# Equilibrium Characterization for Data Acquisition Games [*]

**Jinshuo Dong** , **Hadi Elzayn** , **Shahin Jabbari** , **Michael Kearns** and **Zachary Schutzman**

University of Pennsylvania

{jinshuo, hads}@sas.upenn.edu, {jabbari, mkearns, ianzach}@cis.upenn.edu

## Abstract

We study a game between two firms in which each provides a service based on machine learning. The firms are presented with the opportunity to purchase a new corpus of data, which will allow them to potentially improve the quality of their products. The firms can decide whether or not they want to buy the data, as well as which learning model to build with that data. We demonstrate a reduction from this potentially complicated action space to a one-shot, two-action game in which each firm only decides whether or not to buy the data. The game admits several regimes which depend on the relative strength of the two firms at the outset and the price at which the data is being offered. We analyze the game's Nash equilibria in all parameter regimes and demonstrate that, in expectation, the outcome of the game is that the initially stronger firm's market position weakens whereas the initially weaker firm's market position becomes stronger. Finally, we consider the perspective of the users of the service and demonstrate that the expected outcome at equilibrium is not the one which maximizes the welfare of the consumers.

## 1 Introduction

Recent years have seen explosive growth in the domain of digital data-driven services. Search engines, restaurant recommendations, and social media are among the many products we use day-to-day which sit atop modern data analysis and machine learning (ML). In such markets, firms live and die by the quality of their models; thus success in the 'race for data', whether acquired directly from customers or indirectly via acquisition of rival firms or purchasing data corpuses, is crucial. In this work, we study two questions: whether such markets tend towards monopoly, and how competition affects consumer welfare. Importantly, we consider these questions in light of the modeling choices that firms must make.

In our model, two firms compete for market share (utility) by providing identical services that each rely on an ML

model. The firms' error rates depend on their choices of algorithms, models and the volume of available training data. Each firm's market share is proportional to the error of its model relative to the model built by its competitor. This is motivated by the observation that the services built using ML are highly accurate, so users are more conscientious of the mistakes the service makes, rather than the successes. A competition exponent measures relative ferocity of competition and maps to a plausible Markov model of consumer choice. See Section 2.2 for more details.

The firms initially possess (possibly differing) quantities of data, and are given the opportunity to buy additional data at a fixed price to improve their models. Since data is costly and *relative* (rather than absolute) model quality determines market share, each firm's best course of action may depend on the actions of its rival. Hence, each firm acts strategically and faces two decisions: whether to buy the additional data, and what type of model to build in order to produce the best product given the data it ends up with.

The decision of what model to build seems to complicate the firms' action space greatly; there is a very large set of model classes to select from, and different classes have different efficiencies. For example, when restricting attention to neural networks, the choices of depth and number of nodes per layer produce different hypothesis classes with different optimal models. Thus, in principle, the decisions of what model class to select and whether to purchase additional data must be made jointly. However, learning theory allows us to greatly reduce this large action space. In Section 2.1, we show that the game in which firms jointly choose a model and whether to attempt to buy the additional data reduces to a strategically equivalent game in which firms first choose whether to buy the data and then choose optimal models.

In Section 3, we characterize the Nash equilibria of our game for different parameter regimes. For no combination of parameters does exactly one firm wish to buy the data; unsurprisingly, for very high prices, neither firm buys data, and for very low prices, both firms do. In the middling regime, the competitive aspect of the game imposes a 'prisoners' dilemma'-like flavor: both firms would prefer neither firm buy the data, but each do so in order to prevent the other from strengthening its position. Moreover, the unique mixed strategy Nash equilibrium in this regime involves firms *increasing* their probability of buying data as price *increases*.

---

This counterintuitive result follows from the logic of equilibrium: firms playing mixed strategies must be indifferent to buying and not buying the data, and as the price rises, the probability that a competitor acquires the data must rise in order to make investing in data acquisition a palatable option.

Finally, we study whether any of the dynamics of the game push the market towards a monopoly. Perhaps counter to a 'rich-get-richer' feedback loop that might be expected in data races, we observe that in all equilibria, the data gap (and thus, market share gap) always narrows (in expectation). As measured by consumer welfare, this is actually *undesirable*. Both the direction of the data gap as well as the welfare implication may be counterintuitive, particularly with respect to the well-known stylized fact that market concentration is bad for consumers. However, consumer data that improves a service can be viewed as exhibiting a form of network effects, in which case perfect competition can result in inefficiency and under-provisioning of a good [Katz and Shapiro, 1985]. In other words, a greater data gap would result in more consumers using a less error-prone service. As for the data race, anecdotal evidence, such as GM's acquisition of automated driving startup Cruise, despite Waymo's earlier market entry and research head-start, are suggestive (though not conclusive) that these predictions may be indicative of real-world dynamics [Primack and Korosec, 2016].

