

Profit equitably: An investigation of market maker's impact on equitable outcomes

Kshama Dwarakanath, Svitlana S Vyetrenko and Tucker Balch
 {kshama.dwarakanath,svitlana.s.vyetrenko}@jpmchase.com
 JP Morgan AI Research, USA

ABSTRACT

We look at discovering the impact of market microstructure on equitability for market participants at public exchanges such as the New York Stock Exchange or NASDAQ. Are these environments equitable venues for low-frequency participants (such as retail investors)? In particular, can market makers contribute to equitability for these agents? We use a simulator to assess the effect a market maker can have on equality of outcomes for consumer or retail traders by adjusting its parameters. Upon numerically quantifying market equitability by the entropy of the price returns distribution of consumer agents, we demonstrate that market makers indeed support equitability and that a negative correlation is observed between the profits of the market maker and equitability. We then use multi objective reinforcement learning to concurrently optimize for the two objectives of consumer agent equitability and market maker profitability, which leads us to learn policies that facilitate lower market volatility and tighter spreads for comparable profit levels.

KEYWORDS

Agent based simulations, equitability, multi objective reinforcement learning, volatility

1 BACKGROUND AND RELATED WORK

1.1 Introduction

An estimated 40% to 60% of trading activity across all financial markets (stocks, derivatives, liquid foreign currencies) is due to high frequency trading, which is when trading happens electronically at ultra high speeds [16]. High frequency trading and its implications on market stability and equitability have been in regulatory spotlight after the flash crash of May 6, 2010, when the Dow Jones Industrial Average dropped 9% within minutes [15]. To dissect the impact of high frequency trading, it is necessary to distinguish between intentional manipulation (e.g. front-running, spoofing) and legitimate trading strategies such as market making, pairs trading and statistical arbitrage, new reaction strategies, etc [1]. Designated market makers (denoted MM henceforth), for example, are obliged to continuously provide liquidity to both buyers

and sellers regardless of market conditions - contributing to market stability by providing counterparty to all transactions even in times of market distress. For example, NASDAQ requires MMs "to maintain a continuous two-sided trading interest during regular market hours, at prices within certain parameters expressed as a percentage referenced from the National Best Bid or Offer" [27]. MMs are typically rewarded with lower exchange fees for doing so [11].

MM presence and facilitation of transactions among market participants is associated with better market quality in terms of high frequency of auction clearance, and low variability in returns and trading volume [34]. Classical approaches to MM include modeling the order arrivals and executions, and using control algorithms to find an optimal strategy for the MM [8]. Reinforcement learning (RL) has also been used to learn market making strategies using a replay of historical data to derive MM's policies [9, 30, 31].

Replaying historical data in MM agent training does not take into account trading response from the market environment. In this paper, we study how to learn a MM strategy that optimizes for both equitability and profitability in a responsive agent-based simulated environment. Specifically, we examine market equitability in scenarios where an elementary MM provides liquidity in an equity market including a large number of consumer or retail investors. The MM is defined by a set of parameters that govern the spread and depth of orders it places into the order book. Based on observations of the relationship between the MM profits and market equitability with variations in the MM parameters, we formulate and solve an RL problem to find a MM policy that optimizes for both market equitability as well as MM profits. The contributions of this paper are as follows

- (1) Survey of equitability metrics that are relevant to the domain of finance and equity markets, with the contribution of a new entropy metric that measures individual fairness.
- (2) Formulation of the problem of learning an equitable MM policy as a multi-objective RL problem, and proposal of an effective exploration scheme for training a MM agent.
- (3) Demonstrated that the addition of an equitability reward in learning a policy for the MM helps improve the trade-off between equitability and profitability.
- (4) Detailed analysis of the policy learnt using RL by varying the importance given to equitability and inventory objectives.

1.2 Limit Order Books

More than half of the markets in the financial world today use a limit order book (LOB) mechanism to facilitate trade [14, 25]. An LOB is the set of all orders in the market at the current time, with each order represented by a direction of trade - buy or sell (equivalently bid or ask), price and order size. There are two main order types

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.
 ICAIF'21, November 3–5, 2021, Virtual Event, USA

of interest - market orders and limit orders. A *market order* is a request to buy or sell instantaneously at the current market price. On the other hand, a *limit order* is a request to buy or sell at a price better than or equal to its specified price. Therefore, limit orders face the risk of not being executed instantaneously and joins other limit orders in the queue at its price level. There are a fixed number of price levels in the LOB, with the gap between subsequent prices called the *tick size*. The arithmetic mean of the best bid and best ask prices in the LOB is called the *mid-price* for the stock. And, the difference between the best ask and best bid is called the *spread*. An example snapshot of an LOB is shown in Figure 1.

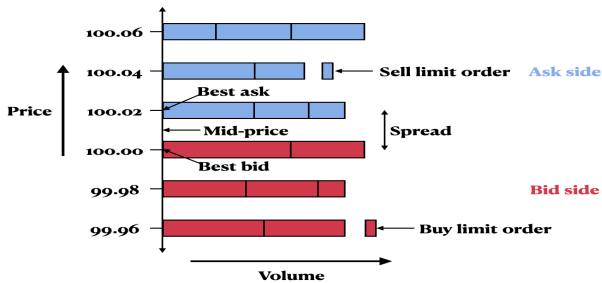


Figure 1: Snapshot of a limit order book

1.3 Simulator

In order to evaluate our MM’s impact on equitability, we employ a multi-agent LOB simulator [citation redacted]. It provides a selection of background trading agent types described in section 1.3.1 and a NASDAQ-like exchange agent which lists securities for trade against an LOB with price-then-FIFO matching rules. It is also equipped with a simulation kernel that manages the flow of time and handles all inter-agent communication. Trading agents may not inspect the state of the exchange directly, but can send messages to request order book depth, obtain last trade prices, or place or cancel limit orders through the kernel, which imposes delays for computational effort and communication latency. Time proceeds in nanoseconds and all securities are priced in cents. This discrete event simulation mechanism permits rapid simulation despite the fine time resolution, as periods of inactivity can be “skipped over” without computational effort.

