

Limited or Biased: Modeling Sub-Rational Human Investors in Financial Markets

Penghang Liu^a, Kshama Dwarakanath^b, Svitlana S Vyetrenko^b, and Tucker Balch^a

^aJ.P.Morgan AI Research, New York, New York, USA; ^bJ.P.Morgan AI Research, Palo Alto, California, USA

ARTICLE HISTORY

Compiled March 12, 2024

ABSTRACT

Human decision-making in real-life deviates significantly from the optimal decisions made by fully rational agents, primarily due to computational limitations or psychological biases. While existing studies in behavioral finance have discovered various aspects of human sub-rationality, there lacks a comprehensive framework to transfer these findings into an adaptive human model applicable across diverse financial market scenarios. In this study, we introduce a flexible model that incorporates five different aspects of human sub-rationality using reinforcement learning. Our model is trained using a high-fidelity multi-agent market simulator, which overcomes limitations associated with the scarcity of labeled data of individual investors. We evaluate the behavior of sub-rational human investors using hand-crafted market scenarios and SHAP value analysis, showing that our model accurately reproduces the observations in the previous studies and reveals insights of the driving factors of human behavior. Finally, we explore the impact of sub-rationality on the investor's Profit and Loss (PnL) and market quality. Our experiments reveal that bounded-rational and prospect-biased human behaviors improve liquidity but diminish price efficiency, whereas human behavior influenced by myopia, optimism, and pessimism reduces market liquidity.

KEYWORDS

Human behavior; Prospect Theory; Bounded Rationality; Reinforcement learning; Multi-agent systems; Market simulations

1. Introduction

In traditional economics and game theoretic studies, humans are often conceptualized as *homo economicus* – fully rational, self-interested agents with unbiased beliefs. However, empirical evidence has consistently demonstrated that real human behavior is far more intricate, and may not always align with perfect decision making (Thaler et al. (1997); Benartzi and Thaler (1999); Chrisman and Patel (2012)). Over the recent years, there has been a growing interest in behavioral economics in response to these complexities (Thaler (2016)). It attempts to incorporate insights from other social sciences, especially psychology, in order to enrich the standard economic models that fail to explain human behavior in real life.

We refer to such behavior as being *sub-rational*, as opposed to perfectly rational decision-making. Many studies have been developed to reveal the various drivers of

human sub-rationality, which can be categorized into the following two classes. First, human sub-rationality is **limited**. Simon (1955) suggested that humans may attempt to take sub-optimal decisions that are satisfactory rather than optimal due to limited access to information and processing power. The second facet of human sub-rationality is that humans are psychologically or cognitively **biased**. For example, Kahneman and Tversky (1979) developed prospect theory to explain human decision making under uncertainty. Benartzi and Thaler (1995); Thaler et al. (1997) discussed that human decision-making can be myopic, with a strong emphasis on short-term gains/losses as opposed to those over longer terms.

Despite the richness of discoveries in aspects of human sub-rationality, there lacks an effective and universal approach to model and predict human trading behavior influenced by these aspects in financial markets. Consequently, existing behavioral finance studies are mostly constrained on real market observations, and it remains unclear how each aspect of human sub-rationality will impact their investment strategy, profit, and the market dynamics. To fill this gap, we employ reinforcement learning (RL) to model market participants. The RL agents are trained in the market environment to learn a trading strategy that optimizes the specified reward function. In one of the first financial market applications of RL, Nevmyvaka, Feng, and Kearns (2006) used RL to train an optimal algorithmic execution agent. Similarly, in subsequent work, Spooner et al. (2018); Dwarakanath, Vyetrenko, and Balch (2021) use RL to design market makers that provide liquidity in the market. By default, the RL agent will obtain a trading strategy that is optimal in maximizing its cumulative rewards, upon sufficient training. Subsequently, we modify the Bellman equation to model each aspect of human sub-rationality as a deviation from these optimal decisions.

In financial markets, labelled data that identifies market participants is typically proprietary, and it is often impossible to retrieve individual level trading activities of human investors (Gutiérrez-Roig et al. (2019)). To address this limitation, we use a high-fidelity market simulator to design, train and test our RL agents. Research in finance is well facilitated by versatile market simulations, which provide feasible experiment control and concrete market observations (Friedman (2018)). Multi-agent market simulators have been applied in financial research to reproduce the scaling laws for returns, assess the benefits of co-location, investigate the impact of large orders, and evaluate trading strategies (Lux and Marchesi (1999); Byrd, Hybinette, and Balch (2019); Balch et al. (2019)). These simulators promote the use of RL algorithms to learn complex trading strategies in a shielded, simulated environment before trying them out in real markets. In particular, we employ ABIDES-Gym (Byrd, Hybinette, and Balch (2019); Amrouni et al. (2021)), which provides discrete event time based Discrete Event Multi-Agent Simulation (DEMAS) for financial markets. It allows us to simulate fine-grained intraday market dynamics, and assess the impacts of human trading behavior on market micro-structure based on the limit order book (LOB). In addition, we build probabilistic neural networks to model the biased internal beliefs of human investors that may differ from true market dynamics given by the simulator.

In this paper, we model and examine the behavior of sub-rational human investors in financial markets resulting from computational limitations or psychological biases. We introduce five types of sub-rational human investors: bounded rational, psychologically myopic, prospect biased, optimistic, and pessimistic. For each type of human investor, we demonstrate the corresponding trading strategy in a hand-crafted market scenario to intuitively explain the strategy. We also investigate the relationship between the investor’s profits and loss (PnL) and their corresponding sub-rationality. In addition, we use an explainability tool called SHAP to investigate the major driving factors of

Table 1. Summary of sub-rational human trading behavior.

	Bounded	Myopic	Prospect Biased	Pessimistic	Optimistic
Method	Boltzmann softmax	temporal discounting	biased reward & internal model	biased internal model	biased internal model
Behavior	random error in decisions	aligned with short-term momentum	risk-averse in gains risk-seeking in losses	small inventory infrequent trading	large inventory aggressive buy/sell
Deciding Factors	N/A	inventory short-term momentum	inventory	N/A	N/A
PnL Variance	high	low	high	low	high
Improve Liquidity	✓	✗	✓	✗	✗
Reduce Volatility	✓	✗	✗	✗	✗
Increase Efficiency	✗	✓	✗	✓	✗

human decision-making. Last but not least, our experimental analysis discovers the impact of sub-rational investors on the market with regard to liquidity, volatility, and market efficiency.

To the best of our knowledge, this is the first work that utilizes the RL framework to model human trading behavior in financial markets as caused by different aspects of sub-rationality. We summarize the main contributions of our work as follows:

- We develop a comprehensive and adaptive RL framework to model human sub-rational trading behavior characterized by bounded rationality, myopia, prospect bias, optimism, and pessimism.
- To address the limited availability of real, individual trader data, we utilize market simulations to train and test the sub-rational RL agents. In addition, we build probabilistic neural networks to model the internal bias in human beliefs.
- We craft specific market scenarios and use SHAP analysis to deeply examine the different types of sub-rational trading behavior. Furthermore, we test the agents in controlled simulated markets to evaluate the impact of sub-rational behavior on market observables.

We believe our models will provide an effective framework to capture and examine human investor behavior while aiding in better understanding of their influence in financial markets.

2. Literature Review

2.1. Human Sub-rationality

One of the most important topics in behavioral finance is to understand and model the decision-making process of humans in real life. Simon (1955) was one of the first critics of modeling economic agents as having unlimited information and processing capabilities. He introduced the idea of “bounded rationality” to describe a more realistic conception of human problem solving capabilities. Lots of the departures from rational choice can be captured by the prospect theory proposed by Kahneman and Tversky (1979), which provides a purely descriptive theory of human decision-making under uncertainty. Ainslie (1975); Thaler and Shefrin (1981) discussed the importance of self-control in human behavior, where temporal discounting is often utilized to explain human’s myopic preference for short-term results rather than long-term results. Researchers have also discovered other aspects of human sub-rationality as resulting

from biased internal beliefs such as optimism (Sharot et al. (2007)) and illusion of control (Barber and Odean (2002); Song et al. (2013)). While we believe that humans are far too complicated to be represented by any mathematical model, the focus of our paper is to develop a unified framework to model the decision-making of humans influenced by different types of sub-rationality discovered in previous studies.

Due to the limitation in data availability and realistic human models, there lacks analytical studies of the impact of human sub-rationality on financial markets. Meanwhile, regulatory literature has examined the impact of electronic traders as compared to that of human traders in markets. Boehmer, Fong, and Wu (2021a) investigated the impact of algorithmic trading on equity market quality measured in terms of quoted spread, price efficiency and volatility. They observed that more algorithmic trading lead to narrower spreads, better price efficiency but higher volatility based on real data, with the effects differing based on asset size. Woodward (2017) investigated the market impact of electronic traders that feature high frequency trading (HFT) techniques, and provided insights to control it from the perspective of regulators. Note that these studies examine the influence of algorithmic and high frequency trading techniques, which are used by a specialized class of electronic investors. In this work, we compare the impact of sub-rational human investors with that of perfectly rational electronic investors (with no change in trading frequency) in financial markets.

2.2. *Human Decision Models*

Reinforcement learning (RL) is an effective approach to obtain a decision-maker’s policy for different rewards/incentives in specific environments. Theoretical efforts have been made to model the human decision-making process using RL. A related field is inverse reinforcement learning (IRL), which aims to bypass the need for reward design by inferring the reward from observed human demonstrations. There are numerous hypotheses behind sub-rationality of humans. Raja and Lesser (2001) attributed human sub-rationality to constraints on resources and address a meta-level control problem for a resource-bounded rational human. Evans and Goodman (2015); Evans, Stuhlmüller, and Goodman (2016) modeled structured deviations from optimality with different hierarchical levels of rationality when inferring preferences and beliefs. In a recent work, Chan, Critch, and Dragan (2021) investigated the effects of human irrationality on reward inference. They introduced a framework to describe different aspects of human irrationality using the Bellman equation by modifying the max operator, the transition function, the reward function, and the sum between reward and future value. While literature in IRL encompasses various interesting models of humans, the goal is to infer the reward function from real human demonstrations. However, in financial markets, it is rarely possible to acquire historical trading data of individual, human investors. The goal of our work is to use behavioral models for human traders to model an investor agent that trades like a human in financial markets, and to subsequently analyze human traders’ impact on the market in a simulated environment.

2.3. *Multi-agent Market Simulations*

Multi-agent simulators have become increasingly prevalent for modeling financial markets, which provide an alternative yet effective approach to study the market in rarely observed scenarios or with limitations of data. Lux and Marchesi (1999) introduced a multi-agent model of financial markets to support the time scaling law from mutual in-

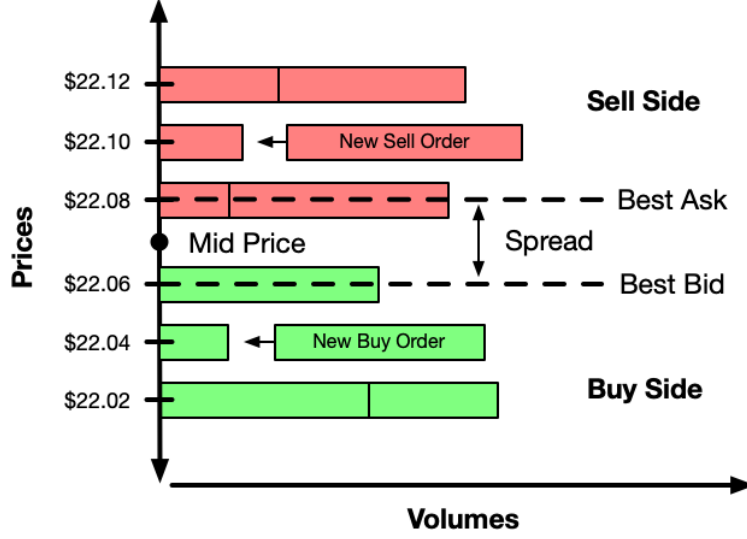


Figure 1. A snapshot of the LOB structure.

teractions of participants. In recent contributions, Byrd, Hybinette, and Balch (2019) developed a discrete event simulator to investigate the market impact of a large market order. Additionally, Vyetrenko et al. (2020) proposed realism metrics to evaluate the fidelity of the simulated markets. While these multi-agent market simulators can be populated with rule based trading agents, they allow for the use of reinforcement learning to develop agents that seek to optimize certain objectives. Amrouni et al. (2021) wrapped market simulations in an OpenAI Gym framework facilitating the use of off-the-shelf RL algorithms, and train investors in various market environments. Dwarakanath, Vyetrenko, and Balch (2021) utilized RL to obtain the policy of market makers and subsequently investigated their impact on market equitability.