We view our work as a first step towards modeling and analyzing competition for data in markets driven by ML. Under our simplifying assumptions, we derive concrete results with relevance both for policymakers analyzing algorithmic actors as well as engineering or business decision-makers considering the tasks of data acquisition and model selection. Our results are qualitatively robust to other natural modeling choices, such as allowing both firms to purchase the data, as well as treating the data seller as a market participant; however, more significant departures may lead to different conclusions. See Sections 4 and 5 for more details.

## 1.1 Related Work

The theory of ML from a single learner's perspective is well developed, but until recently, little work had studied competition between learning algorithms (see e.g [Ben-Porat and Tennenholtz, 2018; Mansour *et al.*, 2018] for notable exceptions). We differ from both works by exploring the comparative statics and welfare consequences of a single decision (data acquisition). Concurrently, Ben-Porat and Tennenholtz [2019] study a game in which learners strategically choose their model to compete for users, but users only care about the accuracy of predictions on their particular data. In contrast, users in our model choose based on the overall model error.

Our work also intersects with several strains of economic literature, including industrial organization and network effects [David, 1987; Economides, 1996; Katz and Shapiro, 1985]. We differ from such models in two key ways. First, in contrast to assuming a static equilibrium [Katz and Shapiro, 1985] or fixing a dynamic but unchanging process at the outset [Farrell and Saloner, 1986], our work can be viewed as an analysis of a shock to a given potentially asymmetric equilibrium in the form of the availability of new data. Second,

the consumers in our model do not behave strategically (see e.g. [Berry *et al.*, 2017; Mansour *et al.*, 2018]).

Finally, our work is related to spectrum auctions, competition with congestion externalities [Berry *et al.*, 2017], and the sale of information or patents [Kamien and Tauman, 1986; Kamien *et al.*, 1992]. Our results primarily share qualitative similarities: the choice of one firm to buy data (spectrum) forces the other to do so to avoid losing market share, though it would not have been profitable absent the rival, and actual outcomes run counter to consumer preferences (see e.g. [Berry *et al.*, 2017]).

## 2 Framework

We formally motivate and model the ML problem of the firms and demonstrate how this reduces to a game in which the firms can either buy or not buy the new data.

## 2.1 Choosing a Model Class

Consider a firm using ML to build a service e.g. a recommendation system. The amount of data available to the firm is a crucial determinant to the effectiveness of the predictive service of the firm. Fixing the amount of data, the firm faces a fundamental tradeoff; it can use a more complex model that can fit the data better, but learning using a complicated model requires more training data to avoid over- or underfitting.

We can formally represent this tradeoff as follows. Let $\mathcal{H}$ denote the hypothesis class from which the firm is selecting its model and assume the data is generated from a distribution $\mathcal{D}$. Then given $m$ i.i.d. draws from $\mathcal{D}$ the error of the firm when learning a hypothesis from $\mathcal{H}$ can be written as $\mathrm{err}_{\mathcal{D}}(\mathcal{H}) = \mathrm{err}(m, \mathcal{H}) + \min_{h \in \mathcal{H}} \mathrm{err}_{\mathcal{D}}(h)$ [Shalev-Shwartz and Ben-David, 2014].

The first term, known as *estimation error*, determines how well in expectation a model learned with $m$ draws from $\mathcal{D}$ can predict compared to the best model in class $\mathcal{H}$. The second term, known as *approximation error*, determines how well the best model in class $\mathcal{H}$ can fit the data generated from $\mathcal{D}$.

The approximation error is independent of the amount of training data, while the estimation error decreases as the volume of training data increases. The choice of $\mathcal{H}$ affects both errors. In particular, fixing the amount of training data, increasing the complexity of $\mathcal{H}$ will increase the estimation error. On the other hand, the additional complexity will decrease the approximation error as more complicated data generating processes can be fit with more complicated models.

Once the amount of data is fixed, the firm can optimize over its choice of model complexity to achieve the best error. We examine a few widely used ML models and their error forms.

As a first example, consider the case where the firm is building a neural network and has to decide how many nodes $d$ to use. $d$ is the measure of the complexity of the model class and given $m$ data points, the error of the model can be written using the following simplification of a result from Baron [1994].

**Lemma 1** ([Barron, 1994]). *Let $\mathcal{H}$ be the class of neural networks with $d$ nodes. Then for any distribution $\mathcal{D}$, with high probability, the error when using $m$ data points to learn a*

*model from $\mathcal{H}$ is at most $c_1 d/m + c_2/d$, for constants $c_1$ and $c_2$.*

Fixing $m$, the choice of $d$ that minimizes the error can be computed by minimizing the bound in Lemma 1 with respect to $d$. This corresponds to $d = c_2\sqrt{m}/c_1$ and we get that the error of the model built by the firm is $\sqrt{c_1 c_2/m}$.