1.3.1 Background trading agents. The simulator includes agents with different trading behaviors and incentives.

Value Agents: The value agents are designed to simulate fundamental traders that trade in line with their belief of the exogenous stock value (which we call *fundamental price*), without any view of the LOB microstructure [17]. The fundamental price represents the agent’s understanding of the outside world (e.g. earnings reports, macroeconomic events, etc) [36], [35]. In this paper, we model the fundamental price of an asset by its historical mid-price series. Note however that significant deviations between the fundamental and the simulated mid-price are possible since agent interactions ultimately define the simulated mid-price. Each value agent arrives to the market multiple times in a trading day according to a Poisson process, and chooses to buy or sell a stock depending on whether

it is cheap or expensive relative to its noisy observation of the fundamental. Upon determining the side of its order, the value agent places a limit order at a random level either inside the spread or deeper into the LOB. Therefore, value agents assist price formation in the LOB by bringing in external information.

MM Agent: MMs play a central role as liquidity providers by continuously quoting prices on both sides of the LOB, and earn the spread if orders execute on both sides. MMs act as intermediaries and essentially eliminate *air pockets* between existing buyers and sellers. They also make markets more liquid and enable investors to trade large quantities with smaller price moves [1]. In this work, we define a MM by the *stylized parameters* that follow from its regulatory definition. While most realistic MM models have adverse selection and/or inventory control mechanisms that are intended to boost the MM’s profitability (e.g., [2]), we however highlight that our definition is model-independent and is intended to explore the effect of its stylized properties on market equitability. Our MM definition is similar in spirit to that of [35] and [8], barring the fact that our MM does not use any reference price series to determine its mid-price. It instead adapts to the mid-price in the LOB.

Momentum Agents: The momentum agents follow a simple momentum strategy, waking up at a fixed rate and observing the mid-price each time. They compare a long-term average of the mid-price with a short-term average. If the short-term average is higher than the long-term average, the agent buys since the price is seen to be rising. On the other hand, if the short-term average is lower than the long-term average, the agent sells since it believes that the price is falling.

Consumer Agents: Consumer agents are designed to emulate consumer agents who trade on demand without any other considerations ([17]). Each consumer agent trades once a day by placing a market order of a random size in a random direction. Consumer agents arrive to the market at times that are uniformly distributed throughout the trading day. Our focus here is to numerically estimate the equitability of outcomes of consumer agents in the simulated environment described above.

1.4 Equitability and metrics

Equitability has conventionally been studied in political philosophy and ethics [21], economics [19], and public resource distribution [38]. In recent times, there has been a renewed interest in quantifying equitability in classification tasks [10]. In finance, risk is synonymous with volatility since volatile markets lead to some investors doing well with others doing poorly [26]. For instance, if a consumer agent places a trade during a period of market instability when asset price moves sharply in the direction opposite to the trade and then rebounds in seconds (called mini-flash crash), the agent’s execution may be perceived to be inequitable as compared to other agents who were luckier to trade in more stable markets. Accordingly, regulators are interested in reducing market volatility to facilitate equitability from the perspective of equality of market outcomes.

For this work, we are interested in the notion of individual fairness which advocates the sentiment that *any two individuals who are similar with respect to a task should be treated similarly*. We now collect some metrics for individual fairness taking inspiration

from those used to quantify income inequality in populations [29], and those from information theory. The Theil index measures the *distance* a population is away from the ideal egalitarian state of everyone having the same income [32]. It belongs to the family of generalized entropy indices (GEI) that satisfy the property of subgroup-decomposability, i.e. the inequality measure over an entire population can be decomposed into the sum of a between-group unfairness component (similar to group fairness metrics in [10], [5]) and a within-group unfairness component. The Gini index is another popular metric of income inequality that measures how far a country's income distribution deviates from a totally equal distribution [12].

Another metric is based on information entropy, which is a measure of uncertainty of a random variable. Entropy has been commonly used as a metric for equitable resource allocation in wireless networks [18], a measure of income inequality and an indicator of population diversity in the social sciences [23], [33], [4], [3]. Since we have finitely many samples y_1, \dots, y_n from the distribution of individual outcomes, we estimate the entropy of individual outcomes by binning them into K bins. Subsequently, we compute the entropy estimate $H_K(Y) = -\sum_{k=1}^K p_k \log p_k$ where p_k denotes the empirical frequency of outcomes in the k^{th} bin. Formally, the equitability metric based on estimating information entropy from samples y_1, \dots, y_n is given by

$$\mathcal{L}_{\text{Ent}}(y_1, \dots, y_n) = 1 - \frac{H_K(Y)}{\log K} = 1 - \frac{-\sum_{k=1}^K p_k \log p_k}{\log K} \quad (1)$$

Note that the scaling factor of $\log K$ corresponds to the entropy of a uniform distribution over the outcomes, and is most inequitable.

1.5 Motivation

Per its regulatory definition, the MM acts as a liquidity provider by placing limit orders on both sides of the LOB with a constant arrival rate. At time t , it starts by cancelling any of its unexecuted orders. Let the levels of the LOB on both bid and ask sides be indexed from 0 to N , with level 0 corresponding to the innermost LOB levels. Let b_t and a_t denote the best bid and best ask respectively. The MM looks at the mid-price $m_t := \frac{a_t + b_t}{2}$, and places new price quotes of constant size K at d price increments around m_t . That is, it places bids at prices $m_t - s_t - d, \dots, m_t - s_t$ and asks at prices $m_t + s_t, \dots, m_t + s_t + d$, where d is the *depth* of placement and s_t is the *half-spread* chosen by the MM at time t . As in NASDAQ, the difference between consecutive LOB levels is one cent. Such a stylised market maker has been shown to be profitable over a sufficiently long time horizon, when the mid-price follows an Ornstein-Uhlenbeck Process in [8].

