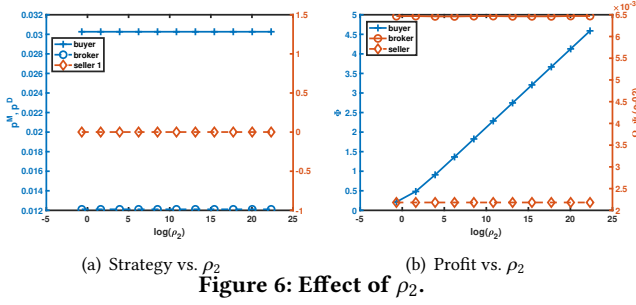
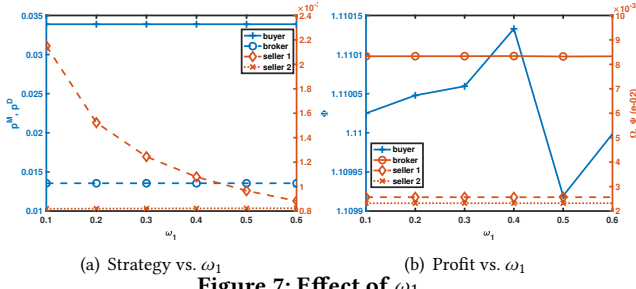


Figure 5: Effect of  $\rho_1$ .

we only choose  $\theta_1$  to make an analysis here. We alter  $\theta_1$  from 0.1 to 0.9 ( $\theta_2$  from 0.9 to 0.1 accordingly) and let other parameters changeless. Fig.4(a) shows that all the strategies boost in a linear rate with increasing  $\theta_1$ , which is intuitive since more concern on data leads to higher prices and better data. Fig.4(b) shows that the profit of the buyer decreases in a linear trend, and the profits of other two parties increase with  $\theta_1$  rising. The change of the buyer's profit indicates the difference between the significance of product performance and manufacturing dataset to the product utility in different tasks, and the descent of  $\Phi(\cdot)$  here suggests that in this task the product plays a more important role in the buyer's profit than the manufacturing dataset.

Fig.5(a) and Fig.5(b) present the effect of  $\rho_1$  on strategies and profits, respectively. Note that  $\rho_1$  is a parameter relevant to the buyer's sensitivity to the dataset quality, which objectively reflects the relationship between the dataset quality and the product utility. Fig.5(a) shows that too small  $\rho_1$  can hardly lead to effective markets due to the buyer's indifference on data. When  $\rho_1$  reaches a certain level, all the strategies keep the same and the market reaches the equilibrium. The influence of  $\rho_1$  is limited within the utility for the

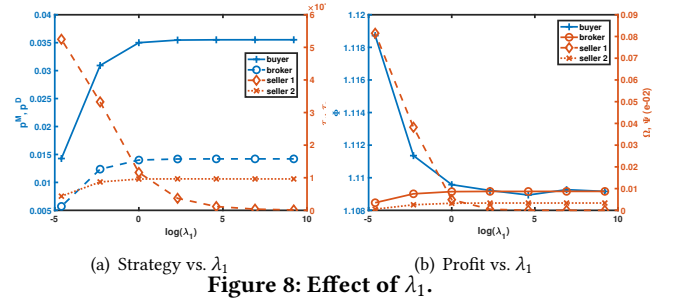
Figure 6: Effect of  $\rho_2$ .Figure 7: Effect of  $\omega_1$ .

buyer and can no longer disturb the market equilibrium, which may be due to the common sense that the dataset quality cannot increase unlimitedly and even with sharper sensitivity to data, higher prices wouldn't bring about better data any more. Fig.5(b) shows that the profit of the buyer surges as  $\rho_1$  increases because she will get more utility from the raise of the dataset quality. When  $\rho_1$  is big enough, the increase of  $\rho_1$  has little effect on the profits of the broker and sellers, which can be explained by the trend of strategies discussed above.

Fig.6(a) and Fig.6(b) present the effect of  $\rho_2$  on strategies and profits, respectively. Note that  $\rho_2$  is related to the sensitivity of the buyer to the product performance. Fig.6(a) shows that  $\rho_2$  influences little on the strategies. Fig.6(b) shows that the buyer gains more as  $\rho_2$  is raised due to the semblable reason in the above analysis about  $\rho_1$ . The profits of the broker and sellers keep unchanged because the broker still aims to achieve the buyer's fixed demanding  $v$  in the production and sellers only deal with the data not the product.