3. Background

In this paper, we analyze and evaluate human investor behavior based on the market microstructure. We now introduce the underlying concepts of market microstructure including the limit order book, and the corresponding market observables.

3.1. Limit Order Book (LOB) structure

Due to the prevalence of high frequency trading, most of the financial market exchanges including NASDAQ and NYSE are now operating on an underlying dynamic mechanism known as the Limit Order Book (LOB) (Roşu (2009); Gould et al. (2013)). The LOB is a dynamic repository of “buy” and “sell” limit orders, which provides a real-time granular view of the supply and demand at various price levels (Bouchaud, Mézard, and Potters (2002); Vyetrenko et al. (2020)). Figure 1 provides a visual representation of the LOB structure. Here we summarize the key concepts of the LOB:

- **Limit Orders:** Unlike the traditional *market orders* that represent the intentions to buy and sell immediately at the current market price, a *limit order* specifies the minimum price that the trader intends to sell at, or the maximum price that the trader is willing to buy at, alongside the volume and direction

(sell/buy).

- **Buy Side and Sell Side:** The LOB maintains limit orders on two sides, the buy and sell sides. The *best bid* and *best ask* represent the highest price a buyer is willing to pay and the lowest price a seller is willing to accept, respectively.
- **Order Queue:** Orders from various market participants are collected and aggregated into queues based on their respective prices and time. The exchange matches buy and sell orders based on a first-in-first-out (FIFO) basis.

The LOB offers a granular and accurate depiction of the underlying mechanisms governing modern financial markets. Moreover, the discrete nature of LOB offers a detailed examination of trader interactions, how liquidity is provided or consumed, and the evolution of order queues over time. Market simulations on LOB facilitate detailed behavioral analyses, enabling the exploration of how various trading strategies impact market dynamics.

3.2. Metrics and Notations

Here we describe the market metrics used for empirical analysis in our study. Following the notation in Coletta et al. (2022), we denote $p_a^i(t), v_a^i(t), p_b^i(t), v_b^i(t)$ as the quoted price and volume of ask and bid orders at i -th level of the LOB, which has n levels on the ask side and m levels on the bid side. We use the below metrics to describe the market dynamics in our study.

- **Mid price** is the average of the best ask and best bid prices, i.e.,

$$m(t) = \frac{p_a^1(t) + p_b^1(t)}{2}$$

- **Traded price** $p(t)$ is the price of the latest transaction before time t , which represents the most recent market price of the asset.
- **Spread** is the difference between the best ask and best bid prices:

$$\Delta(t) = p_a^1(t) - p_b^1(t)$$

- **Market depth** is the price difference between the worst ask and worst bid prices:

$$d(t) = p_a^n(t) - p_b^m(t)$$

- **Traded Price Volatility** represents the degree of variation of the traded price series over the past δ time window:

$$\sigma(t, \delta) = \text{standard deviation}(p(t - \delta), p(t - \delta + 1), \dots, p(t))$$

- **Momentum** represents the trend of traded price over the over the past δ time window:

$$\text{momentum}(t, \delta) = \frac{p(t) - p(t - \delta)}{\delta}$$

- **Volume imbalance** indicates the supply and demand inequality within the

LOB:

$$I(t) = \frac{\sum_{j=1}^m v_b^j(t)}{\sum_{i=1}^n v_a^i(t) + \sum_{j=1}^m v_b^j(t)}$$

Volume imbalance is known to be an informative feature for trading as it can signal private information that can subsequently reduce market liquidity, and can impact the decisions of market makers thereby affecting market returns as well as execution agents Chordia, Roll, and Subrahmanyam (2002); Stoll (2003); Zheng, Moulines, and Abergel (2012).

3.3. Measures of Market Quality

Following the experiments of Boehmer, Fong, and Wu (2021b), we evaluate the market quality on three dimensions: liquidity, volatility, and information efficiency.

3.3.1. Liquidity

Liquidity is a crucial aspect of financial markets, which indicates the ease of trading a security in the market without significantly impacting its market price. High liquidity markets typically have narrow bid-ask spreads and large traded volumes. We measure the liquidity of the stock based on the spread $\Delta(t)$, the relative effective spread RES, and the traded volume in every minute. The relative effective spread is calculated as:

$$\text{RES} = \frac{p(t) - m(t)}{m(t)} \quad (1)$$

where $p(t)$ and $m(t)$ are the traded price and mid-price at time t respectively. RES measures price changes associated with trading. The narrower the relative effective spread is, the more liquid the stock is.

3.3.2. Volatility

Market volatility is a pivotal factor for market analysis, which can exert multifaceted influences on market quality through various channels. The prices of volatile assets have strong fluctuations and are often less predictable. Greater volatility tends to increase the cost of utilizing limit orders, making liquidity provision more expensive for market participants. As volatility rises, the likelihood of prices deviating significantly from the levels set in limit orders increases, making it more challenging for traders to execute orders at their desired prices. Define the price return over a time period δ as $\frac{p(t)}{p(t-\delta)}$. We consider three measures of volatility: the standard deviation of 30-minute price returns (Std-30min-RET), the intraday price range (difference between the lowest and highest prices in a day), and the absolute value of daily returns ($|\text{RET}|$). Markets with low volatility have low values for the three volatility metrics.

3.3.3. Market efficiency

According to the efficient market hypothesis (Fama (1970); Malkiel (1989)), market efficiency (also known as information efficiency or price efficiency) represents the degree to which available information is reflected in market prices. If market prices fail

to fully encompass all available information, opportunities may arise for profit through the collection and processing of such information. An inefficient price also indicates that the markets are not doing an efficient job of allocating resources. We consider the autocorrelation of mid-price as an indicator of price efficiency following Boehmer and Kelley (2009), which suggests that efficient pricing should follow a random walk and the autocorrelation should be close to 0. While it is infeasible to encompass all market information in the history in the real-world, in market simulations, we model and incorporate such information into the simulated fundamental value (see 4.2). Consequently, our measurement of market efficiency extends to $|\text{TP-FP}| / \text{FP}$, representing the relative disparity between the traded price and the fundamental price concurrently. A minimal separation signifies that the market operates with a high degree of information efficiency.

4. Reinforcement Learning for Investor Modeling

In this section we describe how to model investors using reinforcement learning (RL). We choose to use the RL framework to model human investors due to the following reasons: (1) RL is more flexible compared to rule based approaches, and can adapt to different market scenarios. RL allows a more general representation of human behavior as arising from objective maximization rather than rigidly defined rule based behaviors. (2) RL offers a mathematical interpretation of different sub-rational human behaviors as deviations from the optimal policy, and provide controls of the degree of human sub-rationality. (3) RL models can be trained and tested in market simulations efficiently, which helps in evaluating the performance of different trading strategies. It also allows us to study diverse subrational behaviors without the need for real human demonstrations.

4.1. Defining RL Investors

In this work, we utilize reinforcement learning to obtain trading strategies for investors. Formally, we consider the Markov decision process (MDP) (S, A, P_a, R_a) where S and A are the sets of states and actions. $P(s'|s, a)$ is the transition probability from state s to s' by taking the action a , and $R(s, a, s')$ is the immediate reward of the transition. The goal in RL is to maximize the expected sum of discounted rewards. The Bellman equation relates the value of the current state s to the one step reward and the value at the next state s' as

$$V(s) = \max_a \sum_{s' \in S} P(s'|s, a) (R(s, a, s') + \gamma V(s')) \quad (2)$$

where γ is the discount factor. We define a fully rational agent as one that picks the action that maximizes the right hand side of Equation (2) when in state s . The corresponding optimal (rational) policy is

$$\pi(s) = \arg \max_a \sum_{s' \in S} P(s'|s, a) (R(s, a, s') + \gamma V(s')) \quad (3)$$

The RL agent represents an investor that tries to make profits by trading in the market. We model the rational investors and sub-rational human investors using the

same state space, action space, and reward function.

States: The states include the agent’s observation of the market along with its internal states, including

- Cash_t : the amount of cash in the investor’s account at time t . The agent starts the day with 1,000,000 cents.
- Holdings_t : the number of shares of the asset held by the agent at time t . The agent starts the day with no holding.
- $[\text{momentum}(t, 1), \text{momentum}(t, 10), \text{momentum}(t, 30)]$: a vector reflecting the market momentum over the past 1, 10, and 30 minutes.
- $\Delta(t)$: the market spread at time t .
- $d(t)$: the market depth at time t .
- $\sigma(t, 30)$: the traded price volatility in the past 30 minutes.
- Quote history: the information of quoted/placed orders in the past five minutes, including quoted prices and volumes.
- Trade history: the information of the executed trades in the past five minutes, including traded prices and volumes.

Actions: The RL agent wakes up ever minute between 9:30am to 4:00pm, and takes an action defined by the following parameters

- Direction: {BUY, HOLD, SELL}. In the case of BUY or SELL, the agent places a limit order of size 2.
- Limit order price (relative to the mid price): $\{-0.5, -1, -1.5, -2\}$ if the agent takes a BUY action, or $\{+0.5, +1, +1.5, +2\}$ if the agent takes a SELL action.
- Note that RL policy is formulated as a MDP, i.e., the transition from current state s to the next state s' only depends on the action a . We cancel all previous open orders of the RL agent before placing the new order from action a so that the changes of inventory only come from the most recent orders.

Rewards: We define the one step reward as

$$R(s_t, a, s_{t+1}) = \text{PortfolioValue}_t - \text{PortfolioValue}_{t-1} \quad (4)$$

where $\text{PortfolioValue}_t = \text{Cash}_t + \text{Holdings}_t \times p(t)$ is the value of the investor’s portfolio marked to the market, and $p(t)$ is the executed price of the last transaction before time t . Since the reward function measures the change of portfolio values in every minute, the cumulative reward over the horizon equals the agent’s monetary profits at the end of the trading day (assuming there is no transaction cost and with temporal discounting $\gamma = 1$). Here we formally define the profit and loss (PnL) at time t as the change in portfolio value compared to the start of the day:

$$\text{PnL} = \text{PortfolioValue}_t - \text{PortfolioValue}_{t_0} \quad (5)$$

4.2. Training in Multi-Agent Market Simulations

Market simulations are often employed to train RL agents as they offer a controlled environment to learn an optimal trading strategy in presence of reactive background trading agents, while also dealing with the inherent challenge of acquiring proprietary historical LOB data. In this paper, we employ ABIDES-Gym, a discrete event multi-agent simulator that provides a high-fidelity market environment with thousands of trading agents, wrapped in an OpenAI Gym framework Byrd, Hybinette, and Balch

Table 2. The configuration of simulated markets. Except the exchange agent and market makers, the order sizes of the the background agents are randomly sampled from a heavy-tailed distribution with mean equals to 100. The actions of value agent follows a Poisson process with arrival rate $\lambda = 1/5.7 \cdot 10^{-12}$ in nanoseconds, i.e., once per 3 minutes on average.

Agent Type	# of Agents	Order Size	Action Frequency
Value Agent	2	100	Poisson ($\lambda = 1/5.7 \cdot 10^{-12}$) (ns)
Market Maker	1	10	Every minute
Momentum Agent	2	100	Every 5 minutes
Noise Agent	20	100	Once a day

(2019); Amrouni et al. (2021). The simulated market contains the following background agents with different trading strategies and incentives.

- **Value Agents:** The value agents are designed to emulate fundamental traders who trade according to their belief of the exogenous value of an asset, i.e., *fundamental price* Wang et al. (2021). In this paper, we generate the time series of the fundamental price with a discrete-time mean-reverting Ornstein-Uhlenbeck process. Each value agent makes a noisy observation of the fundamental price, and places a sell order if the current mid price is higher than the observed fundamental price, or vice versa.
- **Market Maker Agent:** The market maker agent acts as a liquidity provider in the market, by placing limit orders on both sides of the LOB at regular intervals. The agent places equal volumes of liquidity at various price levels that are pre-defined with respect to the mid-price.
- **Momentum Agents:** The momentum agents trade based on the momentum of the asset price, by comparing the long-term average of the mid-price with the short-term average. Each agent places a buy order if the short-term average is higher, based on the belief that the price will increase in the future. On the other hand, if the long-term average is higher, the agent places a sell order.
- **Noise Agents:** The noise agents mimic retail traders that trade based on their own demand with no consideration of the LOB microstructure. Each noise agent arrives to the market at a random time, and places an order on a random side of the LOB.

We train the RL agent in the simulated market environment along with the background agents. As illustrated in Figure 2, the market simulator receives the action of the RL agent at each step and provides the next state and reward information. The RL agent optimizes its trading policy based on the feedback from the environment.