We consider two versions of our MM agent: one that posts liquidity at a constant half-spread s_t for all times t ; and one that adapts its half-spread to its observation of the LOB at time t , i.e. with $s_t = \frac{a_t - b_t}{2}$. We are interested in examining the effect of the MM parameters of half-spread and depth on equitability for consumer agents. To do so, the half-spread s_t is varied across values in $\{1, 2, \frac{a_t - b_t}{2}\}$ cents, and the depth is varied across values in $\{1, 2\}$ cents. For each (half-spread, depth) pair, we collect 20 samples of the resulting MM profit and equitability as defined by the information entropy metric in equation (1). The green dots in figure

11 are the average values of MM profits and equitability for the above range of spread and depth parameters. We observe a negative correlation between the two variables of interest, as represented by green line fit to the green dots. This discovery inspired us to think of a way to improve MM profits while maintaining a high level of equitability to consumer agents. We formulate this as an RL problem directed at optimizing for both MM profits and market equitability to consumer agents, using our multi agent simulator to play out interactions between traders in the market.

2 PROBLEM SETUP

Consider a market configuration comprising a group of consumer agents trading alongside more intelligent investors, and a MM that provides liquidity. We are interested in using RL to improve the profits made by the MM, while guaranteeing a high degree of equitability to consumer agents. We formulate this as an RL problem by representing the market with different trading agents (including the MM) as a Markov Decision Process (MDP). An MDP is a tuple $(\mathcal{S}, \mathcal{A}, \mathcal{T}, R, \gamma, T)$ comprising the state and action spaces along with a model of the environment (either known or unknown), the reward function, discount factor and the time horizon for the planning/control problem. The state for our MDP captured the states of both the market and the MM as

$$[\text{inventory} \quad \text{imbalance} \quad \text{spread} \quad \text{midprice}] \quad (2)$$

where *inventory* is number of shares held in the MM's inventory, *imbalance* = $\frac{\text{total bid volume}}{\text{total bid volume} + \text{total ask volume}}$ is the volume imbalance in the order book, *spread* is the current market spread, and *midprice* is the current mid-price of the stock. The actions include trading actions of the MM

$$[s_t \quad d]$$

where s_t is the half-spread and d is the depth of orders placed by the MM as described in section 1.5. In order to quantify the reward function, we need a way to numerically quantify equitability alongside MM profits.

Quantifying equitability essentially involves determining if the consumer agents are treated *similarly* to one another, where *similarity* is ascertained with respect to the task at hand. In this paper, we propose a method for quantifying equitability in consumer agent outcomes by looking at the distribution of differences between an agent's traded price and asset price T seconds afterwards. The only incentive consumer agents have to trade is demand, and outcomes should not be contingent on the time when their trades are placed. Formally, let a consumer agent i execute a trade at price p_t at time t . Define the T -period return $r_{i,T}$ of consumer agent i for some $T \geq t$ and future price p_T as

$$r_{i,T} = \begin{cases} p_T - p_t & \text{if agent } i \text{ executes a buy order} \\ p_t - p_T & \text{if agent } i \text{ executes a sell order} \end{cases} \quad (3)$$

Any of the metrics described in 1.4 could be used to quantify individual fairness based on the aforementioned price returns, but the information entropy metric (1) was seen to be more preferable than others. This is because the GEI and Theil index are undefined (for non-even values of the GEI parameter α) when the returns are negative, with the GEI giving spiky values hindering learning for $\alpha = 2$. Since the consumer agents arrive at random times during

a trading day, their returns are accumulated and the equitability reward is computed as the change in information entropy every $N > 1$ time steps as

$$\begin{aligned}
 R_{\text{Equitability}}^t &= \begin{cases} \mathcal{L}_{\text{Ent}}(y_1, \dots, y_{mN}) - \mathcal{L}_{\text{Ent}}(y_1, \dots, y_{(m-1)N}) & \text{if } t = mN + 1 \\ 0 & \text{otherwise} \end{cases} \\
 &= \begin{cases} \frac{H_K(y_1, \dots, y_{(m-1)N})}{\log K} - \frac{H_K(y_1, \dots, y_{mN})}{\log K} & \text{if } t = mN + 1 \\ 0 & \text{otherwise} \end{cases} \quad (4)
 \end{aligned}$$

where $m \in \{1, 2, \dots\}$ with $H(y_0) := 0$. The need for using a change in the entropy metric is justified by looking at the objective of the RL problem over the entire trading day. Before doing that, we quantify the reward for profits of the MM.

The reward from profits and losses (PnL) of the MM comprise a weighted combination of those from holding a certain inventory, and those from making the latest trade

$$R_{\text{PnL}} = \bar{\eta} \cdot \text{inventory PnL} + \text{matched PnL} \quad (5)$$

where $\text{inventory PnL} = \text{inventory} \times \text{change in mid-price}$ and matched PnL refers to profits made by capturing the spread. The weighting factor for the inventory PnL $\bar{\eta}$ (called *inventory-weight*) is chosen to be in $(0, 1)$ due to the ensuing observation. Our MM is intended to make profits by capturing the spread upon executing orders on both sides of the LOB. But, there is always a risk of collecting inventory from one (side) of the orders not being executed, and the inventory PnL incentivizing the MM to exploit price trends. By weighing down the inventory component of PnL, we focus on learning a MM that profits from capturing the spread while also enabling a more stable learning setup as in [30].

Since we have two (potentially competing) objectives of profits and equitability, this becomes a multi-objective RL (MORL) problem similar to those encountered in robotics [20], economic systems [28] and natural resource allocation [6]. The field of MORL involves learning to act in environments with vector reward functions $\mathbf{R} : \mathcal{S} \rightarrow \mathbb{R}^n$. Since the value functions in MORL induce only a partial ordering even for a given state [24], additional information on the relative importance of the multiple objectives is given in the form of a scalarization function. To understand the relationship between profits and equitability in this paper, we found it sufficient to use a linear scalarization function. Thus, the combined reward function for our RL problem is given by

$$R = R_{\text{PnL}} + \eta \cdot R_{\text{Equitability}} \quad (6)$$

where $\eta \geq 0$ is called the *equitability-weight*, and has the interpretation of monetary benefits in \$ per unit of equitability. (6) is also motivated from a constrained optimization perspective as being the objective in the unconstrained relaxation [7] of the problem $\max R_{\text{PnL}}$ s.t. $R_{\text{Equitability}} \geq c$. Usage of a linear scalarization function also helps us analyse the variation in equitability and profits of the learnt policy, as η is varied.