Fig.7(a) and Fig.7(b) present the effect of  $\omega_1$  on strategies and profits, respectively. Note that  $\omega_1, \omega_2, \dots, \omega_m$  are the weights of sellers' datasets and assess the sellers' data in previous transactions. We select  $S_1$  and  $S_2$  as representatives. Fig.7(a) shows that  $\omega_1$  only affects the strategy of the corresponding seller  $S_1$ . The strategies of the buyer and the broker remain the same because  $\omega_1$  only affects the inner competition among sellers. Since the number of sellers is large, varying  $\omega_1$  makes little difference on other sellers, making the strategy of  $S_2$  almost unchanged. Fig.7(b) shows that when  $\omega_1$  varies from 0.1 to 0.6, all profits except the buyer's are stable. Once  $\omega_1$  gets a non-appropriate value, the data of this seller  $S_1$  won't work as expected and affects the profit of the buyer, leading to the unsmooth curve of  $\Phi(\cdot)$ .

Fig.8(a) and Fig.8(b) show the effect of  $S_1$ 's parameter  $\lambda_1$  on strategies and profits, respectively. Note that  $\lambda_i$  is related to seller

Figure 8: Effect of  $\lambda_1$ .

$S_i$ 's privacy sensitivity. Fig.8(a) shows that  $\tau_1$  sinks since seller  $S_1$  will strengthen her data protection if more sensitive to privacy risks.  $p^M$  and  $p^D$  increase possibly because higher prices may be provided to encourage conservative sellers to offer high-fidelity data in spite of heavy privacy risks. Fig.8(b) shows that  $\lambda_1$  mainly influences on the buyer's and the corresponding seller  $S_1$ 's profits. The profit of  $S_1$  decreases because bigger  $\lambda_1$ , more privacy loss  $S_1$  will suffer. The profit of buyer  $B$  dives probably because the seller would enhance the protection on her data when faced with huge privacy risks, thus lowering down the data fidelity and further harming the buyer's profit. The profit of the broker remains unchanged because the broker herself does not rely on data fidelity but just transfers data from the sellers to the buyer.

## 7 CONCLUSION AND FUTURE WORK

We presented the first Stackelberg-Nash game based data markets *Share*. A three-stage Stackelberg-Nash game is proposed to model the data trading mechanism in the buyer-leading data markets with absolute pricing rules. The mutual interaction among three parties (buyers, brokers, and sellers) is considered as a three-stage Stackelberg game, which maximizes the profits of all participants. The inner competition among sellers is considered as a Nash game, which neatly solves the seller selection problem. To derive SNE, the backward induction approach is used. Specifically, to solve the inner Nash game, we propose two methods, the direct derivation and a novel mean-field approximation which can address complex cases with provable approximation guarantees. Our proposed data market model performs well on the real and synthetic datasets in terms of both effectiveness and efficiency.

Our data market model can be easily adapted to a variety of market settings, e.g., broker-leading instead of buyer-leading. Thus we believe the market dynamic described in this work can be a natural and scalable way for data trading. There are also interesting and practical issues to address for deployment. For example, the deficiency of real-world historical trading records brings about the challenge of parameter fitting for each party.

## REFERENCES

- [1] [n.d.]. <https://support.gnip.com/apis/>. ([n.d.]).
- [2] [n.d.]. <https://www.bloomberg.com/professional/product/market-data/>. ([n.d.]).
- [3] [n.d.]. SafeGraph, <https://www.safegraph.com/>.
- [4] Anish Agarwal, Munther Dahleh, and Tuhin Sarkar. 2019. A marketplace for data: An algorithmic solution. In *Proceedings of the 2019 ACM Conference on Economics and Computation*. ACM, 701–726.
- [5] Baoyi An, Mingjun Xiao, An Liu, Xike Xie, and Xiaofang Zhou. 2021. Crowdsensing Data Trading based on Combinatorial Multi-Armed Bandit and Stackelberg