4.3. Training on Internal Beliefs

Many human sub-optimal actions can be considered as the result of internal model error. That is to say that the reason that human actions might deviate from the optimal is that they have incorrect internal beliefs of the rules (dynamics) guiding how their actions affect the environment (Reddy, Dragan, and Levine (2018)). Therefore, the agent obtains a biased policy which is sub-optimal in the real world but near-optimal with respect to its internal model of the dynamics. In reinforcement learning, the environment dynamics are captured by the transition probability $P(s'|s, a)$ in

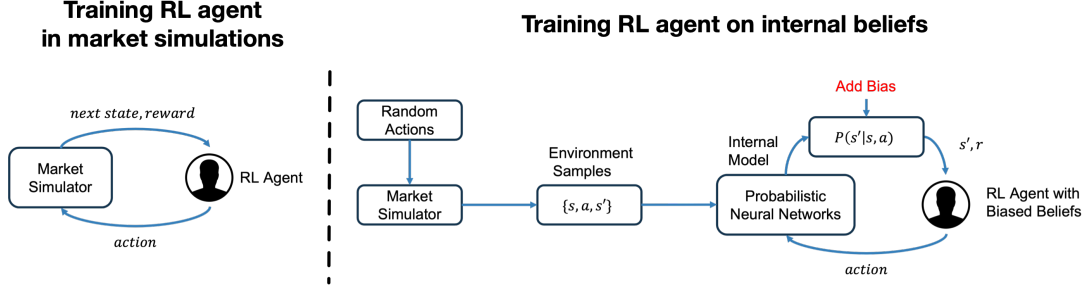


Figure 2. Training RL agents using market simulations or biased internal beliefs. (left) The RL agent learns a trading strategy by directly interacting with the simulated markets. (right) The RL agent learns from the internal model. We first learn a probabilistic internal model from the samples of the environment. We can inject bias to the internal model and use it to train a biased human.

Equation (2), which represents the probability of the next state s' given the current state s and action a . A sub-rational human has a internal model which gives biased estimations of the transition probability. For example, the optimistic/pessimistic agent in Section 5.2.3 has biased estimations of positive outcomes over the negative outcomes. The prospect biased agent in Section 5.2.2 overestimates the probability of unlikely events.

To model the internal beliefs of a human trader, we collect samples of the environment by feeding random actions to the market simulator as illustrated in Figure 2. Each sample (s, a, s') represents the transition of the environment state from s to s' when action a is taken by the agent. We use the environment samples to train the internal model, which predicts the probability distribution of the next states given the current state and an action: $P(s'|s, a)$. In particular, we use a bootstrap ensemble of probabilistic neural network (Janner et al. (2019)) which outputs the prediction as a parameterized Gaussian distribution with diagonal covariance $P(s'|s, a) = \mathcal{N}(\mu(s, a), \Sigma(s, a))$. We validate the prediction results of the internal model in Section 9.1.

Once we obtain the internal model, it can be used as a substitute of the real environment to provide the transition probability $P(s'|s, a)$ for training RL agents, while the rewards are given by the deterministic reward function in Equation (4). Adding bias to the internal model which is used in RL training will yield a biased human policy that is optimal with respect to the biased internal model, but sub-optimal in the real environment. For example, we can bias the internal model to overestimate the probability of positive s' and use it to train an optimistic human agent. Similarly, we can modify the transition probability $P(s'|s, a)$ using Equation (11) to obtain a prospect biased human policy.

5. Models of Human Sub-Rationality

While the Bellman equation (Equation (2)) provides a solution to the optimal policy, it fails to model the deficiencies in human decisions. Based on their deviations from the optimal actions and the underlying driving factors, we classify human sub-rational behavior into two categories: (1) The **bounded (limited) human behavior** which show as random errors in decision-making, and (2) the **biased human behavior** that are systematically deviated from the optimal behavior.

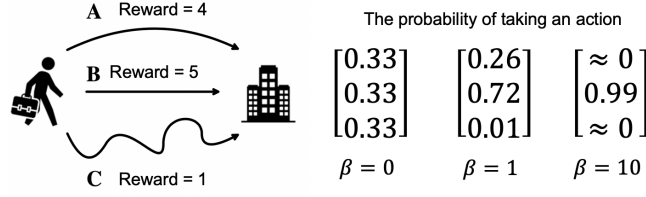


Figure 3. An example of the Boltzmann rationality model from Laidlaw and Dragan (2022) for three actions with different rewards (*left*). For demonstration purpose here we simply consider a one step decision problem in a deterministic environment. The Boltzmann model gives the probability of taking an action using the β parameter that adjusts the degree of rationality (*right*). If $\beta = 0$, each action has the same probability to be selected. When $\beta = 10$, the model becomes more rational and only the action with highest reward is likely to be selected.

5.1. Bounded Human Behavior

To arrive at the optimal decision, one ideally needs comprehensive knowledge about all available choices and the ability to calculate the potential benefits. However, this assumption proves too challenging for humans in daily life due to their limited information access and computational power (Simon (1955)). Consequently, human decision-making tends to exhibit a certain degree of noise, and their decisions often incorporate random errors deviating from the optimal action. Simon (1997, 1990) proposed the idea of bounded rationality, which suggests that humans are limited and tend to select a satisfactory decision (approximately optimal) rather than the optimal decision.

5.1.1. Boltzmann Rationality

The most common approach for modeling bounded rational human decision-making is the Boltzmann rationality model (Baker, Saxe, and Tenenbaum (2005); Ziebart, Bagnell, and Dey (2010); Asadi and Littman (2017)). Unlike the rational policy in Equation (3) that gives the optimal action using $\arg \max$, this model considers a probabilistic policy given by the Boltzmann softmax operator as follows:

$$\pi(a|s) = \frac{e^{\beta Q(s,a)}}{\sum_{a' \in A} e^{\beta Q(s,a')}} \quad (6)$$

where $Q(s, a) = \sum_{s' \in S} P(s'|s, a) (R(s, a, s') + \gamma V(s'))$ (Asadi and Littman (2017)). The

Boltzmann softmax operator introduces a “soft” optimization principle, which relaxes the assumption of unlimited resources and processing power, and allows the agent to take sub-optimal decisions with a preference for high-utility actions. Figure 3 gives an example of modeling human behavior with Boltzmann rationality. The parameter β controls how likely the agent is to select the optimal action. If $\beta = 0$, the agent has zero intelligence and makes uniformly random decisions using the same probability to select every action. As β increases, the agent becomes more intelligent and makes less error in selecting the optimal action. When $\beta \rightarrow \infty$, the Boltzmann operator approaches the $\arg \max$ in Equation (3), and the agent makes fully rational decisions.

An example of bounded rational trading behavior To illustrate the bounded rational human trading behavior, we simulate a market with fundamental price fol-

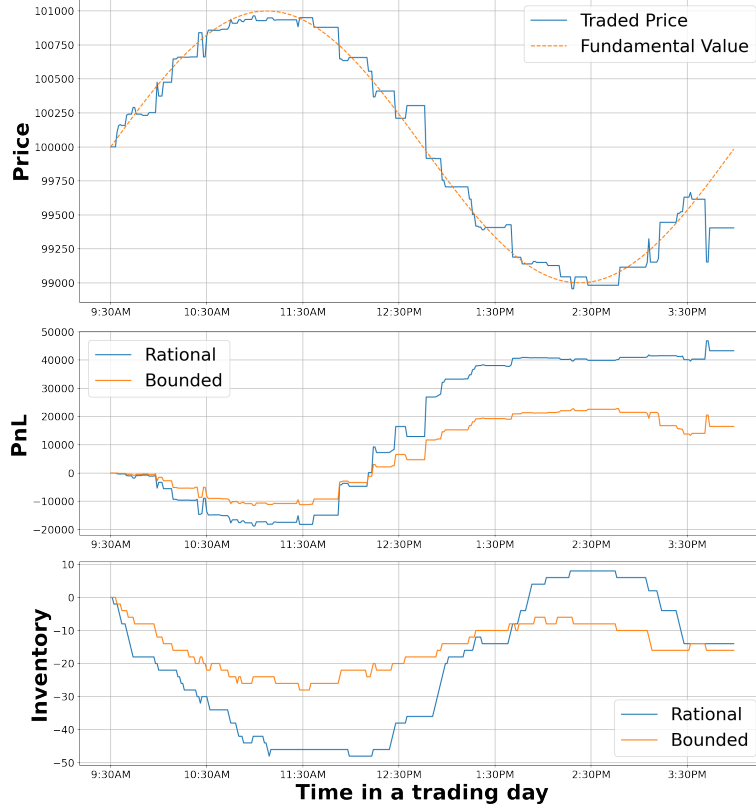


Figure 4. The behavior of bounded rational investors in the simulated market. Compared to the fully rational investor, a bounded rational human investor takes sub-optimal actions that are similar but inferior.

lowing a sine wave and deploy a bounded rational human investor with $\beta = 0.2$ and a fully rational investor with $\beta = \infty$ (see Equation (6)). The policy of the bounded rational human investor is obtained by passing the Q values of the fully rational investor through the Boltzmann softmax operator as in Equation (6). Figure 4 shows the market price and the actions of the two agents through a trading day. Since the market is simulated over a single time period of the sine wave with the same starting and closing price, the optimal strategy is to sell when the price is higher than the closing price (first half of the day), and buy when the price is lower than the closing price (second half of the day).

Overall, we observe similar decisions of rational and bounded rational investors based on the holding position in Figure 4: both of them tend to sell during the first half of the day, and buy during the second half of the day. However, the bounded rational human investor does not sell and buy at the maximum capacity compared to the fully rational investor that has a deeper holding position. As a result, the bounded investor can not achieve the maximum reward even though the decisions are similar to those of the fully rational investor to some extent.

In summary, the bounded rational model emulates human traders with limited information availability and processing capacity. We observe that such limitations lead to sub-optimal decision-making that is similar albeit inferior compared to the homo economicus.

5.2. *Biased Human Behavior*

Here we describe several human cognitive and psychological biases, which lead to systematic deficiencies in human decisions.

5.2.1. *Myopic Behavior*

Humans in reality may be psychologically biased to only care about the short-term results. They make myopic decisions without considering how the actions affect them far into the future. A typical example of myopic human behavior in finance is myopic loss aversion (Thaler et al. (1997); Benartzi and Thaler (1999); Chrisman and Patel (2012)). Investors that focus on short-term return may react too passively to recent losses. As a result, they abandon their existing long-term-oriented investment plan and lose the potential to achieve better benefits in the future. Studies have shown that financial media can often facilitate myopic trading behavior by constantly reporting market news and portraying a sense of urgency to act. Investors who receive such information too frequently tend to avoid investing in riskier assets that may yield better long-term rewards (Larson, List, and Metcalfe (2016)).

Preferences of humans when choosing between short-term and long-term rewards is modeled by discounting the long-term rewards, and have been studied in literature on economics and psychology (Chabris, Laibson, and Schuldt (2010); Grüne-Yanoff (2015)). Two popular ways to model this temporal discounting of rewards include the exponential discounting model (Samuelson (1937)) and the hyperbolic discounting model (Ainslie (1992)). Numerically, the exponential discounting model expresses the present value of a reward r given t time steps from now as

$$V = r\gamma^t \tag{7}$$

where $\gamma \in [0, 1]$ is the exponential discounting factor. In contrast, the hyperbolic discounting model expresses the present value of the reward r given with delay t as

$$V = \frac{r}{1 + kt} \tag{8}$$

where $k \geq 0$ is the hyperbolic discounting factor. Figure 5 shows a comparison of both temporal discounting schemes for median values of discounting parameters observed in real humans in Green, Myerson, and McFadden (1997). For simplicity and amenability to standard RL algorithms, we assume that the discount factors γ and k do not depend on the amount of reward r although there have been empirical studies that demonstrate that the rate of discounting reduces with an increase in the amount of reward (Green, Myerson, and McFadden (1997)).

Empirical evidence from humans and pigeons support the hyperbolic discounting model as it captures reversals in preference between delayed rewards (Green, Fristoe, and Myerson (1994); Mazur (1985)). This is described in the following scenario. Suppose one is given a choice between receiving \$10,000 today versus receiving \$11,000 in a week. Most people would choose the former option to receive \$10,000 today considering the unknown factors that could prevent the receipt of the delayed \$11,000. Now suppose the choice was between receiving \$10,000 in a year versus receiving \$11,000 in a year and one week. In this scenario with rewards delayed further, most people would prefer to wait the extra week to receive \$11,000. This is based on the reasoning that the chances of not receiving \$11,000 after a year and a week are lower conditioned on

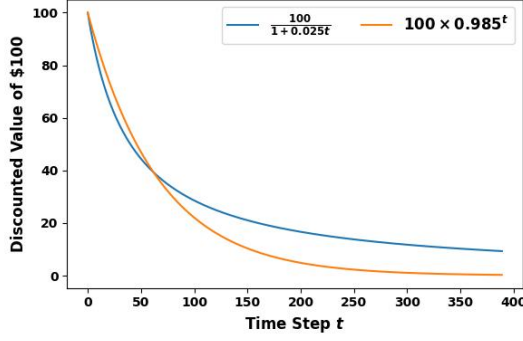


Figure 5. Comparison of exponential discounting to hyperbolic discounting of \$100 over 390 time steps. See the steady drop in discount rate for the orange exponential curve. On the other hand, the blue hyperbolic curve shows steep discounting early on and modest discounting later on.

having received \$10,000 after a year. This change in preference between the smaller and larger rewards when the delay to both is increased is called preference reversal. Exponential discounting does not capture preference reversal since the discount rate remains the same for all values of delay. On the other hand, the discount rate in the hyperbolic discounting scheme is high in the short-term and slows down in the long-term as can be seen in Figure 5.