With the rewards defined in equations (6), (5) and (4), discount factor $\gamma = 1$ and the MM starting from the same initial state every trading day (i.e. same amount of money and shares), the RL objective

under a policy π evaluates to

$$\begin{aligned}
 \bar{\eta} \cdot \mathbb{E} \left[\sum_{t=0}^{T-1} \text{inventory PnL} \right] + \eta \cdot \mathbb{E} \left[\frac{-H_K(y_1, \dots, y_{L-2})}{\log K} \right] \\
 + \mathbb{E} \left[\sum_{t=0}^{T-1} \text{matched PnL} \right] \quad (7)
 \end{aligned}$$

Hence, the objective of the equitable MM is to learn a trading policy that maximizes its PnL over a given time horizon, while minimizing the entropy of consumer agent returns.

Since it is straightforward to compute the optimal policy given the state-action value function or Q function (as the argument that maximizes the Q value in each state), we use the Q learning algorithm that iteratively updates estimates of the optimal Q function [37]. The Q function is defined under a policy π as the long-term value of taking action a in state s and following policy π thereafter

$$\begin{aligned}
 Q^\pi(s, a) = \mathbb{E} \left[\sum_{t=0}^{T-1} \gamma^t R(s_{t+1}) \middle| s_{t+1} \sim \mathcal{T}(\cdot | s_t, a_t), a_t \sim \pi(\cdot | s_t), \right. \\
 \left. s_0 = s, a_0 = a \right] \quad (8)
 \end{aligned}$$

Since we consider a finite set of (discretized) states and actions, the Q function is a table of values for all possible (state,action) pairs. A useful property of learning such a tabular Q function is that the resulting policy is more easily interpretable as opposed to using function approximators such as neural networks.

The Q learning algorithm uses a stochastic approximation method to estimate the optimal Q function from the Bellman equation [13, 22]. It updates the Q values in the n^{th} episode as

$$\begin{aligned}
 Q_n(s, a) &= \begin{cases} (1 - \alpha_n) Q_{n-1}(s, a) + \alpha_n (R(s') + \gamma V_{n-1}(s')) & \text{if } (s_n, a_n) = (s, a) \\ Q_{n-1}(s, a) & \text{otherwise} \end{cases} \quad (9)
 \end{aligned}$$

where s' is the next state when action a is taken in state s , α_n is the learning rate in the n^{th} episode, and $V_{n-1}(s') = \max_{a'} Q_{n-1}(s', a')$. The iterates $Q_n(\cdot, \cdot)$ are guaranteed to converge to the optimal Q values if the rewards are bounded, the learning rates satisfy $0 \leq \alpha_n < 1$, and the sequence of learning rates $\{\alpha_n\}$ satisfy $\sum_n \alpha_n = \infty$ and $\sum_n \alpha_n^2 < \infty$. The optimal policy is then given by $\pi^*(s) = \arg \max_a Q^*(s, a) \forall s \in \mathcal{S}$.

3 EXPERIMENTS

3.1 Non learning experiments

3.1.1 Relationship between profits and equitability for different consumer agent order sizes. We saw a negative correlation between MM profits and market equitability upon varying the MM half spread and depth in Figure 11. Another important parameter is the order size of consumer agents, since equitability for bigger (higher order volume) players potentially differs from that for smaller players, even though they are of the same trading type. Therefore, we now vary the parameters of half-spread and depth of MM orders, as well as consumer agent order size in order to examine their impact on market equitability to consumer agents. The half-spread takes values in $\{2, 5, 10, 20, \frac{a_t - b_t}{2}\}$ cents, the depth takes values in

¹The **bold face** represents a vector.

{1, 2, 5, 10, 15} cents, and consumer agent order size takes values in {5, 10, 30, 50, 100}. For each (half-spread,depth,order size) triple, we collect 200 samples of the resulting MM profit and equitability as defined by the information entropy metric in equation (1).

Figure 2 is a plot of MM profits versus equitability for all order sizes along with a common (black) regression line. On examining the scattered points, we see that equitability increases with order size. Figure 3 is a plot of profits versus equitability with subplots representing each order size with regression lines for their corresponding order sizes. The larger orange stars represent the average MM profit and equitability for a fixed (half-spread,depth). We see an improvement in the correlation between profits and equitability with increase in order size. This confirms our intuition that the stylized MM is more equitable to consumer agents that place large orders than those that place smaller ones.

