# Deep Reinforcement Learning for Portfolio Management

Gang Huang, Xiaohua Zhou and Qingyang Song

*Huang, rockingapple@cqu.edu.cn, Phone number: +8615723205XXX Department of Applied Economics, Chongqing University; Zhou, zhxiaoh@aliyun.com, Department of Applied Economics, Chongqing University; and Song, song211@126.com, School of Big Data and Software, Chongqing University.

**Abstract**

In our paper, we apply deep reinforcement learning approach to optimize investment decisions in portfolio management. We make several innovations, such as adding short mechanism and designing an arbitrage mechanism, and applied our model to make decision optimization for several randomly selected portfolios. The experimental results show that our model is able to optimize investment decisions and has the ability to obtain excess return in stock market, and the optimized agent maintains the asset weights at fixed value throughout the trading periods and trades at a very low transaction cost rate. In addition, we redesigned the formula for calculating portfolio asset weights in continuous trading process which can make leverage trading, that fills the theoretical gap in the calculation of portfolio weights when shorting.

**Key words**：Data driven, Artificial intelligence, Deep reinforcement learning, Portfolio management, Investment decision optimization

## 1.    INTRODUCTION

In contemporary times, multi-factor models are the effective way to get excess return in stock market. However, as the study for multi-factor model progressed, the drawbacks of the factor model became apparent. For example, Hou *et al.* (2020) find 65% of the 452 anomalies cannot clear the single test hurdle of the absolute t-value of 1.96, if the t-value was increased to 2.78, the failure rate will raise to 82%, indicating that the validity of most of the factors had diminished or failed. Chen & Pearl (2013) pointed out that textbook failed to provide coherent mathematical notation that distinguishes causal from statistical concepts. Colquhoun (2014) even indicated that the use of p-values for significance testing is completely wrong, and the findings of this study indirectly suggest that there are irrationalities in the modeling of multi-factor model. In addition, lots of evidence has shown that traditional econometric analysis is not proper approach (Laloux *et al.* 1999; Harvey *et al.* 2016).

We can use a data-driven approach to address these flaws that traditional econometric method has. The traditional method is to first set a relational equation, for instance, the factor model would be set beforehand that asset returns and factors are linearly related to each other and then the parameters of the model are estimated from the overall sample. However, in the data-

driven concept, we look for rules directly from data. Mullainathan & Spiess (2017) suggest: "The real breakthrough came once we stopped trying to deduce these rules. Instead, the problem was turned into an inductive one: rather than hand-curating the rules, we simply let the data tell us which rules work best."

The idea of data-driven has been proposed in the 1990s, With the artificial neural network technology at that time, an academic trend has been set off. However, the artificial neural network at that time could not effectively deal with gradient vanishing gradient and exploding gradient while deepening the network level to improve the fitting accuracy (Bengio *et al.* 1994). Besides, along with the immaturity of the relevant algorithms and the low computing power of hardware, the artificial intelligence (AI) research based on data-driven concept is failed at that time, and then the research process in this field entered a slow development period of more than a decade.

The application of AI in the financial field is basically the application of machine learning (Bartram *et al.* 2020). According to the attributes of the modeling algorithm, machine learning can be divided into supervised learning, unsupervised learning and reinforcement learning. Supervised learning requires additional human effort to set the training labels or objective functions, and the intelligence of the model is directly related to the accuracy of the training labels or objective functions, thus supervised learning is not a fully intelligent machine learning model in a strict sense. Unsupervised learning mainly includes cluster analysis and self-supervised learning. At present, human intervention and good algorithms are still needed to improve the effectiveness of cluster analysis in the financial field. Self-supervised learning is an intelligent algorithmic model that has emerged in recent years, Heaton *et al.* (2017) proposed an asset allocation model with no short-selling mechanism, but self-supervised learning is not a mature algorithmic model. Therefore, unsupervised learning is not a fully intelligent machine learning model too. The modeling idea of reinforcement learning is much more different. It generates data through the interaction between agent and the environment, and the agent uses these generated data to learn the optimal strategy. Therefore, reinforcement learning is considered to be the closest type of machine learning algorithm to general AI. In the early days of asset management, applications of reinforcement

learning mainly used Q-learning and Policy Gradient (PG) algorithms. The Q-learning is a discrete data analysis algorithm, however, once the number of actions and environments increases to a certain level, the computing power requirements of this model increase dramatically, and even with the current computing power standards, the Q-learning can hardly satisfy the current investment market requirements. Moody & Saffell (1999) first applied the PG algorithm for individual asset management, but the drawback of PG algorithm is that the optimized strategy often falls into local optimality. Later, many scholars made improvements on the basis of John Moodys model , such as Dempster & Leemans (2006) and Deng *et al.* (2016), however, their model is still essentially based on the traditional reinforcement learning algorithms for individual asset management, which do not involve the scope of portfolio management.