- Game. In *2021 IEEE 37th International Conference on Data Engineering (ICDE)*. IEEE, 253–264.
- [6] Gaurang Bansal and Biplob Sikdar. 2021. Security Service Pricing Model for UAV Swarms: A Stackelberg Game Approach. In *2021 IEEE Conference on Computer Communications Workshops, INFOCOM Workshops 2021, Vancouver, BC, Canada, May 10–13, 2021*. IEEE, 1–6. <https://doi.org/10.1109/INFOCOMWKSHSPS51825.2021.9484577>
  - [7] Tamer Başar and Geert Jan Olsder. 1998. *Dynamic noncooperative game theory*. SIAM.
  - [8] Stephen Boyd, Stephen P Boyd, and Lieven Vandenbergh. 2004. *Convex optimization*. Cambridge university press.
  - [9] Javier Castro, Daniel Gómez, and Juan Tejada. 2009. Polynomial calculation of the Shapley value based on sampling. *Computers & OR* 36, 5 (2009), 1726–1730. <https://doi.org/10.1016/j.cor.2008.04.004>
  - [10] Lingjiao Chen, Paraschos Koutris, and Arun Kumar. 2019. Towards Model-based Pricing for Machine Learning in a Data Marketplace. In *SIGMOD*, Peter A. Boncz, Stefan Manegold, Anastasia Ailamaki, Amol Deshpande, and Tim Kraska (Eds.). ACM, 1535–1552. <https://doi.org/10.1145/3299869.3300078>
  - [11] Laurits R Christensen, Dale W Jorgenson, and Lawrence J Lau. 1975. Transcendental logarithmic utility functions. *The American Economic Review* 65, 3 (1975), 367–383.
  - [12] Zicun Cong, Xuan Luo, Jian Pei, Feida Zhu, and Yong Zhang. 2022. Data pricing in machine learning pipelines. *Knowl. Inf. Syst.* 64, 6 (2022), 1417–1455. <https://doi.org/10.1007/s10115-022-01679-4>
  - [13] Vincent Conitzer and Tuomas Sandholm. 2003. Complexity Results about Nash Equilibria. In *IJCAI-03, Proceedings of the Eighteenth International Joint Conference on Artificial Intelligence, Acapulco, Mexico, August 9–15, 2003*, Georg Gottlob and Toby Walsh (Eds.). Morgan Kaufmann, 765–771. <http://ijcai.org/Proceedings/03/Papers/111.pdf>
  - [14] Constantinos Daskalakis, Paul W Goldberg, and Christos H Papadimitriou. 2009. The complexity of computing a Nash equilibrium. *SIAM J. Comput.* 39, 1 (2009), 195–259.
  - [15] DAWEX [n.d.]. <https://www.dawex.com/en/>
  - [16] Dheeru Dua and Casey Graff. 2017. UCI Machine Learning Repository. <http://archive.ics.uci.edu/ml>
  - [17] John C. Duchi, Michael I. Jordan, and Martin J. Wainwright. 2013. Local Privacy and Statistical Minimax Rates. In *54th Annual IEEE Symposium on Foundations of Computer Science, FOCS 2013, 26–29 October, 2013, Berkeley, CA, USA*. IEEE Computer Society, 429–438. <https://doi.org/10.1109/FOCS.2013.53>
  - [18] Cynthia Dwork. 2006. Differential Privacy. In *Automata, Languages and Programming, 33rd International Colloquium, ICALP 2006, Venice, Italy, July 10–14, 2006, Proceedings, Part II (Lecture Notes in Computer Science, Vol. 4052)*, Michele Bugliesi, Bart Preneel, Vladimiro Sassone, and Ingo Wegener (Eds.). Springer, 1–12. [https://doi.org/10.1007/11787006\\_1](https://doi.org/10.1007/11787006_1)
  - [19] Cynthia Dwork and Jing Lei. 2009. Differential privacy and robust statistics. In *Proceedings of the 41st Annual ACM Symposium on Theory of Computing, STOC 2009, Bethesda, MD, USA, May 31 - June 2, 2009*, Michael Mitzenmacher (Ed.). ACM, 371–380. <https://doi.org/10.1145/1536414.1536466>
  - [20] Cynthia Dwork and Aaron Roth. 2014. The Algorithmic Foundations of Differential Privacy. *Found. Trends Theor. Comput. Sci.* 9, 3–4 (2014), 211–407. <https://doi.org/10.1561/04000000042>
  - [21] S. Shaheen Fatima, Michael J. Wooldridge, and Nicholas R. Jennings. 2008. A linear approximation method for the Shapley value. *Artif. Intell.* 172, 14 (2008), 1673–1699. <https://doi.org/10.1016/j.artint.2008.05.003>