While it maybe more pragmatic to consider hyperbolic discounting schemes for our myopic human model, non-exponential discounting poses computational challenges to standard RL algorithms that are based on the validity of the Bellman equation. There is work on formulating temporal difference learning schemes for hyperbolic discounting (Alexander and Brown (2010); Redish and Kurth-Nelson (2010)) Redish and Kurth-Nelson (2010) and later Fedus et al. (2019) propose approximating the policy for hyperbolic discounting using a combination of many exponentially discounting policies. In particular, Fedus et al. (2019) propose a solution to approximate the hyperbolic Q function using a weighted sum of exponential Q functions that are each computed for a different exponential discount factor. This work allows for the use of deep neural networks to approximate the underlying exponential Q functions. This adds a computational overhead of learning many RL policies that use different exponential discount factors alongside the instability issues that arise with the use of off-policy RL algorithms such as Deep Q-Learning that estimate the Q function (Mnih et al. (2013); Haarnoja et al. (2018)). Therefore, for the purpose of this paper, we consider exponential discounting as a model for myopicity of human traders.

Recall that the discount factor γ in equation Equation (2) models decay in the value of rewards with time delay (Chabris, Laibson, and Schuldt (2010)). To model myopic human investors, we decrease the γ in the Bellman equation Equation (2) and corresponding policy Equation (3). When $\gamma \approx 1$, the agent is fully rational and considers both short-term and long-term rewards. As $\gamma \rightarrow 0$, the agent becomes myopic and only acts to maximize the immediate reward $R(s, a, s')$.

An example of myopic trading behavior

Here we give an example to distinguish the behavior of myopic investors from that of fully rational investors, using a hand-crafted market simulation. The human investor is trained in the default market simulations with Ornstein-Uhlenbeck fundamental with $\gamma = 0.1$, while the rational investor is trained with $\gamma = 0.99$ (see Equation (3)). The

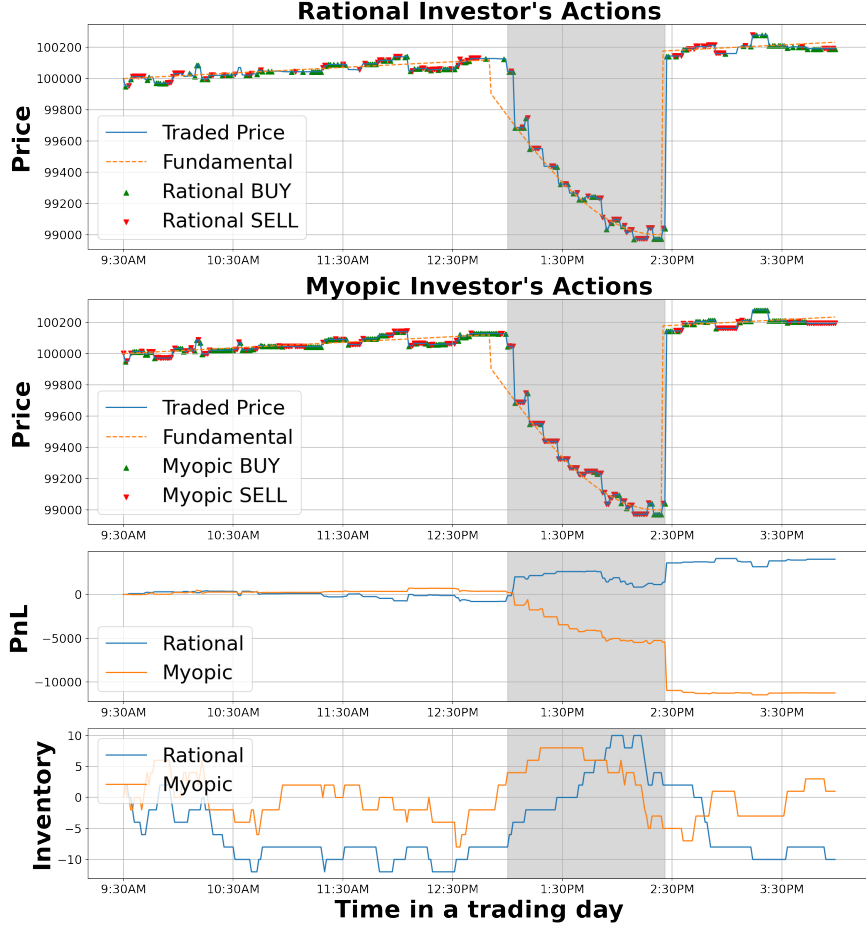


Figure 6. The behavior of myopic investors in the simulated market. Compared to the fully rational investor, a myopic human investor reacts too passively to the short-term losses and abandons its long-term investment plan.

trained investors are then tested in a hand-crafted market scenario with a different fundamental price as described here. The trading day can be decomposed into three stages. In the first and last stages, the price increases with a linear trend, while the market is under a shock in the second stage with the price temporarily dropping by 1%. Figure 6 shows the market prices and the actions of the agents through the trading day. Table 3 displays the ratio of buy orders to sell orders for both investors, alongside the percentage of hold decisions (do not buy or sell) relative to all decisions in each of the three stages of the market. Since the asset price increases continuously in stage 1, both rational and myopic investors place more buys than sells (Table 3), and achieve similar rewards. When the market is under a shock in stage 2, the myopic human investor tries to sell out their shares to mitigate the loss from the temporary price drop. They place more sell orders than buy orders, and their inventory drops significantly during the second stage. The rational investor, on the other hand, shows better vision of the future and is less affected by the shock. They place similar amount of buy and sell orders, and their inventory increases as the buy orders are easier to be executed during the shock in stage 2. By the end of the day, the rational investor obtains a significantly better PnL as they bought a lot of shares at low prices, compared to the myopic human investor who sold a lot of shares during the shock.

Overall, the myopic investor demonstrates a short-term trading strategy with a

Table 3. The decisions of investors during the three stages in Figure 6.

	Buy / Sell Ratio			Hold		
	Stage 1	Stage 2 (shock)	Stage 3	Stage 1	Stage 2 (shock)	Stage 3
Rational	1.38	0.94	1.33	38%	34%	38%
Myopic	1.49	0.41	1.26	2%	1%	7%

much smaller tendency to hold (and not trade) as compared to the rational investor especially during the shock (Table 3). Subsequently, their inventory fluctuates between long and short positions more frequently than that of the rational investor (Figure 6). In summary, the hand-crafted market scenario demonstrates that our model of psychologically myopic behavior successfully reproduces the myopic investing behavior observed in financial markets, i.e., myopic loss aversion (as described in Section 5.2.1).

5.2.2. Prospect Theory

A well-known model of decision making under uncertainty is the expected utility theory, which measures the utility of the outcomes as:

$$V = \sum_{i=1}^n p_i \cdot v(x_i) \quad (9)$$

where $v(x_1), v(x_2), \dots, v(x_n)$ are the values of the potential outcomes and p_1, p_2, \dots, p_n are their respective probabilities. The objective function of an RL agent is in line with the expected utility theory, as in Equation (2) the value of a state is measured by the sum of utilities of the potential next states weighted by their transition probabilities.

However, a substantial body of evidence has shown that human decision makers systematically violate the expected utility theory. To model such deviation, Kahneman and Tversky introduced prospect theory (Kahneman and Tversky (1979)) which consists of two key elements:

- A value function that is concave for gains, convex for losses, and steeper for losses than gains.
- A nonlinear transformation of the probability scale, which overweights small probabilities and underweights moderate and high probabilities.

Prospect Biased Utility

The first element of prospect theory stems from loss aversion, which refers to the fact that humans feel losses more than two times greater than the equivalent gains. Empirical studies (Kahneman and Tversky (1979)) have shown that human’s value function is concave for gains and convex for losses. For example, human prefer to choose 100% chance to win \$490 rather than 50% chance to win \$1000, even though the expected utility of the latter is higher. As shown in Figure 7 (left), the prospect biased utility of \$490 gain is larger than a half of the \$1000 gain. Therefore the certain gain is more preferable as $v(+\$490) > 0.5v(+\$1000) \rightarrow 100\% \times v(+\$490) > 50\% \times v(+\$1000)$. On the other hand, humans prefer to choose a 50% chance to avoid a \$1000 loss rather than accepting a certain loss of \$490 with higher expected utility. As shown in Figure 7 (right), the prospect biased utility of a \$490 loss is lower than a half of the \$1000 loss, i.e., $v(-\$490) < 0.5v(-\$1000) \rightarrow 100\% \times v(-\$490) < 50\% \times v(-\$1000)$.

Note that loss aversion results in both risk averse behavior in gains and risk seeking

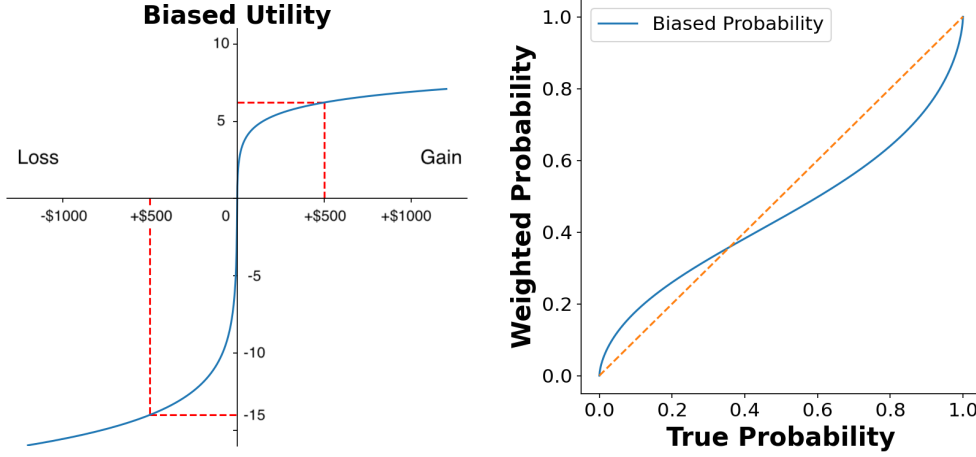


Figure 7. (*left*) The prospect biased utility function. The x-axis represents the amount of gain/loss and the y-axis is the corresponding weighted value. (*right*) The prospect biased probability function. The x-axis is the true probability and the y-axis shows the weighted probability.

behavior in losses. This differs from the risk-sensitive RL approaches (Mihatsch and Neuneier (2002); Vyetrenko and Xu (2019)) that set a risk preference parameter for the agent to be either risk-averse, risk-neutral, or risk-seeking, as well as other economic analyses that commonly only assume risk aversion.

Inspired by Chan, Critch, and Dragan (2021), we modify the reward function in the Bellman equation to $V(s) = \max_a \sum_{s' \in S} P(s'|s, a) (f_c(R(s, a, s')) + \gamma V(s'))$, where

$$f_c(r) = \begin{cases} \log(1 + |r|) & \text{if } r > 0 \\ 0 & \text{if } r = 0 \\ -c \log(1 + |r|) & \text{if } r < 0 \end{cases} \quad (10)$$

In the equation, c is the coefficient of loss aversion ($c > 1$ as human naturally feel losses greater than gains). While the agent is still risk-seeking in losses and risk-averse in gains regardless of the coefficient, an agent with larger c value will be more risk averse when comparing the same amount of potential gains and losses.