As another example, consider the very *simple* setting of *realizable PAC learning* where the data points are generated by some hypothesis in a fixed hypothesis class.

**Lemma 2** ([Kearns and Vazirani, 1994]). *Any algorithm for PAC learning a concept class of VC dimension $d$ must use $\Omega(d/\epsilon)$ examples in the worst case.*

Thus in this setting, in the worst case, firms need $\Theta(1/\epsilon)$ training data points to achieve error $\epsilon$. A similar bound gives that with high probability, the firms can guarantee error of $\Theta(1/m)$ (see e.g. [Kearns and Vazirani, 1994]).

In the examples above the error of a firm with $m$ data points takes the form of either $\Theta(m^{-1/2})$ or $\Theta(m^{-1})$ after the firm optimizes over the choice of model complexity. Importantly, the error in both cases (and more generally) degrades as the number of data points increases. The rate at which the error degrades is commonly known as the *learning rate*.

There are other learning tasks with learning rates different than the examples above. Consider a stylized model of a search engine where the set of queries is drawn from a fixed and discrete distribution over a *very large* or even *infinite* set, and the search engine can only correctly answer queries that it has seen before. If, as is often assumed, the query distribution is heavy-tailed, then the search engine will require a large training set to return accurate answers.

In this framework, the probability that a search engine incorrectly answers a query drawn from the distribution is exactly the expectation of the *unobserved* mass of the distribution given the queries observed so far. This quantity is known as the *missing mass* of a distribution (see e.g. [Berend and Kontorovich, 2011; Decrouez *et al.*, 2018; Good, 1953; Orlitsky *et al.*, 2003]). Lemma 3 shows how to bound the expected missing mass for the class of polynomially decaying query distributions.

**Lemma 3** ([Decrouez *et al.*, 2018]). *Let $P^k$ for $k > 1$ be a discrete distribution with polynomial decay defined over $i \in \mathbb{N}_{\geq 0}$ such that $\Pr_{x \sim P^k}[x = i] = i^{-k}/\Sigma_{j=0}^{\infty} j^{-k}$. Then the expected missing mass given $m$ draws from $P^k$ is $\Theta(m^{1/k-1})$.*

By varying $k$ in the query distribution of Lemma 3, the learning rate in the search problem can take the form of $\Theta(m^{-i})$ for any $i \in (0, 1)$. Thus, the learning rate for search may be much faster or slower compared to the previous examples, and the exact rate depends on the value of $k$.

We saw that given a fixed amount of data, a firm using ML can optimize over its learning decisions to get the best possible error guarantee. Furthermore, while error decays as more data becomes available, the rate of decay can vary widely depending on the task. We next see how various learning rates can be incorporated into the parameters of our game.
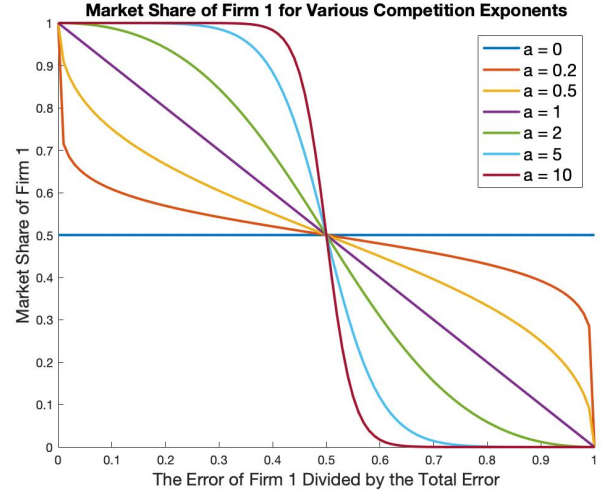


Figure 1: Plot of $f = (1-r)^a/(r^a + (1-r)^a)$ for various $a$ values. When $r = \text{err}_1/(\text{err}_1 + \text{err}_2)$, $f$ is the market share of Firm 1.