  - [22] Raul Castro Fernandez. 2022. Protecting Data Markets from Strategic Buyers. In *SIGMOD '22: International Conference on Management of Data, Philadelphia, PA, USA, June 12 - 17, 2022*, Zachary Ives, Angela Bonifati, and Amr El Abbadi (Eds.). ACM, 1755–1769. <https://doi.org/10.1145/3514221.3517855>
  - [23] D. Gale and L. S. Shapley. 1962. College Admissions and the Stability of Marriage. *The American Mathematical Monthly* 69, 1 (1962), 9–15. <http://www.jstor.org/stable/2312726>
  - [24] Hui Gao, Chi Harold Liu, Jian Tang, Dejun Yang, Pan Hui, and Wendong Wang. 2019. Online Quality-Aware Incentive Mechanism for Mobile Crowd Sensing with Extra Bonus. *IEEE Trans. Mob. Comput.* 18, 11 (2019), 2589–2603. <https://doi.org/10.1109/TMC.2018.2877459>
  - [25] Amirata Ghorbani and James Zou. 2019. Data shapley: Equitable valuation of data for machine learning. In *International Conference on Machine Learning*. PMLR, 2242–2251.
  - [26] John C Harsanyi. 1967. Games with incomplete information played by “Bayesian” players, I–III Part I. The basic model. *Management science* 14, 3 (1967), 159–182.
  - [27] Ken-Ichi Inada. 1963. On a two-sector model of economic growth: Comments and a generalization. *The Review of Economic Studies* 30, 2 (1963), 119–127.
  - [28] IOTA [n.d.]. <https://data.iota.org/>
  - [29] Chandra K. Jaggi, Mamta Gupta, Amrina Kausar, and Sunil Tiwari. 2019. Inventory and credit decisions for deteriorating items with displayed stock dependent demand in two-echelon supply chain using Stackelberg and Nash equilibrium solution. *Ann. Oper. Res.* 274, 1–2 (2019), 309–329. <https://doi.org/10.1007/s10479-018-2925-9>
  - [30] Ruoxi Jia, David Dao, Boxin Wang, Frances Ann Hubis, Nezihe Merve Gurel, Bo Li, Ce Zhang, Costas Spanos, and Dawn Song. 2019. Efficient task-specific data valuation for nearest neighbor algorithms. *Proceedings of the VLDB Endowment* 12, 11 (2019), 1610–1623.
  - [31] Changkun Jiang, Lin Gao, Lingjie Duan, and Jianwei Huang. 2018. Data-Centric Mobile Crowdsensing. *IEEE Trans. Mob. Comput.* 17, 6 (2018), 1275–1288. <https://doi.org/10.1109/TMC.2017.2763956>
  - [32] Paraschos Koutris, Prasang Upadhyaya, Magdalena Balazinska, Bill Howe, and Dan Suciu. 2012. Query-based data pricing. In *PODS*. ACM, 167–178.
  - [33] Paraschos Koutris, Prasang Upadhyaya, Magdalena Balazinska, Bill Howe, and Dan Suciu. 2013. Toward practical query pricing with QueryMarket. In *SIGMOD*. ACM, 613–624.
  - [34] Jean-Michel Lasry and Pierre-Louis Lions. 2007. Mean field games. *Japanese journal of mathematics* 2, 1 (2007), 229–260.
  - [35] Chao Li, Daniel Yang Li, Jerome Miklau, and Dan Suciu. 2014. A Theory of Pricing Private Data. *ACM Trans. Database Syst.* 39, 4 (2014), 34:1–34:28. <https://doi.org/10.1145/2691190.2691191>
  - [36] Henger Li, Wen Shen, and Zizhan Zheng. 2020. Spatial-Temporal Moving Target Defense: A Markov Stackelberg Game Model. In *Proceedings of the 19th International Conference on Autonomous Agents and Multiagent Systems, AAMAS '20, Auckland, New Zealand, May 9–13, 2020*, Amal El Fallah Seghrouchni, Gita Sukthankar, Bo An, and Neil Yorke-Smith (Eds.). International Foundation for Autonomous Agents and Multiagent Systems, 717–725. <https://dl.acm.org/doi/abs/10.5555/3398761.3398847>
  - [37] Man Li, Jiahui Qin, Qichao Ma, Wei Xing Zheng, and Yu Kang. 2021. Hierarchical Optimal Synchronization for Linear Systems via Reinforcement Learning: A Stackelberg-Nash Game Perspective. *IEEE Trans. Neural Networks Learn. Syst.* 32, 4 (2021), 1600–1611. <https://doi.org/10.1109/TNNLS.2020.2985738>