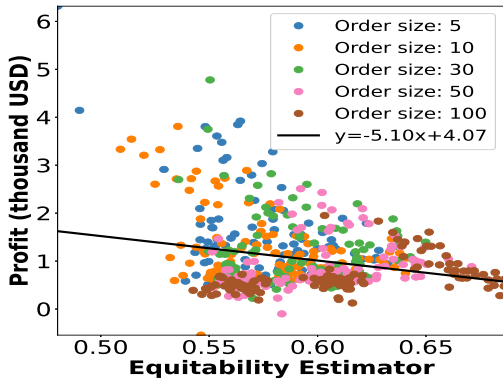


Figure 2: Profits versus equitability for all order sizes

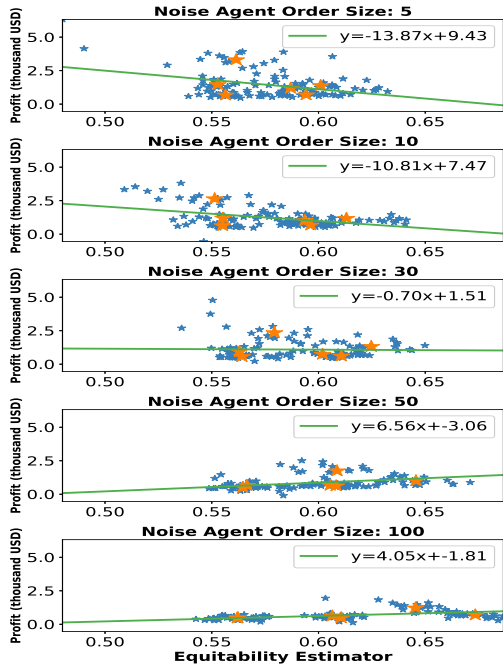


Figure 3: Profits versus equitability for each order size

3.1.2 Relationship between profits and equitability for different sources of liquidity. We note that the consumer agents we consider in this paper are meant to refer to individual traders that place orders purely based on demand. They do not use external information or any intelligence in trading, unlike value agents and momentum agents described previously. Once the value agents bring in fundamental information into the market through their trading strategies, other background agents can observe their actions to infer the fundamental value. Thus, intelligent traders cause information flow in the market from their access to external information. We now look at the influence of external information on market equitability by varying the amount of liquidity provided by value agents as a fraction of that provided by our stylized MM. Figure 4 is a plot of MM profits versus equitability to consumer agents for decreasing amounts of liquidity provided by value agents as a fraction of that provided by our MM. Each sub-figure also has regression lines corresponding to different order sizes for consumer agents. We see that the correlation improves with decrease in value agent liquidity for each consumer agent order size, i.e. a market with more liquidity coming from our stylized MM is more equitable to consumer agents, as opposed to one where more liquidity comes from value agents. Hence, it is more equitable to have less proprietary information flowing into the market since the consumer agents (from the way we have defined them) don’t use that information themselves.

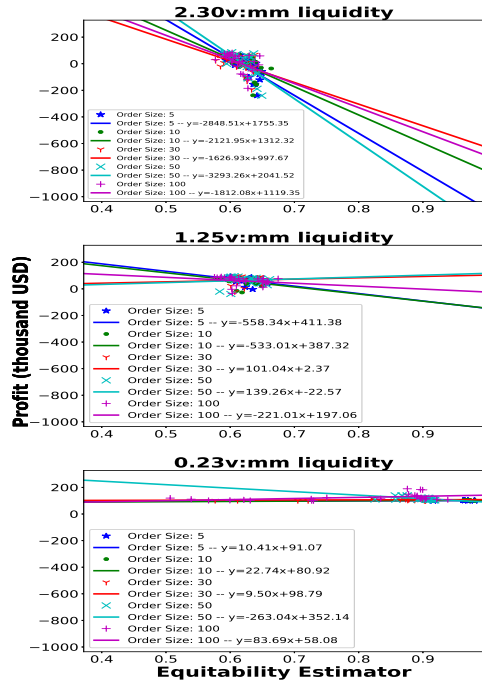


Figure 4: Varying liquidity from value agents

3.2 Learning experiments

We first discretize our states and actions so as to make the problem suited to using RL methods for finite MDPs. An important

requirement for the functioning of the tabular Q learning algorithm is that there be enough visitations of each (state,action) pair. Accordingly, we pick our state and action discretization bins by observing the range of values taken in a sample experiment. The discrete state computed after measuring the state in (2) is $s = [s_{\text{inventory}} \ s_{\text{imbalance}} \ s_{\text{spread}} \ s_{\text{midprice}}]$ where

$$s_{\text{inventory}} = \begin{cases} 0 & \text{if inventory} < -10I \\ 1 & \text{if } -10I \leq \text{inventory} < -5I \\ 2 & \text{if } -5I \leq \text{inventory} < 0 \\ 3 & \text{if } 0 \leq \text{inventory} < 5I \\ 4 & \text{if } 5I \leq \text{inventory} < 10I \\ 5 & \text{if inventory} \geq 10I \end{cases}$$

$$s_{\text{imbalance}} = \begin{cases} 0 & \text{if } 0 \leq \text{imbalance} < 0.25 \\ 1 & \text{if } 0.25 \leq \text{imbalance} < 0.5 \\ 2 & \text{if } 0.5 \leq \text{imbalance} < 0.75 \\ 3 & \text{if } 0.75 \leq \text{imbalance} \leq 1 \end{cases}$$

$$s_{\text{spread}} = \begin{cases} 0 & \text{if } 0 \leq \text{spread} < 2 \\ 1 & \text{otherwise} \end{cases}$$

$$s_{\text{midprice}} = \begin{cases} 0 & \text{if midprice} < m_0 \\ 1 & \text{if midprice} \geq m_0 \end{cases}$$

where I and m_0 are constants. The discrete actions taken by the MM are of the form $a = [a_{\text{mm-spread}} \ a_{\text{mm-depth}}]$ to encode the spread and depth of orders placed by the MM as

$$a_{\text{mm-spread}} = \begin{cases} 0 & \leftrightarrow \text{mm-spread} = \text{current spread} \\ 1 & \leftrightarrow \text{mm-spread} = 1 \\ 2 & \leftrightarrow \text{mm-spread} = 2 \end{cases}$$

$$a_{\text{mm-depth}} = \begin{cases} 0 & \leftrightarrow \text{mm-depth} = 1 \\ 1 & \leftrightarrow \text{mm-depth} = 2 \end{cases}$$

The discretization bins for consumer agent returns are chosen based on their observed values in a sample experiment with $K = 12$ as $(-\infty, -10^5) \cup [-10^5, -10^4) \dots \cup [10^4, 10^5) \cup [10^5, \infty)$ with intervals corresponding to bins for computing empirical probabilities.

3.2.1 Training and convergence. The Q learning algorithm (9) is used to compute estimates of the optimal Q function for the MDP described above. In order to balance exploration and exploitation, we use an ϵ -greedy approach to choose the next action in a given state. Hence, in episode n , the next action is chosen to be that which maximizes the current Q function estimate Q_n with a probability of $1 - \epsilon_n$, and is chosen randomly with probability ϵ_n . We call ϵ_n the exploration parameter that is decayed as episodes progress. An important assumption underlying the convergence results for the tabular Q learning algorithm is that we have *adequate* visitation of each (state, action) pair. We found that having a dedicated exploration phase where the exploration parameter ϵ and learning rate α are kept constant at high values helps us in obtaining these visitations.