## 2.2 Error-Based Market Share

Consider two competing firms (denoted by Firm 1 and 2) that provide identical services e.g. search engines. We assume the market shares of the firms depend on their ability to make accurate predictions e.g. responding to search queries. As discussed above, the quality of their models is determined ultimately by the size of their training data with a task-dependent learning rate. Each firm trains a model on its data and uses its model to provide the service. Let $\text{err}_1$ and $\text{err}_2$ denote the *excess* error of the firms for the corresponding models. Intuitively, these errors measure the quality of the firms' services, so a firm with smaller error should have higher market share. We assume each firm captures a *market share* proportional to the relative errors of the two models. Formally, we define Firm 1 and 2's *error-based* market share as

$$\mu_1 = 1 - \frac{\text{err}_1^a}{\text{err}_1^a + \text{err}_2^a} = \frac{\text{err}_2^a}{\text{err}_1^a + \text{err}_2^a} \text{ and } \mu_2 = 1 - \mu_1. \quad (1)$$

The constant $a \in \mathbb{N}$, which we call the *competition exponent* (following Tullock [2001]), indicates the *ferocity* of the competition, or how strongly a relative difference in the errors of the firms' models translates to a market advantage. As $a$ approaches 0, the tendency is towards each firm capturing half of the market, and thus a large difference in the models' errors is needed for one firm to gain a significant advantage in the market share. Conversely, as $a$ grows larger, even tiny differences in the models' errors translate to massive differences in the market share. (See Figure 1.)

An error-based model reflects markets for services which demand extremely low errors, such as vision systems for self-driving cars. Under the error-based model, if Firm 1 has $99.99\%$ accuracy and Firm 2 has $99\%$ accuracy, Firm 1 will capture $99\%$ of the market share. By contrast, an accuracy-based model (i.e. when the market share of Firm 1 is defined as $\text{acc}_1^a/(\text{acc}_1^a + \text{acc}_2^a)$) would suggest much less realistic near-even split.
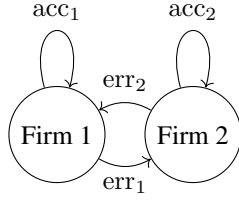
Figure 2: Vertices denote the firms and the directed arrows denote the probability of transition. acc is shorthand for accuracy and err is shorthand for error.

| Firm 1/Firm 2 | Buy (B) | Not Buy (NB) |
|---|---|---|
| Buy (B) | $\frac{1}{2}\left(\mu_1(m_1+n,m_2,b) + \mu_1(m_1,m_2+n,b) - p\right)$ | $\mu_1(m_1+n,m_2,b) - p$ |
| Not Buy (NB) | $\mu_1(m_1,m_2+n,b)$ | $\mu_1(m_1,m_2,b)$ |

Table 1: $u_1(s)$ in all of the strategy profiles of the game.

We provide another justification suggesting that an error-based market can arise even when the learned model is used to provide an everyday service in which high accuracy is not a strict requirement. Consider a customer who, each day, uses the service. She begins by choosing the service of one of the firms uniformly at random. As long as the answers she receives are correct, she has no reason to switch to the other firm's service, and uses the same firm's service tomorrow. However, once the firm makes an error, the customer switches to the other firm's service. The transition probabilities are therefore given by the accuracy and error of each firm. See Figure 2 for the Markov process representing this example.

We can think of the market share captured by each firm as the proportion of the days on which each firm saw the customer. This is exactly the stationary distribution of the associated Markov process as stated in Lemma 4.

**Lemma 4.** *Let $\mu_1$ and $\mu_2$ denote the probability mass that the stationary distribution of the Markov process in Figure 2 assigns to Firms 1 and 2. Then $\mu_1 = \mathrm{err}_2/(\mathrm{err}_2 + \mathrm{err}_1)$, and $\mu_2 = \mathrm{err}_1/(\mathrm{err}_1 + \mathrm{err}_2)$.*

We defer proof of this result (and others) to the full version. By Lemma 4, each firms' market share in the Markov process is exactly the error-based market share defined in Equation 1 for $a = 1$. Similar arguments can motivate error-based market shares for values of $a \in \mathbb{N}$, where the customer switches firms after experiencing $a$ mistakes in a row. The probability of making $a$ mistakes in a row is just $\mathrm{err}_i^a$ for Firm $i$, so the stationary distribution of the Markov process is exactly the two error-based market shares defined in Equation 1.

Using our observations from Section 2.1, we can write the error-based market share in the large-data regime as follows.

**Theorem 1.** *Let $m_1$ and $m_2$ denote the number of data points of Firm 1 and 2, respectively. Then for some $b \in \mathbb{R}^+$, the market share of Firm 1 can be written (asymptotically) as $\mu_1 = m_1^b/(m_1^b + m_2^b)$.*

Because $a$ can be any integer and there exists a corresponding learning problem for any learning rate in $(0, 1]$, Theorem 1 implies that the combined competition exponent in our game can be *any* positive real number, motivated by the initial choice of $a$ and the learning rate of the firms' ML algorithms.

The reductions and derivations in Sections 2.1 and 2.2 allow us to simplify the acquisition games as follows. We first simplify the actions of each firm to only decide whether to buy the data or not, since model choice can be optimized once the number of available data points is known. Moreover,

Theorem 1 not only allows us to simplify the form of market share, but also provides us with a meaningful interpretation for any positive (combined) competition exponent.