Prospect Biased Probability

Another aspect of prospect theory is that humans do not perceive probability in a linear fashion. When addressing any uncertainty, humans consider the natural boundaries of probability, impossibility and certainty, as the two reference points. The impact of a change in probability diminishes with an increase in its distance to the reference points. Considering a 0.1 increase in the probability of winning a prize, a change from 0 to 0.1 chance to win has more impact than a change from 0.45 to 0.55. To model this distortion in probability seen in humans, Tversky and Kahneman (1992) introduce a non-linear weighting function as follows,

$$w(p) = \frac{p^\delta}{(p^\delta + (1 - p)^\delta)^{1/\delta}} \quad (11)$$

where p is the probability of a potential outcome, and δ is estimated to be 0.61 and

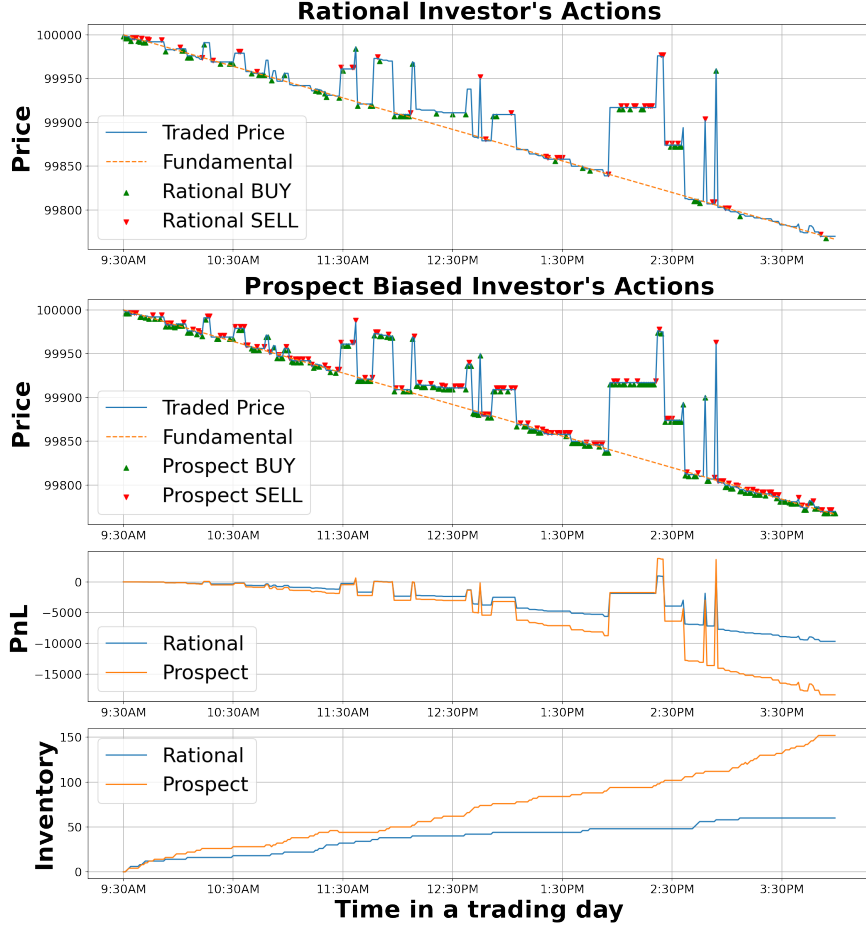


Figure 8. The behavior of prospect biased investors ($c = 2.5$ and $\delta = 0.65$) in the simulated market. When experiencing losses (negative PnL, see period around 2pm), the rational investor decides to hold more often and does not increase its inventory. In contrast, the prospect biased investor is more willing to take risks to avoid losses. They buy more shares with the hope of reducing their average purchase price, which allows them to avoid losses when the stock price temporarily increases around 2:20pm.

0.69 for gains and losses respectively (here for simplicity we set $\delta = 0.65$ for all types of outcomes). As shown in Figure 7 (right), humans overweight events with low probabilities and underweight events with high probabilities. We apply the above weighting function in the Bellman equation $V(s) = \max_a \sum_{s' \in S} w(P(s'|s, a)) (R(s, a, s') + \gamma V(s'))$ to model prospect biased probability in the RL setting.

An example of prospect biased trading behavior.

Figure 8 demonstrates how prospect biased behavior deviates from the optimal trading strategy. The market price continuously decreases in the entire day, and both investor agents are not able to gain by the end of the day. Since they are experiencing losses throughout the day, it allows us to investigate how rational and prospect biased investors handle losses differently. In the first half of the day, the rational and prospect biased investor have similar holding positions. After recognizing the long-term market trend, the rational investor stops trading frequently and does not increase the holding in the second half of the day. When the market becomes extremely volatile in a short time (between 1:50PM and 2:40PM), the prospect biased investor addresses

the uncertainty with a risk-seeking strategy. Due to the convex utility function Equation (10) for losses and over-weighted low probabilities (Equation (3)), they prefer to gamble on actions that can mitigate or avoid losses even if the expected utility is low and the success probability is small. As a result, they decide to buy more shares after the price drop, which reduces the average price at which the stock is purchased. This allows them to offset their losses once the stock return to the average purchased price. In fact, the prospect biased investor is able to achieve positive PnL during 2:20 pm. However, such strategy increases the overall risk exposure, and the PnL is highly volatile.

In summary, our model successfully capture human bias described in the prospect theory literature (Kahneman and Tversky (1979)). In the handcrafted market scenario, the prospect biased investor tend to be risk-seeking when facing choices between losses.

5.2.3. Optimism/Pessimism

Humans may have biased expectations of the future. Sharot et al. (2007) show that humans often systematically overestimate or underestimate their chances of experiencing positive and negative events. Inspired by Chan, Critch, and Dragan (2021), we model the optimistic and pessimistic humans by modifying the Bellman equation to:

$$V(s) = \max_a \sum_{s' \in S} P^\omega(s'|s, a) (R(s, a, s') + \gamma V(s')) \quad (12)$$

$$\text{where } P^\omega(s'|s, a) = \frac{P(s'|s, a)e^{\omega(R(s, a, s') + \gamma V(s'))}}{\sum_{s' \in S} P(s'|s, a)e^{\omega(R(s, a, s') + \gamma V(s'))}}. \quad (13)$$

When $\omega = 0$, Equation (12) considers the original transition probability and the agent is rational. As $\omega \rightarrow +\infty$ (or $\omega \rightarrow -\infty$), the agent becomes optimistic (pessimistic) by overestimating the probability of positive (negative) transitions.

An example of optimistic/pessimistic trading behavior

We illustrate the behavior of optimistic and pessimistic investors in simulated markets by adopting a fundamental price which follows a sine wave as shown in Figure 9. Similar to observations in Figure 4, the rational investor takes a optimal strategy by selling at the higher prices (first half of the day) and buying at the lower prices (second half of the day). The optimistic agent has a more aggressive strategy compared to the fully rational investor. While they have similar holding positions between 9:30am and 12:00pm, the optimistic investor builds a strong belief that the market trend will remain the same and continues short selling even when the price is lower than the starting price. The portfolio value of the optimistic investor is extremely volatile due to their over-confident trading strategy: they obtain massive gains when they are right (from 12:00pm to 2:30pm) and massive losses when they are wrong (from 2:30pm to 4:00pm). On the other hand, the pessimistic investor trades very little as can be seen by the large number of hold decisions due to their under-confidence, and do not have any hope of profitable investments. Additionally, the pessimistic investor has a small inventory and is less affected by the market trends.

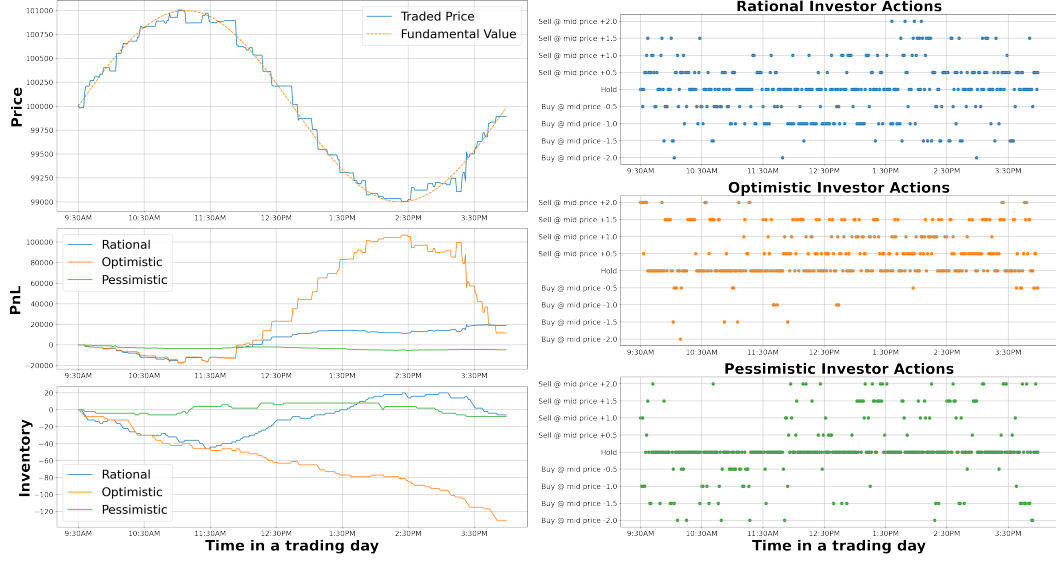


Figure 9. The behavior of optimistic and pessimistic investors ($\omega = 1$ and $\omega = -1$ respectively) in the simulated market. (left) The market prices and agents’ portfolio values and inventories. (right) The actions of the rational, optimistic, and pessimistic traders. Compared to the rational agent (blue), the optimistic agent (orange) tends to place orders aggressively on one side of the LOB, while the pessimistic (green) agent prefer to hold in most of the time.

6. Results

In this section, we examine the impact of sub-rationality of human strategies on: (1) their daily profit and loss (PnL) in simulated markets, (2) the importance given to various state features through an explainability analysis of their policy, and (3) market observables. We train and test the RL investors in markets simulated with ABIDES-gym using the same configuration described in Section 4.2.

6.1. Profit and Loss of Sub-Rational Investors

We first assess the profit and loss (PnL) of the sub-rational human investors in the simulated markets. In particular, we test each type of sub-rational trading strategy over 100 simulated trading days and show the distribution of the single-day PnL in Figure 10 and Table 4. Formally, we define the single-day PnL as the change in portfolio value by the end of a trading day (390 minutes in total from 9:30am to 4:00pm) in cents:

$$\text{PnL} = \text{PortfolioValue}_{\text{end}} - \text{PortfolioValue}_{\text{start}}$$

Table 4. The mean and standard deviation of the PnL of sub-rational human investors.

Metrics	Rational	Prospect Biased $c = 2.5$ $\delta = 0.65$	Bounded (zero-intelligence) $\beta = 0$	Bounded $\beta = 1$	Myopic (fully-myopic) $\gamma = 0$	Myopic $\gamma = 0.8$	Pessimistic $\omega = -1$	Optimistic $\omega = 1$
Mean	1060.42	117.3	29.75	1056.35	162.76	575.98	174.11	454.93
Std	1835.58	3740.26	2711.04	1927.67	867.96	1732.98	561.98	5487.65

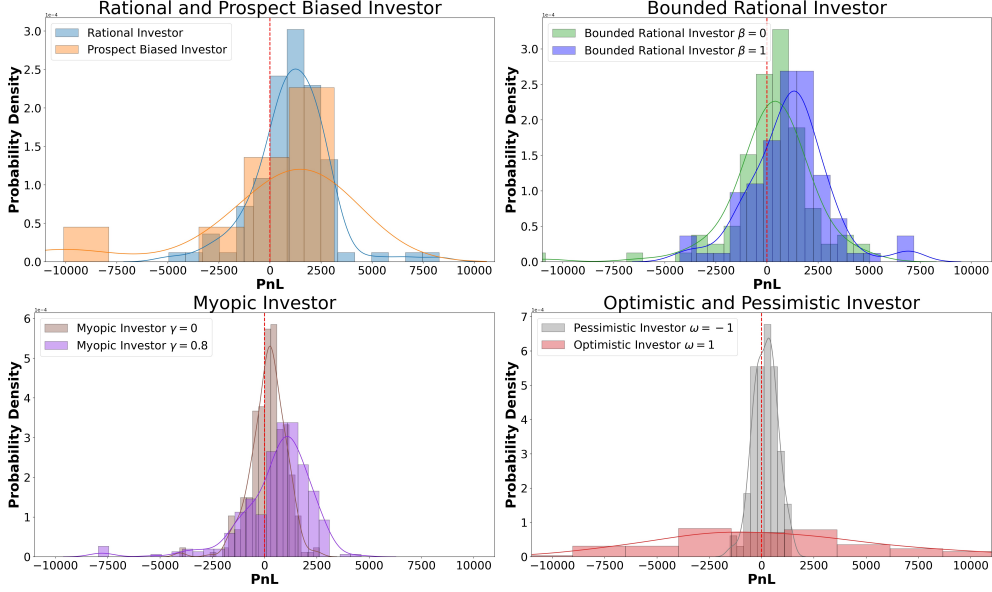


Figure 10. The PnL of sub-rational human investors in simulated markets.