  - [38] Jinfei Liu, Qiongqiong Lin, Jiayao Zhang, Kui Ren, Jian Lou, Junxu Liu, Li Xiong, Jian Pei, and Jimeng Sun. 2021. Demonstration of Dealer: An End-to-End Model Marketplace with Differential Privacy. *Proc. VLDB Endow.* 14, 12 (2021), 2747–2750. <http://www.vldb.org/pvldb/vol14/p2747-zhang.pdf>
  - [39] Jinfei Liu, Jian Lou, Junxu Liu, Li Xiong, Jian Pei, and Jimeng Sun. 2021. Dealer: An End-to-End Model Marketplace with Differential Privacy. *Proc. VLDB Endow.* 14, 6 (2021), 957–969. <http://www.vldb.org/pvldb/vol14/p957-liu.pdf>
  - [40] Frank McSherry and Kunal Talwar. 2007. Mechanism Design via Differential Privacy. In *48th Annual IEEE Symposium on Foundations of Computer Science (FOCS 2007), October 20–23, 2007, Providence, RI, USA, Proceedings*. IEEE Computer Society, 94–103. <https://doi.org/10.1109/FOCS.2007.41>
  - [41] Oskar Morgenstern and John Von Neumann. 1953. *Theory of games and economic behavior*. Princeton university press.
  - [42] Anna Nagurney and Pritha Dutta. 2019. Supply chain network competition among blood service organizations: a Generalized Nash Equilibrium framework. *Ann. Oper. Res.* 275, 2 (2019), 551–586. <https://doi.org/10.1007/s10479-018-3029-2>
  - [43] John Nash. 1951. Non-cooperative games. *Annals of mathematics* (1951), 286–295.
  - [44] John F Nash Jr. 1950. Equilibrium points in n-person games. *Proceedings of the national academy of sciences* 36, 1 (1950), 48–49.
  - [45] Chaoyue Niu, Zhenzhe Zheng, Fan Wu, Shaojie Tang, Xiaofeng Gao, and Guihai Chen. 2018. Unlocking the Value of Privacy: Trading Aggregate Statistics over Private Correlated Data. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining, KDD 2018, London, UK, August 19–23, 2018*, Yike Guo and Faisal Farooq (Eds.). ACM, 2031–2040. <https://doi.org/10.1145/3219819.3220013>
  - [46] Jian Pei. 2021. A Survey on Data Pricing: from Economics to Data Science. *IEEE Trans. Knowl. Data Eng.* (2021). <https://doi.org/10.1109/TKDE.2020.3045927>
  - [47] Jian Pei, Feida Zhu, Zicun Cong, Xuan Luo, Huiwen Liu, and Xin Mu. 2021. Data Pricing and Data Asset Governance in the AI Era. In *KDD '21: The 27th ACM SIGKDD Conference on Knowledge Discovery and Data Mining, Virtual Event, Singapore, August 14–18, 2021*, Feida Zhu, Beng Chin Ooi, and Chunyan Miao (Eds.). ACM, 4058–4059. <https://doi.org/10.1145/3447548.3470818>
  - [48] Alvin E. Roth. 2016. Lloyd Shapley (1923–2016). *Nat.* 532, 7598 (2016), 178. <https://doi.org/10.1038/532178a>
  - [49] Reinhard Selten. 1965. Spieltheoretische behandlung eines oligopolmodells mit nachfragertragheit: Teil i: Bestimmung des dynamischen preisgleichgewichts. *Zeitschrift für die gesamte Staatswissenschaft/Journal of Institutional and Theoretical Economics* H. 2 (1965), 301–324.
  - [50] Lloyd S Shapley. 1953. A value for n-person games. *Contributions to the Theory of Games* 2, 28 (1953), 307–317.
  - [51] Arunesh Sinha, Fei Fang, Bo An, Christopher Kiekintveld, and Milind Tambe. 2018. Stackelberg Security Games: Looking Beyond a Decade of Success. In *Proceedings of the Twenty-Seventh International Joint Conference on Artificial Intelligence, IJCAI 2018, July 13–19, 2018, Stockholm, Sweden*, Jérôme Lang (Ed.). ijcai.org, 5494–5501. <https://doi.org/10.24963/ijcai.2018/775>
  - [52] Ankur Sinha, Pekka Malo, and Kalyanmoy Deb. 2018. A Review on Bilevel Optimization: From Classical to Evolutionary Approaches and Applications. *IEEE Trans. Evol. Comput.* 22, 2 (2018), 276–295. <https://doi.org/10.1109/TEVC.2017.2712906>