### 2.3 The Structure of the Game

Given the reductions so far, we model our game as a two-player, one-shot, simultaneous move game. Firms 1 and 2 begin the game endowed with an existing number of data points, denoted by $m_1$ and $m_2$, respectively. Without loss of generality, we assume $m_1 \geq m_2$. Each firm must decide whether or not to purchase an additional corpus of $n$ data points[1] at a fixed price of $p$. The firm can either Buy (denoted by $B$) or Not Buy (denoted by $NB$) the new data. If both firms attempt to buy the data, the tie is broken uniformly at random (Section 4 discusses relaxing the assumption that only one firm may buy the data). After the purchase, each firm uses its data to train an ML model for its service.

We assume the particular form of the market share of Firm 1 using the reduction in Theorem 1. The market share of Firm 2 is defined to be one minus the market share of Firm 1.

A *strategy profile* $s$ is a pair of strategies, one for each of the firms. Fixing $s$, the utility of Firm $i$ (denoted by $u_i(s)$) is its market share less any expenditure. The utility of Firm 1 in all of the strategy profiles of the game is summarized in Table 1 (rows and columns correspond to the actions of Firm 1 and 2). The utility of Firm 2 is defined symmetrically.

A strategy profile is a *pure strategy Nash equilibrium* (pure equilibrium) if no firm can improve its utility by taking a different action, fixing the other firm's action. A *mixed strategy Nash equilibrium* (mixed equilibrium) is a pair of distributions over the actions (one for each firm) where neither firm can improve its expected utility by using a different distribution over the actions, fixing the other firm's distribution. We are interested in analyzing the Nash equilibria (equilibria).

## 3 Equilibria of the Game

We now turn to finding and analyzing the equilibria. Let

$$A = \frac{(m_1+n)^b}{(m_1+n)^b + m_2^b} - \frac{m_1^b}{m_1^b + (m_2+n)^b},$$

$$C = \frac{(m_1+n)^b}{(m_1+n)^b + m_2^b} - \frac{m_1^b}{m_1^b + m_2^b},$$

$$D = \frac{(m_2+n)^b}{m_1^b + (m_2+n)^b} - \frac{m_2^b}{m_1^b + m_2^b}.$$

These parameters have intuitive interpretations. $A/2$ is the expected change in Firm 1's (or Firm 2's) market share when moving the outcome from $(NB, B)$ (or similarly $(B, NB)$)

---

[1] For simplicity we assume this data is independent of and identically distributed to the data in possession of the firms.

to $(B, B)$. $C$ is the change in market share that Firm 1 receives if it moves from $(NB, NB)$ to $(B, NB)$, and $D$ is the symmetric relation from the perspective of Firm 2. We observe that $A = C + D$. Moreover, since $C$ and $D$ are nonnegative, it is immediately clear that $A > \max\{C, D\}$.

Finally, when $m_1 > m_2$ (i.e. Firm 1 starts with strictly more data), Firm 2 experiences a larger *absolute* change in market share moving from $(NB, NB)$ to $(NB, B)$ than to $(B, NB)$. In other words, Firm 2 experiences a larger *increase* in market share when it buys the data compared to the *decrease* it experiences when Firm 1 receives the data.

**Lemma 5.** *If $m_1 > m_2$, then for all $n$ and $b$, we have $C < D$.*

### 3.1 Characterization of the Equilibria

The equilibria of the game clearly depend on the values of the parameters $m_1$, $m_2$, $n$, $p$ and $b$. For example, if $p > 1$ ($p \leq 0$), then neither firm should ever (not) buy the data. We observe that, fixing the values of $m_1$, $m_2$, $n$ and $b$, there is a range of values for $p$ where the data is *too expensive* (*too cheap*) and $NB$ ($B$) is a dominant strategy for both firms. There is also an intermediate range where more interesting behaviors emerge, as formally characterized in Theorem 2.

**Theorem 2.**

1. *When $p \leq \max\{C, D\}$, $(B, B)$ is the unique equilibrium.*

2. *When $p \geq A$, $(NB, NB)$ is the unique equilibrium.*

3. *When $\max\{C, D\} < p < A$, $(B, B)$ and $(NB, NB)$ are both equilibria. Furthermore, there exists a (unique) mixed equilibrium $((\alpha, 1 - \alpha), (\beta, 1 - \beta))$ such that*

$$\frac{\alpha}{2(1 - \alpha)} = \frac{p - D}{A - p} \text{ and } \frac{\beta}{2(1 - \beta)} = \frac{p - C}{A - p},$$

*where $\alpha$ and $\beta$ denote the probabilities that Firms 1 and 2 select the action $B$, respectively.*

Theorem 2 allows us to make several key observations about the market structure of this game. First, since $C$ and $D$ represent the maximum market share increase achievable by buying the data, the fact that the only equilibrium when $p \in [\min\{C, D\}, \max\{C, D\}]$ is $(B, B)$ means that both firms buy the data despite the fact that the *best-case* market share improvement is less than what they pay. This 'race for data' thus mirrors the prisoner's dilemma – if both firms could agree not to buy the data, they would be better off, but either would be tempted to buy the data and improve market share.