Training is divided into three phases - pure exploration, pure exploitation and convergence phases. During the pure exploration phase, α_n and ϵ_n are both held constant at high values in order to facilitate the visitation of as many state-action discretization bins

Number of pure exploration episodes	400
Number of pure exploitation episodes	200
Number of convergence episodes	400
Total number of training episodes	1000
γ	1.0
T	5040
$\alpha_0 = \alpha_1 = \dots = \alpha_{399} = \dots = \alpha_{599}$	0.9
α_{999}	10^{-5}
$\epsilon_0 = \epsilon_1 = \dots = \epsilon_{399}$	0.9
ϵ_{599}	0.1
ϵ_{999}	10^{-5}
η	{0, 0.05, 0.1, 0.15, 0.2, 0.3}
$\bar{\eta}$	{0, 3, 6, 10, 20, 50}
I	500
m_0	10^5
K	12

Table 1: Numerics for experiments

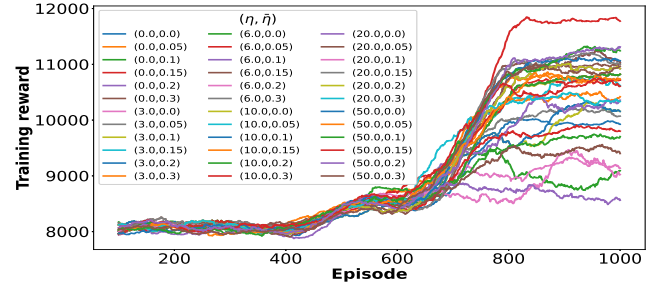


Figure 5: Training rewards for $(\eta, \bar{\eta})$

as possible. During the pure exploitation phase, ϵ_n is decayed to an intermediate value while α_n is held constant at its exploration value so that the Q Table is updated to reflect the one step optimal actions. After the pure exploration and pure exploitation phases, we have the learning phase where both α_n and ϵ_n are decayed to facilitate convergence of the Q learning algorithm. Note that each training episode corresponds to one trading day from 9:30am until 4:30pm. Since the MM wakes up every five seconds, this gives $T = 5040$ steps per episode. The precise numerics of our learning experiments are given in Table 1.

We vary the equitability weight η and the inventory weight $\bar{\eta}$ over the range of values given in Table 1. The training rewards for all $(\eta, \bar{\eta})$ are plotted as a function of training episodes in Figure 5. We are able to achieve convergence for the tabular Q learning algorithm to estimate optimal MM actions in our multi-agent simulator for the range of weights $(\eta, \bar{\eta})$ considered.

3.2.2 Understanding learnt policies for different $(\eta, \bar{\eta})$. We now attempt intuitively explaining the effect of varying the weights $(\eta, \bar{\eta})$ on the policy learnt using the objective (7). Figure 6 shows the cumulative equitability reward $\mathbb{E} \left[\frac{-H_K(y_1, \dots, y_{L-2})}{\log K} \right]$ under the learnt policy for $\eta \in \{3, 10, 20\}$, and $\bar{\eta} = 0.3$. As expected, we see that the equitability reward for the learnt policy increases as a function of the equitability weight η (for fixed $\bar{\eta}$).

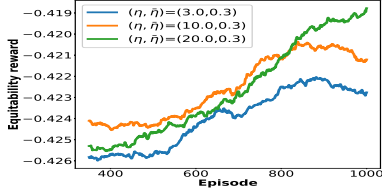


Figure 6: Equitability reward for different η

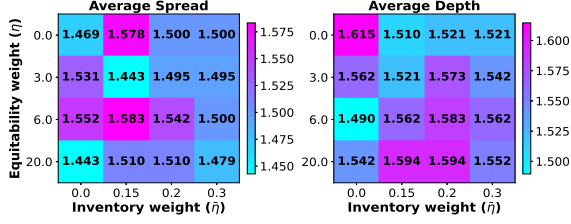


Figure 7: Average spreads and depths of orders placed by MM for different $(\eta, \bar{\eta})$

Recall that the learnt policy is a map from the current state to the (estimated) optimal action to be taken in that state. In our case, the learnt policy specifies the spread and depth of orders to be placed by our stylized MM based on current market conditions described by the state. We now look at the average spread and depth of orders placed by the learnt MM policy (averaged using a uniform distribution on the states) as functions of the weights $(\eta, \bar{\eta})$. Figure 7 shows heat maps of the average spread and depth for various values of η and $\bar{\eta}$. We see that as we increase η , the MM learns to place orders at smaller average spread, and larger depths resulting in more liquid markets. And, this confirms our intuition that more liquid markets are more equitable.

3.2.3 Effect of $(\eta, \bar{\eta})$ on volatility of returns and PnL. To further explore the effects of the weights $(\eta, \bar{\eta})$ on the learnt policy, we examine the price returns defined in (3) and the cumulative MM PnL $\mathbb{E} \left[\sum_{t=0}^{T-1} (\text{inventory PnL} + \text{matched PnL}) \right]$. Figure 8 is a histogram of the resulting price returns from using the learnt policy for different $(\eta, \bar{\eta})$. Each sub-figure corresponds to the distribution of returns for a fixed $\bar{\eta}$. Within each sub-figure, we see that the distribution corresponding to the highest equitability weight $\eta = 50$ has the lowest variance. Thus, the return volatility is seen to reduce with an increase in η , even though we do not explicitly optimize for volatility in our objective.

Likewise, we plot a histogram of the cumulative MM PnL that results from using the learnt policy with different $(\eta, \bar{\eta})$ in Figure 9. Each sub-figure corresponds to the distribution of the PnL for a fixed η . Within each sub-figure, we see that the distribution corresponding to the highest inventory weight $\bar{\eta} = 0.3$ has the lowest variance. Therefore, the PnL volatility is seen to reduce with an increase in $\bar{\eta}$, suggesting that a high $\bar{\eta}$ results in a more risk sensitive MM.