- [53] Mohammad Amin Tajeddini, Hamed Kebriaei, and Luigi Glielmo. 2020. Decentralized Hierarchical Planning of PEVs Based on Mean-Field Reverse Stackelberg Game. *IEEE Trans Autom. Sci. Eng.* 17, 4 (2020), 2014–2024. <https://doi.org/10.1109/TASE.2020.2986374>
- [54] Heinrich Von Stackelberg. 2010. *Market structure and equilibrium*. Springer Science & Business Media.
- [55] Omar Abdel Wahab, Jamal Bentahar, Hadi Otrouk, and Azzam Mourad. 2021. Resource-Aware Detection and Defense System against Multi-Type Attacks in the Cloud: Repeated Bayesian Stackelberg Game. *IEEE Trans. Dependable Secur. Comput.* 18, 2 (2021), 605–622. <https://doi.org/10.1109/TDSC.2019.2907946>
- [56] Kaidi Wang, Zhiguo Ding, Daniel K. C. So, and George K. Karagiannis. 2021. Stackelberg Game of Energy Consumption and Latency in MEC Systems With NOMA. *IEEE Trans. Commun.* 69, 4 (2021), 2191–2206. <https://doi.org/10.1109/TCOMM.2021.3049356>
- [57] Tianhao Wang, Johannes Rausch, Ce Zhang, Ruoxi Jia, and Dawn Song. 2020. A principled approach to data valuation for federated learning. In *Federated Learning*. Springer, 153–167.
- [58] Shuyue Wei, Yongxin Tong, Zimu Zhou, and Tianshu Song. 2020. Efficient and Fair Data Valuation for Horizontal Federated Learning. In *Federated Learning - Privacy and Incentive*, Qiang Yang, Lixin Fan, and Han Yu (Eds.). Lecture Notes in Computer Science, Vol. 12500. Springer, 139–152. [https://doi.org/10.1007/978-3-030-63076-8\\_10](https://doi.org/10.1007/978-3-030-63076-8_10)
- [59] Detlof Von Winterfeldt and Gregory W Fischer. 1975. Multi-attribute utility theory: models and assessment procedures. *Utility, probability, and human decision making* (1975), 47–85.
- [60] Hui Yin, Ye-Hwa Chen, and Dejie Yu. 2020. Stackelberg-Theoretic Approach for Performance Improvement in Fuzzy Systems. *IEEE Trans. Cybern.* 50, 5 (2020), 2223–2236. <https://doi.org/10.1109/TCYB.2018.2883729>
- [61] Jinsung Yoon, Serkan Arik, and Tomas Pfister. 2020. Data valuation using reinforcement learning. In *International Conference on Machine Learning*. PMLR, 10842–10851.
- [62] Jin Zhang and Qian Zhang. 2009. Stackelberg game for utility-based cooperative cognitiveradio networks. In *Proceedings of the 10th ACM International Symposium on Mobile Ad Hoc Networking and Computing, MobiHoc 2009, New Orleans, LA, USA, May 18-21, 2009*, Edward W. Knightly, Carla-Fabiana Chiasserini, and Xiaojun Lin (Eds.). ACM, 23–32. <https://doi.org/10.1145/1530748.1530753>