Second, Theorem 2 illustrates how several features of equilibrium depend on the ferocity of competition, as determined by the exponent $a$; as $a$ varies, the frontiers of the regimes described in Theorem 2 shift too. For example, if $a = 0$, market share is split evenly between the two firms, regardless of error or accuracy; unsurprisingly, as $a \to 0$ (which implies $b \to 0$), $A$, $C$, and $D$ also approach 0, so the payoff difference between strategy profiles becomes negligible. As a consequence, regimes (1), (2), and (3) collapse, and all but very small $p$ induce regime (4), where $(NB, NB)$ is the only equilibrium. We observe similar behavior when $a$ is large. Assuming that $m_1 > m_2 + n$, then $a \to \infty$ implies that

$A \to 0$ (and hence $b \to \infty$), again implying that regimes (1), (2), and (3) collapse. Thus again, unless $p$ is very close to 0, $(NB, NB)$ is the unique equilibrium. This is for a different reason than the small $a$ case, however: Firm 1 now has no incentive to buy, since it is guaranteed almost the whole market share using its current model, and Firm 2's initial disadvantage is too great to be overcome by buying the data.

For $a$ between these two extremes, many choices of $m_1$ and $m_2$ lead to a non-empty interval $(\max\{C, D\}, A)$, with endpoints far from 0 and 1. When $p$ falls in this interval, regime (2) holds, so a mixed equilibrium exists; we solve for mixed equilibria in Section 3.2. The complete equilibrium characterization for all regimes in Theorem 2 allows us to pin down the optimal fixed price from the seller revenue perspective. However, in full generality, the seller's problem encompasses further possibilities like auction pricing; hence, we defer this calculation to future work. See Section 5 for a discussion.

### 3.2 Mixed Equilibrium and Monotonicity Analysis

Next, we carefully examine the mixed equilibrium and study the relationship between the weights each firm places on each action and the parameters of the game.

Recall that $\alpha$ and $\beta$ in Theorem 2 denote the probability that Firms 1 and 2 purchase the data in the mixed equilibrium. When $m_1 > m_2$, then $\alpha < \beta$ which implies that the smaller firm will succeed more often in purchasing the data in the mixed equilibrium. The relationship of $\alpha$ and $\beta$ with the number of data points $n$ and the price $p$ is as follows.

**Lemma 6.** *Let $((\alpha, 1-\alpha), (\beta, 1-\beta))$ denote the mixed equilibrium in the the regime where $\max\{C, D\} < p < A$. Then $\alpha$ and $\beta$, both increase when $p$ increases or $n$ decreases.*

Lemma 6 may seem counterintuitive, as it implies that as the price $p$ *rises* through the range in which a mixed equilibrium exists, the probability that any of the firms want to buy the data *also increases*. However, once the price $p$ crosses the threshold $A$, the unique equilibrium is the pure strategy $(NB, NB)$. This gives rise to a discontinuity. See Figure 3.

Of course, this all says nothing about the equilibrium utilities for the firms; as long as the equilibrium utilities are not identical, players will naturally have ordinal preferences over the set of equilibria. We analyze these preferences in Lemma 7, which elucidates the discontinuity at $p = A$.

**Lemma 7.** *When $p \in (\max\{C, D\}, A)$, $u_1(NB, NB) \geq u_1(s)$ and $u_2(NB, NB) \geq u_2(s)$ for all strategy profiles $s$. However, $u_1(B, NB) \geq u_1(B, B) \geq u_1(NB, B)$, while $u_2(NB, B) \geq u_2(B, B) \geq u_2(B, NB)$.*

While both firms agree that $(NB, NB)$ is the most preferred outcome, their preferences over other outcomes are discordant. In particular, given that at least one firm will try to buy the data, each firm would prefer itself to be the buyer. If either firm believes the other may try to buy the data, it will put positive weight on the action $B$ in the mixed equilibrium. But for $p > A$, both firms know that it would be irrational for the other to buy, so the unique equilibrium is $(NB, NB)$.