We observe that the rational investor has the best overall performance in simulated markets. Their PnL follows a normal distribution with a mean of 1060.42 and standard deviation of 1835.58. Compared to the rational investor, the prospect biased investor (with $c = 2.5$ and $\delta = 0.65$) has lower average daily PnL (117.3) and larger variance ($std = 3740.26$). The PnL distribution of the prospect biased investor is also heavily left skewed: there is no gain larger than 3000 and there is a substantial amount of losses around -10000. This is because the prospect biased humans are risk-averse in gains and risk-seeking in losses, which leads to small variance in the positive domain and large variance in the negative domain. Since β in the bounded rationality model Equation (6) decides how well the agent optimizes their objective function, the PnL increases as β increases. We observe that bounded rational agent with zero intelligence ($\beta = 0$) has the overall worst average daily PnL = 29.75 and high variance. When β increases to 1, the human becomes more rational: the average PnL approaches that of the rational investor and the variance reduces. The myopic investor ($\gamma = 0$), on the other hand, has smaller variance (867.96) compared to the fully rational investor, but worse average daily PnL (162.76). The small variance of myopic investor PnL comes from the fact that they only make short-term investments through the trading day, which result in gains and losses in a smaller magnitude compared to long-term investments. The PnL of the myopic investor becomes similar to the fully rational investor when $\gamma = 0.8$. The pessimistic investor ($\omega = -1$) has the smallest variance in the PnL distribution (561.98) but low average daily PnL (174.11). This is because they are extremely cautious in investing and often hold a small inventory. Therefore, they are unlikely to have large gains or losses by the end of the day. However, the optimistic investor has the largest variance in the PnL distribution due to their over-confident and aggressive trading strategy: they can win a large reward if the market goes according to their expectation, or lose a massive amount of money if they made a wrong expectation.

Summary. We observe that different sub-rational trading strategies have unique daily-PnL distributions due to the noise and bias in the decision-making process. Our experiments also demonstrate that market simulation is an efficient tool to test and evaluate

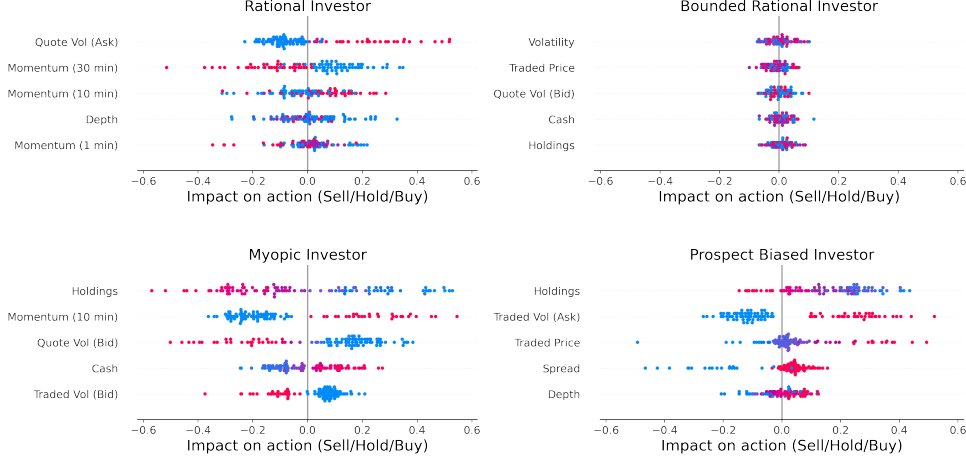


Figure 11. The impact of state features on the actions of investor (Sell/Hold/Buy) shown for the top 5 features. The color of the scatter point represents the feature value: blue \rightarrow low and red \rightarrow high.

various trading strategies in a controlled environment.

6.2. Investigation of Sub-rational Human Decision-making

To explain the decision-making process of sub-rational human investors, we perform an explainability analysis of the RL policy network that maps the current state to actions of the investor. We use the explainability tool called SHapley Additive exPlanations (or SHAP in short) to score the different state features based on their importance towards determining the trader’s actions (Lundberg and Lee (2017)). These scores, called SHAP values, are computed using cooperative game theory by decomposing the network output locally into a sum of effects attributed to each input feature. Hence, every (state, action) pair has a vector of SHAP values with elements corresponding to each state feature, and that accumulate to give the action per action dimension.

We compute SHAP values for every entry in a dataset comprising (state, action) pairs for the different sub-rational investor policies. Figure 11 shows the distribution of SHAP values for each state feature towards determining the order direction, ordered by the corresponding feature importance (sum of absolute values of SHAP values). The x-axis indicates the impact on the model output ($-1 \rightarrow$ sell, $0 \rightarrow$ hold, $1 \rightarrow$ buy) while the color represents the observed value of the state (blue \rightarrow low, red \rightarrow high). For example, the 30-minute momentum is one of the most impactful state features for the fully rational policy. If the momentum is low (blue), it has a positive impact on the investor’s action and they place a buy order, and vice-versa when the momentum is high (red). So, the rational agent will place buy/sell orders if they observe a strong decrease/increase in the past 30 minutes, indicating that the current price is at a low/high position. The short-term momentum is less important to the rational investor, especially for the 1-minute momentum which shows a hardly distinguishable SHAP value distribution.

For the zero-intelligence bounded rational investor ($\beta = 0$), we do not observe a distinguishable contribution of the state features to the agent’s actions. The distribution of SHAP values is not correlated with the value of the state features. Compared to the rational investor, the myopic investor ($\gamma = 0$) focuses more on the short-term market observations, as the feature importance of 10-minute momentum is higher than

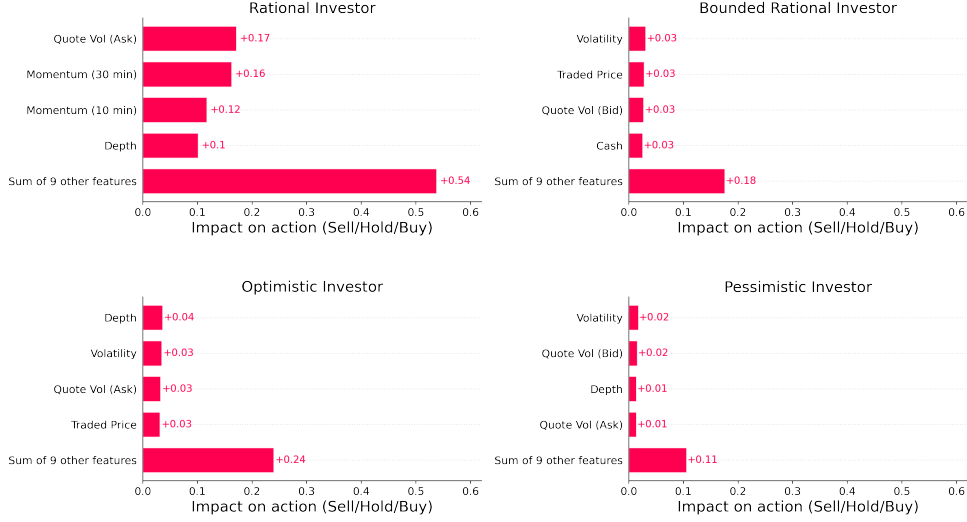


Figure 12. Importance of state features towards actions of the sub-rational investor as measured by the average of absolute SHAP values. State features are ordered in decreasing order of their importance on the y-axis. Human investors with bounded rationality or extreme internal beliefs (optimistic/pessimistic) pay less attention to market observations when making decisions.

the 30-minute momentum. They are more careful on their current portfolio situation (holding, cash) opting to sell when holdings are high, and buy when holdings are low, which explains the myopic loss aversion behavior. The SHAP-value analysis of the prospect biased investors shows that holding is the most important state feature for their decision-making. This indicates that the prospect biased investors are very sensitive to the inventory risk (Avellaneda and Stoikov (2008)). When experiencing gains in the market, they incline to reduce the inventory if they have a large amount of holdings.

We also notice that investors with bounded rationality ($\beta = 0$), pessimistic bias ($\omega = -1$), or optimistic bias ($\omega = 1$) do not build their decisions based on any market observations. As shown in Figure 12, none of the state features gives average SHAP values greater than 0.05 for bounded rational, pessimistic, and optimistic investors, which suggests that they are not strong contributing factors for the investors' decisions. **Summary.** The SHAP value analysis of the state-action relationship gives an explicit explanation of the sub-rational human strategy. We observe that the zero-intelligence bounded rational investors as well as those with extreme internal bias (optimism/pessimism) do not make decisions based on their observations of the markets, the myopic investors pay more attention to the short-term momentum, and the decisions of prospect biased investors strongly base decisions on their inventory – which is aligned with one's expectations about these types of sub-rational behaviors.

6.3. Impact of Sub-Rational Investors on Market Observables

Here, we investigate how sub-rational trading behaviors affect the financial market. In particular, we examine and compare the statistics of simulated markets with different types of human investors while the background agents configuration (see Table 2) remains the same. We perform this exercise for seven different market settings: no RL agents, rational (10 rational RL agents), bounded-rational (10 RL agents with $\beta = 0$),

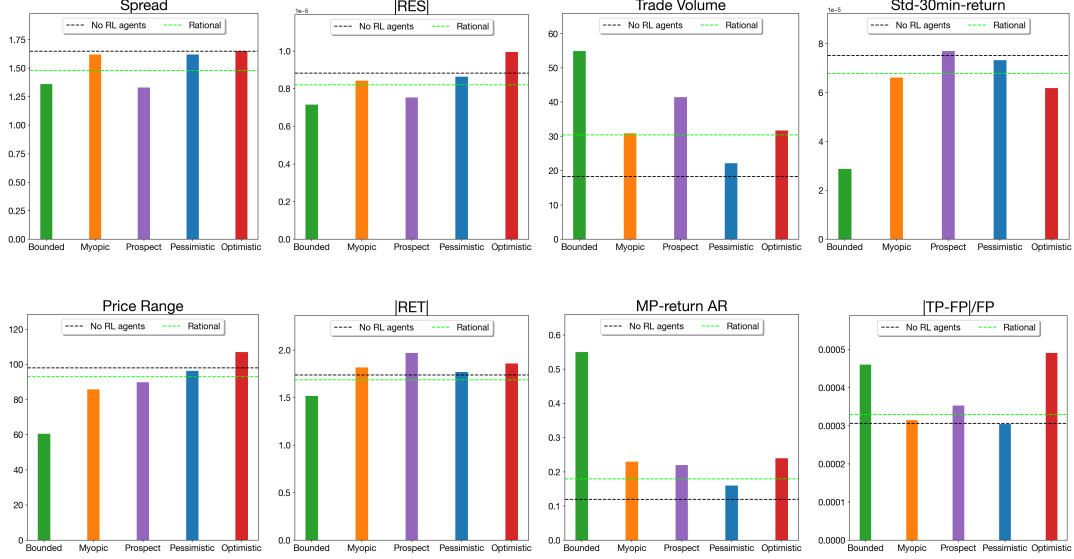


Figure 13. Impact of different types of RL investors on market quality. As mentioned in Section 3.3, we evaluate market quality on three aspects: liquidity (spread, $|RES|$, trader volume), volatility (std-30min-return, price range, $|RET|$), and efficiency (MP-return AR, $|TP-FP|$).

psychologically myopic (10 RL agents with $\gamma = 0.2$), prospect biased (10 RL agents with $c = 2.5$ and $\delta = 0.65$), optimistic (10 RL agents with $\omega = 1$), and pessimistic (10 RL agents with $\omega = -1$). Figure 13 shows average metrics per trading day over 100 simulations for each market setting described earlier in this section. The simulation random seeds are different per simulation run within each setting and same across different settings. We evaluate market quality on dimensions of liquidity, volatility, and information efficiency as described in Section 3.3. For all metrics except trade volume, smaller values indicate better market quality. We observe that market liquidity almost always improves with the addition of RL investors except in the optimistic setting for $|RES|$.

Bounded rational investor. We observe that the bounded rational investor leads to markets with better liquidity and volatility, but worse information efficiency. Compared to simulated markets with rational investors (rational) and without any RL investor, markets with bounded rational investors have larger trade volume and smaller spread, $|RES|$, Std-30min-RET, price range, and $|RET|$. This arises from the fact that a collection of zero-intelligence bounded rational investors essentially mirrors the functions of a market maker Venkataraman and Waisburd (2007). While a collection of bounded rational investors tend to place orders uniformly at random on both sides of the LOB due their limited computational capacity, market makers are required by regulation to place orders on both sides regularly. Additionally, as illustrated in Figure 14, bounded rational investors exhibit a lower likelihood of placing orders on the side of the LOB with larger quoted volume, and the orders they execute are in proximity to the mid-price. Accordingly, due to their limited utilization of market information as seen in Figure 11, the bounded rational investors noticeably diminish market efficiency.