3.2.4 Effect of competing MM on learnt policy. Having understood the policy learnt by our stylized MM when acting solo in the market, a natural next step would be to add a competing MM into the market and contrast the new learnt policy to the previous one. The

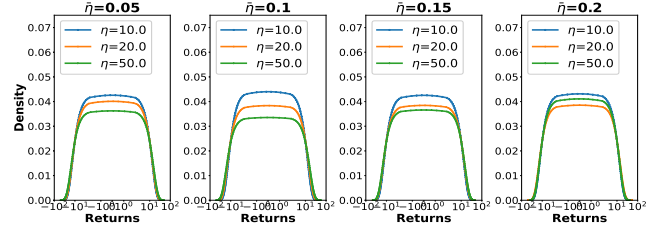


Figure 8: Distribution of consumer agent returns

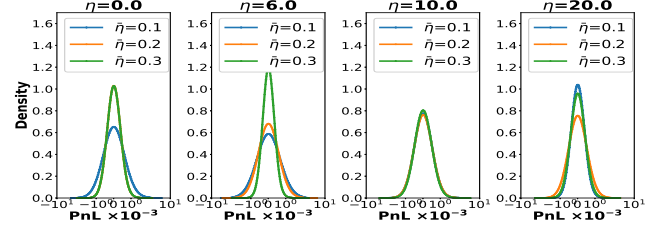


Figure 9: Distribution of MM PnL

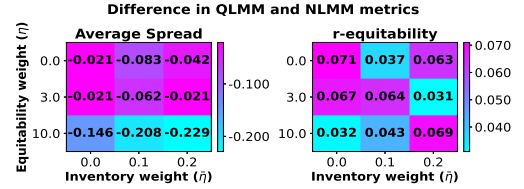


Figure 10: Effect of competing MM on learnt policy

competing MM is a non-learning based stylized MM (denoted by NLMM for non-learning MM) that posts liquidity on both sides of the spread at a fixed number of levels. Note that both MMs are made to post liquidity at the same frequency and of equal order volumes into our market. The training steps described previously are repeated for the learning based MM (denoted by QLMM for Q Learning MM), and the learnt policy is compared to that of NLMM. Figure 10 is a plot of the difference in the average spread of orders and the equitability metric between QLMM and NLMM, for different $(\eta, \bar{\eta})$. We see that QLMM learns to place at smaller spreads than NLMM in order to achieve a competitive advantage. We also see that QLMM is more equitable as per our entropy based equitability metric than the competing NLMM.

The last experiment is to check if learning helped solved the problem that motivated our work in section 1.5. In order to understand the advantages of using learning to dynamically pick the spread and depth of MM orders based on current market conditions, we augment the motivating observations with the resulting MM profit and equitability to consumer agents for the learnt policies as in Figure 11. We see that learning helps improve the profits of the MM alongside preserving market equitability.

4 CONCLUSION AND REMARKS

In this paper, we analyzed the connection between a MM's profitability and the equitability of outcomes and conclude that market

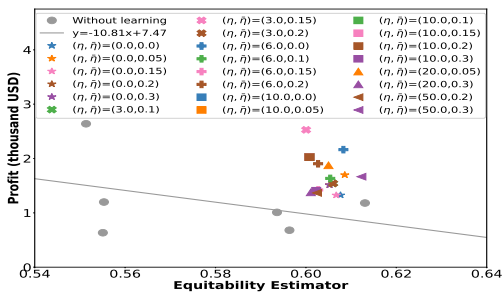


Figure 11: MM profits versus equitability to consumer agents

configurations that are more equitable are less profitable for the MM. Not only do MMs have the capacity to strongly affect market equitability, they also lose money for enabling it. Accordingly, regulators and exchanges may consider incentives for MM behavior to encourage less volatile market outcomes that are more equitable. Note that our findings do not rely on a specific MM model, and only use stylized properties from the regulatory definition of a MM. We further demonstrated the ability to derive MM policies using RL that improve upon equitability outcomes of consumer agents in a simulated market. This opens up the possibility of improving on MM profits by using equitability as part of the MM's objectives.

ACKNOWLEDGMENTS

This paper was prepared for informational purposes by the Artificial Intelligence Research group of JPMorgan Chase & Co and its affiliates ("JP Morgan"), and is not a product of the Research Department of JP Morgan. JP Morgan makes no representation and warranty whatsoever and disclaims all liability, for the completeness, accuracy or reliability of the information contained herein. This document is not intended as investment research or investment advice, or a recommendation, offer or solicitation for the purchase or sale of any security, financial instrument, financial product or service, or to be used in any way for evaluating the merits of participating in any transaction, and shall not constitute a solicitation under any jurisdiction or to any person, if such solicitation under such jurisdiction or to such person would be unlawful.

REFERENCES

- [1] James Angel and Douglas McCabe. 2013. Fairness in Financial Markets: The Case of High Frequency Trading. *Journal of Business Ethics* (2013).
- [2] Marco Avellaneda and Sasha Stoikov. 2008. High-frequency trading in a limit order book.
- [3] Tucker Balch. 1998. *Behavioral diversity in learning robot teams*. Ph.D. Dissertation, Georgia Institute of Technology.
- [4] Tucker R. Balch. 2000. Hierarchic Social Entropy: An Information Theoretic Measure of Robot Group Diversity. *Autonomous Robots* 8 (2000), 209–238.
- [5] Rachel KE Bellamy, Kuntal Dey, Michael Hind, Samuel C Hoffman, Stephanie Houde, Kalapriya Kannan, Pranay Lohia, Jacquelyn Martino, Sameep Mehta, Aleksandra Mojsilovic, et al. 2018. AI Fairness 360: An extensible toolkit for detecting, understanding, and mitigating unwanted algorithmic bias. *arXiv preprint arXiv:1810.01943* (2018).
- [6] Christopher Bone and Suzana Dragičević. 2009. GIS and Intelligent Agents for Multiobjective Natural Resource Allocation: A Reinforcement Learning Approach. *Transactions in GIS* 13, 3 (2009), 253–272. <https://doi.org/10.1111/j.1467-9671.2009.01151.x> arXiv:<https://onlinelibrary.wiley.com/doi/pdf/10.1111/j.1467-9671.2009.01151.x>
- [7] Stephen Boyd, Stephen P Boyd, and Lieven Vandenbergh. 2004. *Convex optimization*. Cambridge university press.
- [8] Tanmoy Chakraborty and Michael Kearns. 2011. Market making and mean reversion. In *Proceedings of the 12th ACM conference on Electronic commerce*.