### 3.3 Change in the Market Share
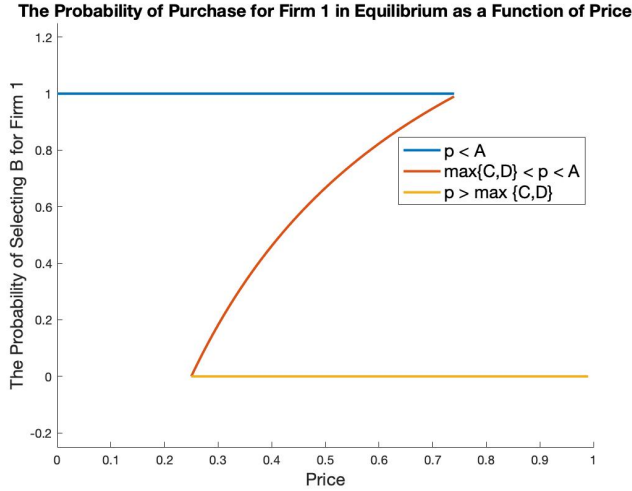
We now analyze the change in the market shares.

Figure 3: Firm 1 selection probability of $B$ by $p$. Blue, red and yellow lines correspond to $(B, B)$, mixed, and $(NB, NB)$ equilibria.

**Lemma 8.** *When $m_1 \geq m_2$, the only strategy profile that strictly increases the market share of Firm 1 is $(B, NB)$.*

So while only $(B, NB)$ leads to an increase in the market share of Firm 1, it is not a pure equilibrium. We show that even when firms play according to the mixed equilibrium, the expected market share of Firm 1 does not strictly increase.

**Theorem 3.** *When $m_1 \geq m_2$ and $p \in (\max\{C, D\}, A)$, the expected market share of Firm 1 does not strictly increase if both firms play according to the mixed equilibrium.*

Together, Lemma 8 and Theorem 3 demonstrate that the natural forces of the interaction on the market are, perhaps surprisingly, antimonopolistic. Since we assume that Firm 1 enters the game with a greater market share than Firm 2, but that no equilibrium allows Firm 1 to increase its market share, the game *disfavors* the concentration of market power. We analyze the implications for consumers below.

### 3.4 Consumer Welfare in Equilibrium

In this section, we show that consumers prefer the outcome $(B, NB)$, in which the initially stronger firm concentrates its market power. This is not supported by a pure equilibrium in any regime, nor is it favored by mixed equilibrium; hence, the interests of the firms do not align with the interests of the consumers. We define the consumer welfare as follows.

**Definition 1.** *Let $m_1(s)$ and $m_2(s)$ denote the (expected) number of data points that Firm 1 and 2 posses when playing according to strategy profile $s$. Then the* consumer welfare *is*
$$CW(s) = \mu_1(s)(1 - \mathrm{err}_1(m_1(s))) + \mu_2(s)(1 - \mathrm{err}_2(m_2(s))).$$

The welfare definition arises from assuming consumers receive 1 unit of utility for accurate predictions and 0 for erroneous ones. Notice that maximizing this definition of consumer welfare is exactly equivalent to minimizing the market-share weighted error probability. This leads to Theorem 4.

**Theorem 4.** *Suppose $m_1 > m_2$. Then the consumers have the following preferences over the strategy profiles.*
$$CW(B, NB) > CW(B, B) > CW(NB, B) > CW(NB, NB).$$

Note that consumers' preference for the outcome in which Firm 1 concentrates its market power is *not* the same as saying that the consumers prefer a monopoly. Rather, the consumers prefer higher quality services. When $m_1 > m_2$, Firm 1's model before acquiring the data has a lower error rate than that of Firm 2, and so, of all the possible outcomes, the one which leads to a product with the lowest error rate is the one in which Firm 1 is able to improve on its already superior product. But if Firm 2 were not a player at all, then a monopolistic Firm 1 would *have no incentive to* buy the new data. Therefore, a monopoly without the threat of competition will not lead to the best outcome from the consumer's perspectives.

## 4 Extensions and Robustness

Next, we consider robustness to two simple extensions.

**Firm Acquisition** We treat the data seller as a market participant with its own customers and market share. This allows us to model firm acquisition: buying the data translates to acquiring the firm and its customers, and neither firm buying the data corresponds to the third firm remaining in the market.

**Simultaneous Sale** Rather than the data being exclusively sold to one firm in the case that both firms buy, we allow the seller to sell the data to both firms at the same fixed price.

In both of these extensions, we can again derive the quantities $A$, $C$, and $D$; while the precise quantities change, their rankings and relationships do not. Thus the general phenomenon of three regimes, with mixing over the middle regime, remains unchanged. Moreover, in expectation, the market share becomes less asymmetric in both extensions.