Myopic investor. As we discovered in the SHAP analysis (Figure 11), the myopic investor’s decisions tend to follow the short-term market momentum. As a result, the myopic investor typically positions orders on the side of the LOB with larger

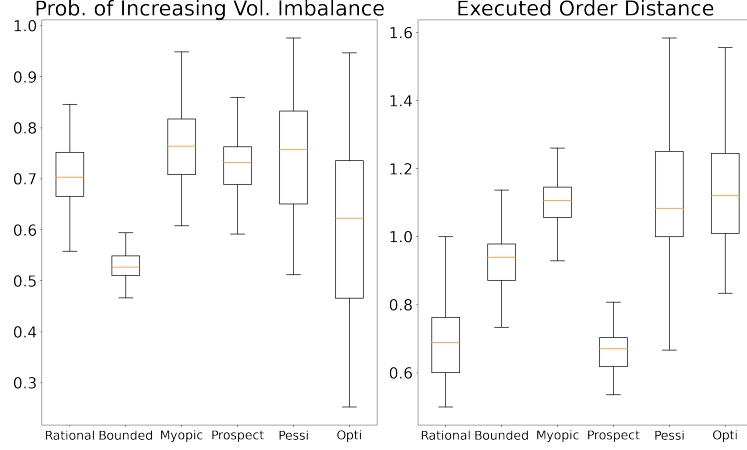


Figure 14. The distribution of agents’ actions in simulated markets. (left) The probability of the agent placing orders on the heavier side of the LOB, which increases the volume imbalance. (right) The distance between the price of executed orders and the mid price.

quoted volume, thereby increasing the LOB imbalance. Our observations in Figure 13 indicate that myopic trading behavior has negative impacts on market liquidity when contrasted with rational behavior. However, it enhances market efficiency with regard to the $|TP-FP| / FP$ since it elevates the propagation of fundamental information by following the market momentum. Meanwhile, there is no significant discernible difference in the impact of myopic investors on volatility compared to their rational counterparts.

Prospect biased investor. The prospect biased investors significantly improve market liquidity by narrowing the spread and $|RES|$, and increasing the trade volume. A major reason being that among all types of RL investors, the prospect biased investor has executed orders closest to the mid price as seen in Figure 14. This is a result of them placing orders close to the mid-price so as to reduce any uncertainty in order execution. This subsequently reduces the spread and facilitates trade executions.

Pessimistic investor. Since pessimistic investors are not willing to trade most of the time, simulated markets with pessimistic investors have less liquidity than those with other types RL investors (excepting the optimistic). On the other hand, they do not exhibit strong influences on the market, resulting in the market statistics being very close to those of markets without any RL investors.

Optimistic investor. As shown in Figure 12, optimistic investors rely more on their internal beliefs than market variables, which undermines market efficiency. The aggressive trading strategy of optimistic investors also leads to simulated markets with largest spread, price range and absolute daily returns $|RET|$. As illustrated in Figure 9, the actions of optimistic investors are often heavily skewed to buy (or sell). Therefore, they tend to increase the LOB imbalance and reduce the liquidity of the market.

Summary. By comparing with market simulations without any RL investors, we show that bounded rational and prospect biased investors improve liquidity, but reduce efficiency; myopic investors reduce liquidity but improve efficiency; optimistic investors increase volatility and reduce efficiency; and pessimistic investors do not have strong impacts on the market. All of our aforementioned findings on the different sub-rationality models are summarized in Table 1.

7. Discussion

One of the largest obstacles in studying sub-rational trading behavior is the scarcity of historical trading data at an individual level, which is usually not public-available (Gutiérrez-Roig et al. (2019); Cont et al. (2023)). While several studies have collected sub-rational behavior from controlled lab experiments on human subjects (Mischel and Ebbesen (1970); Rothblum, Solomon, and Murakami (1986); Green, Myerson, and McFadden (1997)), such data are very limited at scale and may raise ethical and security concerns in the financial domain. Even with access to labeled trading data, factorizing the impact of human traders on complex market environments remains a challenge. Recent studies have employed large language models (LLM) to generate synthetic human demonstration data Coletta et al. (2024), which have not been successfully applied for generating complex investing decisions in financial markets. As a result, there lacks studies which can identify sub-rational trading behaviors and evaluate their impacts on financial markets. Among the few that exist, Olsen (1998) gives behavioral finance explanations of market volatility, and Levy, Levy, and Solomon (2000) discover that human traders influenced by prospect bias generate price deviation, high traded volume, and excess volatility. Our experimental results in Figure 13 align with their findings. By utilizing RL and market simulations, our model can generate sub-rational trading decisions with or without using ground truth historical data, which liberates research from the scarcity of sub-rational trading data.

While we provide a mathematical framework to capture various aspects of human sub-rationality, it is crucial to specify the sub-rationality parameters for real humans (e.g. β for bounded rational humans, or γ for myopic humans). A common approach is to estimate these parameters using human demonstration data. For example, Tversky and Kahneman (1992) used a nonlinear regression model to estimate the parameters in Equations (10) and (11) from human subject studies. Similarly, Green, Myerson, and McFadden (1997) collected decisions of human subjects and fit exponential and hyperbolic discounting factors (Equations (7) and (8)) using nonlinear regression. However, these sub-rationality parameters vary between humans with different degrees of bias and intelligence, making it expensive to conduct human studies more broadly. Consequently, existing works on estimating parameters for human models show disagreement among each other (Frederick, Loewenstein, and O’donoghue (2002)). Hence, we do not aim to estimate a single subrationality parameter for humans using trading data, and instead propose a holistic framework to model human investors with different types and degrees of sub-rationality. And, investigate the impact of such subrationality on trading behaviors and market environments.

8. Conclusion

In this work, we employ reinforcement learning (RL) techniques to model sub-rational human investors in financial markets. We consider sub-rational behavior as explained by two driving factors: limited information/computational capacity and biased internal beliefs. Our approach involves utilizing sophisticated market simulations for the training and evaluation of diverse human investor models. Additionally, we construct probabilistic neural networks to accurately represent the biased internal beliefs guiding human investors in their decision-making process.

In particular, we consider five sub-rational human investor models: (1) the bounded rational investors that make sub-optimal decisions due to limitations in information

access and computational power, (2) the psychologically myopic investors that maximize only short-term rewards regardless of the future, (3) the prospect biased investors that are risk-averse in gains and risk-seeking in losses, (4) the optimistic investors that have exaggerated hope in receiving positive outcomes, (5) and the pessimistic investors that tend to underestimate the likelihood of receiving positive outcomes over negative outcomes. We first craft specific market scenarios to assess and illustrate the behavior of each type of sub-rational human investors. We then examine the relationship between the Profit and Loss (PnL) of these investors and their respective computational limitations and psychological biases. Furthermore, we provide an intuitive exploration of the sub-rational human trading strategies through SHAP value analysis. We show that our models successfully reproduce sub-rational trading behavior identified in the previous literature. Upon populating simulated markets with different types of sub-rational human investors, we evaluate their impact on market quality in terms of liquidity, volatility, and price efficiency. Our findings indicate that a market populated with bounded rational or prospect biased investors has improved liquidity but worse price efficiency. Additionally, we observe that bounded rational investors contribute to a reduction in stock volatility. On the other hand, myopic investors show adverse impacts on liquidity but improve information efficiency. For the pessimistic trading behavior, we do not identify a significant impact on the market quality. Meanwhile, the optimistic investors cause negative impact on market liquidity and price efficiency. We note that our comparative study is different from previous regulatory work on algorithmic and high-frequency trading Boehmer, Fong, and Wu (2021a); Woodward (2017) in that we compare investors with different types of rationality, while not varying their trading frequencies in simulated markets.

A future direction of this work would be to consider other aspects of human sub-rationality. For example, economic studies show that investors make sub-optimal decisions due to overconfidence, which can be described as resulting from an illusion of knowledge or an illusion of control Barber and Odean (2002); Song et al. (2013). In addition, we are interested in estimating the degree of sub-rationality for real human investors, using historical data and (potentially) demonstration data generated by large language models as proxy human subjects Coletta et al. (2024).

Acknowledgements

This paper was prepared for informational purposes by the CDAO group of JPMorgan Chase & Co and its affiliates (“J.P. Morgan”) and is not a product of the Research Department of J.P. Morgan. J.P. 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

- Ainslie, George. 1975. "Specious reward: a behavioral theory of impulsiveness and impulse control." *Psychological bulletin* 82 (4): 463.
- Ainslie, George. 1992. *Picoeconomics: The strategic interaction of successive motivational states within the person*. Cambridge University Press.
- Alexander, William H, and Joshua W Brown. 2010. "Hyperbolically discounted temporal difference learning." *Neural computation* 22 (6): 1511–1527.
- Amrouni, Selim, Aymeric Moulin, Jared Vann, Svitlana Vyetrenko, Tucker Balch, and Manuela Veloso. 2021. "ABIDES-gym: gym environments for multi-agent discrete event simulation and application to financial markets." In *Proceedings of the Second ACM International Conference on AI in Finance*, 1–9.
- Asadi, Kavosh, and Michael L Littman. 2017. "An alternative softmax operator for reinforcement learning." In *International Conference on Machine Learning*, 243–252. PMLR. <http://proceedings.mlr.press/v70/asadi17a/asadi17a.pdf>.
- Avellaneda, Marco, and Sasha Stoikov. 2008. "High-frequency trading in a limit order book." *Quantitative Finance* 8 (3): 217–224.
- Baker, Chris, Rebecca Saxe, and Joshua Tenenbaum. 2005. "Bayesian models of human action understanding." *Advances in neural information processing systems* 18. <https://proceedings.neurips.cc/paper/2005/file/f5b1b89d98b7286673128a5fb112cb9a-Paper.pdf>.
- Balch, Tucker Hybinette, Mahmoud Mahfouz, Joshua Lockhart, Maria Hybinette, and David Byrd. 2019. "How to evaluate trading strategies: Single agent market replay or multiple agent interactive simulation?" *arXiv preprint arXiv:1906.12010*.
- Barber, Brad M, and Terrance Odean. 2002. "Online investors: do the slow die first?" *The Review of financial studies* 15 (2): 455–488.
- Benartzi, Shlomo, and Richard H Thaler. 1995. "Myopic loss aversion and the equity premium puzzle." *The quarterly journal of Economics* 110 (1): 73–92.
- Benartzi, Shlomo, and Richard H Thaler. 1999. "Risk aversion or myopia? Choices in repeated gambles and retirement investments." *Management science* 45 (3): 364–381.
- Boehmer, Ekkehart, Kingsley Fong, and Juan (Julie) Wu. 2021a. "Algorithmic Trading and Market Quality: International Evidence." *Journal of Financial and Quantitative Analysis* 56 (8): 2659–2688. <https://doi.org/10.1017/S0022109020000782>.
- Boehmer, Ekkehart, Kingsley Fong, and Juan Julie Wu. 2021b. "Algorithmic trading and market quality: International evidence." *Journal of Financial and Quantitative Analysis* 56 (8): 2659–2688.
- Boehmer, Ekkehart, and Eric K Kelley. 2009. "Institutional investors and the informational efficiency of prices." *The Review of Financial Studies* 22 (9): 3563–3594.
- Bouchaud, Jean-Philippe, Marc Mézard, and Marc Potters. 2002. "Statistical properties of stock order books: empirical results and models." *Quantitative finance* 2 (4): 251.
- Byrd, David, Maria Hybinette, and Tucker Hybinette Balch. 2019. "Abides: Towards high-fidelity market simulation for ai research." *arXiv preprint arXiv:1904.12066*.
- Chabris, Christopher F, David I Laibson, and Jonathon P Schuldt. 2010. "Intertemporal choice." In *Behavioural and experimental economics*, 168–177. Springer.
- Chan, Lawrence, Andrew Critch, and Anca Dragan. 2021. "Human irrationality: both bad and good for reward inference." *arXiv preprint arXiv:2111.06956*.
- Chordia, Tarun, Richard Roll, and Avanidhar Subrahmanyam. 2002. "Order imbalance, liquidity, and market returns." *Journal of Financial economics* 65 (1): 111–130.
- Chrisman, James J, and Pankaj C Patel. 2012. "Variations in R&D investments of family and nonfamily firms: Behavioral agency and myopic loss aversion perspectives." *Academy of management Journal* 55 (4): 976–997.
- Coletta, Andrea, Kshama Dwarakanath, Penghang Liu, Svitlana Vyetrenko, and Tucker Balch. 2024. "LLM-driven Imitation of Subrational Behavior: Illusion or Reality?" *arXiv preprint arXiv:2402.08755*.