- ACM, 307–314.
- [9] Nicholas Tung Chan and Christian Shelton. 2001. An electronic market-maker. (2001).
- [10] Cynthia Dwork, Moritz Hardt, Toniann Pitassi, Omer Reingold, and Richard Zemel. 2012. Fairness through awareness. In *Proceedings of the 3rd innovations in theoretical computer science conference*. 214–226.
- [11] Thierry Foucault, Ohad Kadan, and Eugene Kandel. 2009. *Liquidity cycles and make/take fees in electronic markets*. Technical Report.
- [12] Corrado Gini. 1936. On the measure of concentration with special reference to income and statistics. *Colorado College Publication, General Series* 208, 1 (1936), 73–79.
- [13] EG Gladyshev. 1965. On stochastic approximation. *Theory of Probability & Its Applications* 10, 2 (1965), 275–278.
- [14] Martin D Gould, Mason A Porter, Stacy Williams, Mark McDonald, Daniel J Fenn, and Sam D Howison. 2013. Limit order books. *Quantitative Finance* 13, 11 (2013), 1709–1742.
- [15] Andrei Kirilenko, Albert Kyle, MEHRDAD SAMADI, and Tugkan Tuzun. 2017. The Flash Crash: High Frequency Trading in an Electronic Market. *The Journal of Finance* (06 2017).
- [16] Andrei Kirilenko and Andrew Lo. 2013. Moore's Law vs. Murphy's Law: Algorithmic Trading and Its Discontents. *Journal of Economic Perspectives* (03 2013).
- [17] Albert S Kyle. 1985. Continuous auctions and insider trading. *Econometrica: Journal of the Econometric Society* (1985), 1315–1335.
- [18] Tian Lan, David Kao, Mung Chiang, and Ashutosh Sabharwal. 2010. An Axiomatic Theory of Fairness in Network Resource Allocation. In *Proceedings of the 29th Conference on Information Communications (INFOCOM'10)*.
- [19] Hervé Moulin. 2003. *Fair Division and Collective Welfare*. The MIT Press.
- [20] Y. Nojima, F. Kojima, and N. Kubota. 2003. Local episode-based learning of multi-objective behavior coordination for a mobile robot in dynamic environments. In *The 12th IEEE International Conference on Fuzzy Systems, 2003. FUZZ '03*, Vol. 1. 307–312 vol.1. <https://doi.org/10.1109/FUZZ.2003.1209380>
- [21] JOHN RAWLS. 1999. *A Theory of Justice*. Harvard University Press. <http://www.jstor.org/stable/j.ctvkjb25m>
- [22] Herbert Robbins and Sutton Monro. 1951. A stochastic approximation method. *The annals of mathematical statistics* (1951), 400–407.
- [23] Elizabeth Roberto. 2016. Measuring inequality and segregation. *arXiv preprint arXiv:1508.01167* (2016).
- [24] Diederik M Roijers, Peter Vamplew, Shimon Whiteson, and Richard Dazeley. 2013. A survey of multi-objective sequential decision-making. *Journal of Artificial Intelligence Research* 48 (2013), 67–113.
- [25] Ioanid Rosu. [n.d.]. A dynamic model of the limit order book. *Review of Financial Studies* (n. d.), 2009.
- [26] Jeff Schwartz. 2010. Fairness, utility and market risk. *Oregon Law Review* (2010).
- [27] Securities and Exchange Commission. 2012. Order Instituting Proceedings to Determine Whether to Approve or Disapprove Proposed Rule Changes Relating to Market Maker Incentive Programs for Certain Exchange-Traded Products. (2012).
- [28] Christian Robert Shelton. 2001. Importance sampling for reinforcement learning with multiple objectives. (2001).
- [29] Till Speicher, Hoda Heidari, Nina Grgic-Hlaca, Krishna P Gummedi, Adish Singla, Adrian Weller, and Muhammad Bilal Zafar. 2018. A unified approach to quantifying algorithmic unfairness: Measuring individual & group unfairness via inequality indices. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. 2239–2248.
- [30] Thomas Spooner, John Fearnley, Rahul Savani, and Andreas Koukorinis. 2018. Market Making via Reinforcement Learning. In *Proceedings of the 17th International Conference on Autonomous Agents and MultiAgent Systems*.
- [31] Richard S Sutton and Andrew G Barto. 2018. *Reinforcement learning: An introduction*. MIT press.
- [32] Henri Theil. 1967. *Economics and information theory*. Amsterdam: North-Holland 1967.
- [33] P. Twu, Y. Mostofi, and M. Egerstedt. 2014. A measure of heterogeneity in multi-agent systems. In *2014 American Control Conference*.
- [34] Kumar Venkataraman and Andrew C. Waisburd. 2007. The Value of the Designated Market Maker. *The Journal of Financial and Quantitative Analysis* 42, 3 (2007), 735–758. <http://www.jstor.org/stable/27647318>
- [35] Elaine Wah, Mason Wright, and Michael P Wellman. 2017. Welfare effects of market making in continuous double auctions. *Journal of Artificial Intelligence Research* 59 (2017), 613–650.
- [36] Xintong Wang and Michael P Wellman. 2017. Spoofing the limit order book: An agent-based model. In *Proceedings of the 16th Conference on Autonomous Agents and MultiAgent Systems*. 651–659.
- [37] Christopher JCH Watkins and Peter Dayan. 1992. Q-learning. *Machine learning* 8, 3-4 (1992), 279–292.
- [38] H. Peyton Young. 1994. *Equity: In Theory and Practice*. Princeton University Press. <http://www.jstor.org/stable/j.ctv10crf7>