## 5 Future Directions

We view our work as a first step towards modeling and analyzing competition for data in markets driven by ML. There are several directions for further investigation. First, we modeled the data to be acquired as having a fixed size and a fixed price, but real datasets are divisible. Second, we can consider a strategy space expanded to include buying any number of data points at a fixed price *per data point* or nonlinear function of the number of data points purchased. More generally, treating the seller of the data as an additional player in the game allows for further questions, such as: what is the optimal revenue-generating mechanism to sell the data? And does the optimal mechanism maximize social welfare?

Additionally, many firms that provide learning-based services acquire their data through their customers that use the service. In this way, capturing a larger market share induces a feedback loop which allows a firm to iteratively improve its product. What can be said about our game in a repeated setting with dynamic feedback effects? Furthermore, firms that provide digital services often operate in a secondary market in which other firms pay for advertising spots in their product. Improving one's market share should in principle allow a firm to charge advertisers a higher price, but we do not know to what extent this affects the analysis of the equilibria of the game. Incorporating advertiser behavior would greatly complicate the model but provide potentially interesting results.

# References

[Barron, 1994] Andrew Barron. Approximation and estimation bounds for artificial neural networks. *Machine Learning*, 14(1):115–133, 1994.

[Ben-Porat and Tennenholtz, 2018] Omer Ben-Porat and Moshe Tennenholtz. A game-theoretic approach to recommendation systems with strategic content providers. In *Proceedings of the 32nd Annual Conferenceon Neural Information Processing Systems*, pages 1118–1128, 2018.

[Ben-Porat and Tennenholtz, 2019] Omer Ben-Porat and Moshe Tennenholtz. Regression equilibrium. In *Proceedings of the 2019 ACM Conference on Economics and Computation*, 2019.

[Berend and Kontorovich, 2011] Daniel Berend and Aryeh Kontorovich. The missing mass problem. *CoRR*, abs/1111.2328, 2011.

[Berry et al., 2017] Randall Berry, Michael Honig, Thanh Nguyen, Vijay Subramanian, and Rakesh Vohra. The value of sharing intermittent spectrum. *CoRR*, abs/1704.06828, 2017.

[David, 1987] Paul David. Some new standards for the economics of standardization in the information age. *Economic Policy and Technological Performance*, pages 206–239, 1987.

[Decrouez et al., 2018] Geoffrey Decrouez, Michael Grabchak, and Quentin Paris. Finite sample properties of the mean occupancy counts and probabilities. *Bernoulli*, 24(3):1910–1941, 2018.

[Economides, 1996] Nicholas Economides. The economics of networks. *International Journal of Industrial Organization*, 14(6):673–699, 1996.

[Farrell and Saloner, 1986] Joseph Farrell and Garth Saloner. Installed base and compatibility: Innovation, product pre-announcements, and predation. *The American Economic Review*, 76(5):940–955, 1986.

[Good, 1953] Irving Good. The population frequencies of the species and the estimation of population parameters. *Biometrika*, 40(3-4):237–264, 1953.

[Kamien and Tauman, 1986] Morton Kamien and Yair Tauman. Fees versus royalties and the private value of a patent. *The Quarterly Journal of Economics*, 101(3):471–491, 1986.

[Kamien et al., 1992] Morton Kamien, Shmuel Oren, and Yair Tauman. Optimal licensing of cost-reducing innovation. *Journal of Mathematical Economics*, 21(5):483–508, 1992.

[Katz and Shapiro, 1985] Michael Katz and Carl Shapiro. Network externalities, competition, and compatibility. *The American Economic Review*, 75(3):424–440, 1985.

[Kearns and Vazirani, 1994] Michael Kearns and Umesh Vazirani. *An Introduction to Computational Learning Theory*. MIT Press, 1994.

[Mansour et al., 2018] Yishay Mansour, Aleksandrs Slivkins, and Zhiwei Steven Wu. Competing bandits: Learning under competition. In *Proceedings of the 9th Innovations in Theoretical Computer Science Conference*, pages 48:1–48:27, 2018.

[Orlitsky et al., 2003] Alon Orlitsky, Narayana Santhanam, and Junan Zhang. Always good turing: Asymptotically optimal probability estimation. In *Proceedings of the 44th Symposium on Foundations of Computer Science*, pages 179–188, 2003.

[Primack and Korosec, 2016] Dan Primack and Kirsten Korosec. GM buying self-driving tech startup for more than $1 billion. *Fortune*, 2016.

[Shalev-Shwartz and Ben-David, 2014] Shai Shalev-Shwartz and Shai Ben-David. *Understanding Machine Learning: From Theory to Algorithms*. Cambridge University Press, 2014.

[Tullock, 2001] Gordon Tullock. *Efficient Rent Seeking*, pages 3–16. Springer, 2001.