- Coletta, Andrea, Aymeric Moulin, Svitlana Vyetenko, and Tucker Balch. 2022. “Learning to simulate realistic limit order book markets from data as a World Agent.” In *Proceedings of the Third ACM International Conference on AI in Finance*, 428–436.
- Cont, Rama, Mihai Cucuringu, Vacslav Glukhov, and Felix Prezel. 2023. “Analysis and modeling of client order flow in limit order markets.” *Quantitative Finance* 23 (2): 187–205.
- Dwarakanath, Kshama, Svitlana S Vyetenko, and Tucker Balch. 2021. “Profit equitably: an investigation of market maker’s impact on equitable outcomes.” In *Proceedings of the Second ACM International Conference on AI in Finance*, 1–8.
- Evans, Owain, and Noah D Goodman. 2015. “Learning the preferences of bounded agents.” In *NIPS Workshop on Bounded Optimality*, Vol. 6, 2–1.
- Evans, Owain, Andreas Stuhlmüller, and Noah Goodman. 2016. “Learning the preferences of ignorant, inconsistent agents.” In *Thirtieth AAAI Conference on Artificial Intelligence*, .
- Fama, Eugene F. 1970. “Efficient capital markets: A review of theory and empirical work.” *The journal of Finance* 25 (2): 383–417.
- Fedus, William, Carles Gelada, Yoshua Bengio, Marc G Bellemare, and Hugo Larochelle. 2019. “Hyperbolic discounting and learning over multiple horizons.” *arXiv preprint arXiv:1902.06865* .
- Frederick, Shane, George Loewenstein, and Ted O’donoghue. 2002. “Time discounting and time preference: A critical review.” *Journal of economic literature* 40 (2): 351–401.
- Friedman, Daniel. 2018. “The double auction market institution: A survey.” In *The Double Auction Market Institutions, Theories, and Evidence*, 3–26. Routledge.
- Gould, Martin D, Mason A Porter, Stacy Williams, Mark McDonald, Daniel J Fenn, and Sam D Howison. 2013. “Limit order books.” *Quantitative Finance* 13 (11): 1709–1742.
- Green, Leonard, Nathanael Fristoe, and Joel Myerson. 1994. “Temporal discounting and preference reversals in choice between delayed outcomes.” *Psychonomic Bulletin & Review* 1: 383–389.
- Green, Leonard, Joel Myerson, and Edward McFadden. 1997. “Rate of temporal discounting decreases with amount of reward.” *Memory & cognition* 25: 715–723.
- Grüne-Yanoff, Till. 2015. “Models of temporal discounting 1937–2000: An interdisciplinary exchange between economics and psychology.” *Science in context* 28 (4): 675–713.
- Gutiérrez-Roig, Mario, Javier Borge-Holthoefer, Alex Arenas, and Josep Perelló. 2019. “Mapping individual behavior in financial markets: synchronization and anticipation.” *EPJ Data Science* 8 (1): 1–18.
- Haarnoja, Tuomas, Aurick Zhou, Kristian Hartikainen, George Tucker, Sehoon Ha, Jie Tan, Vikash Kumar, et al. 2018. “Soft actor-critic algorithms and applications.” *arXiv preprint arXiv:1812.05905* .
- Janner, Michael, Justin Fu, Marvin Zhang, and Sergey Levine. 2019. “When to trust your model: Model-based policy optimization.” *Advances in neural information processing systems* 32.
- Kahneman, Daniel, and Amos Tversky. 1979. “Prospect Theory: An Analysis of Decision under Risk.” *Econometrica* 47 (2): 263–292.
- Laidlaw, Cassidy, and Anca Dragan. 2022. “The Boltzmann Policy Distribution: Accounting for Systematic Suboptimality in Human Models.” <https://arxiv.org/abs/2204.10759>.
- Larson, Francis, John A List, and Robert D Metcalfe. 2016. *Can myopic loss aversion explain the equity premium puzzle? Evidence from a natural field experiment with professional traders*. Technical Report. National Bureau of Economic Research.
- Levy, Haim, Moshe Levy, and Sorin Solomon. 2000. *Microscopic simulation of financial markets: from investor behavior to market phenomena*. Elsevier.
- Lundberg, Scott M, and Su-In Lee. 2017. “A unified approach to interpreting model predictions.” *Advances in neural information processing systems* 30.
- Lux, Thomas, and Michele Marchesi. 1999. “Scaling and criticality in a stochastic multi-agent model of a financial market.” *Nature* 397 (6719): 498–500.
- Malkiel, Burton G. 1989. “Efficient market hypothesis.” In *Finance*, 127–134. Springer.
- Mazur, James E. 1985. “Probability and delay of reinforcement as factors in discrete-trial

- choice.” *Journal of the Experimental Analysis of Behavior* 43 (3): 341–351.
- Mihatsch, Oliver, and Ralph Neuneier. 2002. “Risk-sensitive reinforcement learning.” *Machine learning* 49: 267–290.
- Mischel, Walter, and Ebbe B Ebbesen. 1970. “Attention in delay of gratification.” *Journal of personality and social psychology* 16 (2): 329.
- Mnih, Volodymyr, Koray Kavukcuoglu, David Silver, Alex Graves, Ioannis Antonoglou, Daan Wierstra, and Martin Riedmiller. 2013. “Playing atari with deep reinforcement learning.” *arXiv preprint arXiv:1312.5602*.
- Nevmyvaka, Yuriy, Yi Feng, and Michael Kearns. 2006. “Reinforcement Learning for Optimized Trade Execution.” ICML ’06, New York, NY, USA, 673–680. Association for Computing Machinery. <https://doi.org/10.1145/1143844.1143929>.
- Olsen, Robert A. 1998. “Behavioral finance and its implications for stock-price volatility.” *Financial analysts journal* 54 (2): 10–18.
- Raja, Anita, and Victor Lesser. 2001. “Towards Bounded-Rationality in Multi-Agent Systems: A Reinforcement-Learning Based Approach.” *University of Massachusetts Computer Science Technical Report* 34: 2001.
- Reddy, Sid, Anca Dragan, and Sergey Levine. 2018. “Where do you think you’re going?: Inferring beliefs about dynamics from behavior.” *Advances in Neural Information Processing Systems* 31.
- Redish, A David, and Zeb Kurth-Nelson. 2010. “Neural models of delay discounting.”
- Roşu, Ioanid. 2009. “A dynamic model of the limit order book.” *The Review of Financial Studies* 22 (11): 4601–4641.
- Rothblum, Esther D, Laura J Solomon, and Janice Murakami. 1986. “Affective, cognitive, and behavioral differences between high and low procrastinators.” *Journal of counseling psychology* 33 (4): 387.
- Samuelson, Paul A. 1937. “A note on measurement of utility.” *The review of economic studies* 4 (2): 155–161.
- Sharot, Tali, Alison M Riccardi, Candace M Raio, and Elizabeth A Phelps. 2007. “Neural mechanisms mediating optimism bias.” *Nature* 450 (7166): 102–105.
- Simon, Herbert A. 1955. “A behavioral model of rational choice.” *The quarterly journal of economics* 99–118.
- Simon, Herbert A. 1990. “Bounded rationality.” In *Utility and probability*, 15–18. Springer.
- Simon, Herbert Alexander. 1997. *Models of bounded rationality: Empirically grounded economic reason*. Vol. 3. MIT press.
- Song, Reo, Sungha Jang, Dominique Hanssens, and Jaebeom Suh. 2013. “When Overconfidence Meets Reinforcement Learning.”
- Spooner, Thomas, John Fearnley, Rahul Savani, and Andreas Koukorinis. 2018. “Market making via reinforcement learning.” *arXiv preprint arXiv:1804.04216*.
- Stoll, Hans R. 2003. “Market microstructure.” In *Handbook of the Economics of Finance*, Vol. 1, 553–604. Elsevier.
- Thaler, Richard H. 2016. “Behavioral economics: Past, present, and future.” *American economic review* 106 (7): 1577–1600.
- Thaler, Richard H, and Hersh M Shefrin. 1981. “An economic theory of self-control.” *Journal of political Economy* 89 (2): 392–406.
- Thaler, Richard H, Amos Tversky, Daniel Kahneman, and Alan Schwartz. 1997. “The effect of myopia and loss aversion on risk taking: An experimental test.” *The quarterly journal of economics* 112 (2): 647–661.
- Tversky, Amos, and Daniel Kahneman. 1992. “Advances in prospect theory: Cumulative representation of uncertainty.” *Journal of Risk and uncertainty* 5: 297–323.
- Venkataraman, Kumar, and Andrew C Waisburd. 2007. “The value of the designated market maker.” *Journal of Financial and Quantitative Analysis* 42 (3): 735–758.
- Vyetrenko, Svitlana, David Byrd, Nick Petosa, Mahmoud Mahfouz, Danial Dervovic, Manuela Veloso, and Tucker Balch. 2020. “Get real: Realism metrics for robust limit order book market simulations.” In *Proceedings of the First ACM International Conference on AI in*

- Finance*, 1–8.
- Vyetrenko, Svitlana, and Shaojie Xu. 2019. “Risk-sensitive compact decision trees for autonomous execution in presence of simulated market response.” *arXiv preprint arXiv:1906.02312* .
- Wang, Xintong, Christopher Hoang, Yevgeniy Vorobeychik, and Michael P. Wellman. 2021. “Spoofing the Limit Order Book: A Strategic Agent-Based Analysis.” *Games* 12 (2). <https://doi.org/10.3390/g12020046>, <https://www.mdpi.com/2073-4336/12/2/46>.
- Woodward, Megan. 2017. “The need for speed: regulatory approaches to high frequency trading in the United States and the European Union.” *Vand. J. Transnat’l L.* 50: 1359.
- Zheng, Ban, Eric Moulines, and Frédéric Abergel. 2012. “Price Jump Prediction in Limit Order Book.” .
- Ziebart, Brian D, J Andrew Bagnell, and Anind K Dey. 2010. “Modeling interaction via the principle of maximum causal entropy.” .

9. Appendices

9.1. Performance of Internal Belief Models

As illustrated in Section 4.3, we use probabilistic neural networks (PNN) as the internal model to incorporate biased human beliefs. To evaluate the performance of the internal model, we generate sequences of random trading actions and feed them to the true environment (ABIDES) and internal model (PNN) respectively. Figure 15 shows the distributions of 30-minute traded price volatility and the step reward with mean and kernel density estimation curve. In addition, Table 5 shows the Earth Mover’s Distance (EMD, also known as Wasserstein distance) and the root mean square error (RMSE) between the internal model prediction and the true environment outputs. The results show that the internal model can be considered as an alternative of the true environment for injecting the biased human beliefs. Note that the error for spread and volume is slightly higher. This discrepancy may arise from the fact that PNN generates the predictions from Gaussian distributions, which might not be an ideal fit for variables following heavy-tailed distributions. As a part of our future work, we intend to enhance this aspect by developing a more sophisticated world model.

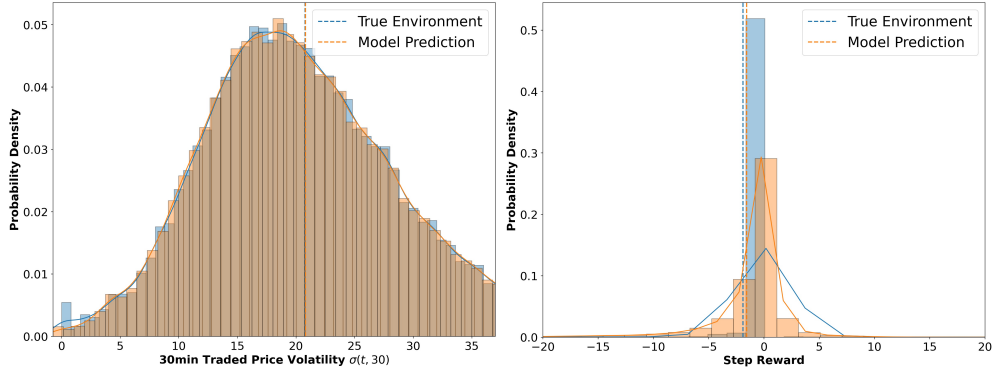


Figure 15. The distributions of 30-minute traded price volatility and step rewards from the internal model and the true environment.

Table 5. The distribution of the internal model vs. the true environment. Each variable is normalized by the min-max values of the combined distributions of PNN and the true environment.

Variable	EMD	RMSE
Quote Volume _{t+1}	0.159	0.229
Spread _{t+1}	0.169	0.242
Depth _{t+1}	0.064	0.102
Holdings _{t+1}	0.005	0.001
Cash _{t+1}	0.005	0.001
Trade Volume _{t+1}	0.113	0.178
Traded Price _{t+1}	0.069	0.098
Momentum($t + 1, 30$)	0.141	0.177
Volatility $\sigma(t + 1, 30)$	0.001	0.002
Reward $R(s_t, s_{t+1})$	0.0082	0.0004