Not until in 2015, A major technological breakthrough occurred, Mnih *et al.* (2015) applied deep reinforcement learning (DRL) to play computer games and prove that human can obtain AI by using DRL, which is a method that combines reinforcement learning with deep learning. Deep neural networks in deep learning can superimpose more network layers to improve fitting accuracy without vanishing gradient or exploding gradient, and their ability to fit functions far exceeds that of traditional artificial neural networks, which can easily crush humans in many intelligent tasks with improved reinforcement learning algorithms, such as Go, computer games StarCraft and League of Legends, etc. Therefore, this breakthrough inspires a lot of scholars to apply DRL in solving financial problem. Zhang *et al.* (2020) improved the model of John Moody *et al.* and used the improved reinforcement learning algorithm combined with LSTM networks to study equity index, commodities, fixed income and foreign exchange markets. However, Zhang *et al.*s model is not a deep reinforcement learning portfolio model in a strict sense for two reasons. First, LSTM was invented at the end of the AI research boom in the mid-1990s, and it has good performance in processing two-dimensional structured time series data with low computing power requirements, and LSTM can handle vanishing gradient or exploding gradient to some extent effectively, however, LSTM cannot improve the fitting accuracy of the network while superimposing more network layers, thus, if the reinforcement learning model only uses the LSTM as the network

structure, then the model cannot be considered as a true deep reinforcement learning model. Second, Zhang *et al.*s model is still a single asset management model in essence, which does not consider the portfolio as a whole and does not analyze how the weight of each asset changes during the trading process. In other words, the authors fix the investment decision of each asset in portfolio, and no longer let the agent to find the optimal investment decision, so the model is not a truly intelligent portfolio management model.

Jiang *et al.* (2017) pioneered a framework for applying deep reinforcement learning for portfolio management and applied it to the allocation of digital currency assets. In our opinion Jiang *et al.*s innovative framework has two salient features that make it more promising for application. First, unlike the commonly used two-dimensional data structure, the state space in this model has a three-dimensional structure, which is a high-dimensional data model, and such a data structure is very suitable for processing using deep neural networks. For example, image data is a three-dimensional data structure, and current deep neural network architectures with excellent performance are already comparable to humans in image recognition, which looks within reach of surpassing humans. But Jiang *et al.* did not take advantage of deep neural networks in their model, which making the reproduction of their model in other capital markets unsatisfactory (Liang *et al.* 2018). Second, Although the model only has a long mechanism, it regards the weight change of the portfolio as a continuous process of change, this way of setting up model inspired us to design shorting and arbitrage mechanism that fits the actual trading process.

The reason why the shorting mechanism in a continuous trading process is considered as an innovation is that in the literatures (Chincarini & Kim 2018; Guo *et al.* 2017), the sum weights of each asset in a portfolio is set to 1, i.e., $\sum_i^m \omega_i = 1$, regardless of whether it is long or short. Such a formula is fine if only long, but if shorting is considered, such a way of setting would exaggerate the portfolio's return and create the illusion of high returns. For example, considering two assets in a portfolio, if an investor sets the weight of the shorted asset to -1/3, the other assets weight will be 4/3. This setting is completely wrong and cannot be applied in a continuous trading process. In order to comply with the $\sum_i^m \omega_i = 1$ constraint and not to overstate the portfolio

returns, Jacobs *et al.* (2005); Kim *et al.* (2016) proposed models with a bound on the shorting asset weights, but such models are too cumbersome and make idealized assumptions on the asset weights, therefore, these models are rarely applied. Gu *et al.* (2020, 2021) applied the long-short decile method to set long-short asset weights, which avoids the problem of exaggerated return, but it is still essentially a fixed trading decision without analyzing the changes in asset weights, besides, the authors avoid the calculation of transaction costs. In addition, we found in many quantitative investment papers involving shorted assets that the authors did not explain how the weights of the shorted assets were set by investors (Song *et al.* 2017) or used the wrong asset weight settings (Jegadeesh & Titman 2001). Therefore, there is no theory to explain how shorted asset's weight changes in a continuous trading process that is consistent with the actual transaction process.

It has been well documented that no scholars have used DRL to make short in continuous action space for portfolio management with clearly explaining the shorting principle and the calculation of portfolio market value. Our innovations are mainly in the following four areas. First, we add a short mechanism to the original model. Second, we design an arbitrage mechanism according to Arbitrage Pricing Theory. Third, we redesign the activation function that be used to acquire action vector. Fourth, we reconstruct the neural networks in DRL based on deep learning network used in image processing, and apply DDPG algorithm to obtain more accurate results. According these four innovations, we can apply DRL in portfolio management to get AI in practice. In order to verify whether our model has a certain level of AI, we randomly select several stocks from the constituents of CSI300, and associate these stocks with CSI300 to form a portfolio. We repeat such a randomly selected portfolio for several sets, and verify whether these portfolios' rate of return (ROR) can outperform the market's average ROR which is denoted by CSI300's ROR. For the selection of indicators of strategy performance, we draw on the investment experience of Chan (2009) and use the Sharpe ratio as the main assessment indicator. If a randomly selected portfolio is able to outperform the CSI 300 in Sharpe ratio, the DRL portfolio model presented in this paper is considered to automatically optimize the asset weights in the portfolio, which allowing the investor to reap excess returns from the optimized decision, in other words, our model is able

to automatically optimize investment decisions. Besides, In the strategy performance comparison section, we also compared the stock picking strategy based on the factor model with the DRL strategy introduced in this paper under the short-selling mechanism.

## 2.    Theory and Methodology

### 2.1.    Reinforcement Learning

The core of current AI is reinforcement learning, which is a branch of machine learning. Reinforcement learning has four basic elements: agent, state or environment, action and reward. The agent takes actions to obtain rewards in the state and gets to the next state. The agent will repeat this process again and again until to the end of states. Such a process of iterative operation is called Markov Decision Process (MDP), which constitutes a complete trajectory from beginning to end. The trajectory could be defined as $\tau = (S_0, A_0, R_1, S_1, A_1, R_2, \cdots)$. The trading process of a portfolio can be regarded as an MDP, and a trader who owns a certain portfolio is the agent in reinforcement learning.

In a trajectory $\tau$, the average rewards obtained from an interaction between the agent and environment is $\bar{R}$. In our paper, the $\bar{R}$ is defined as

(1)
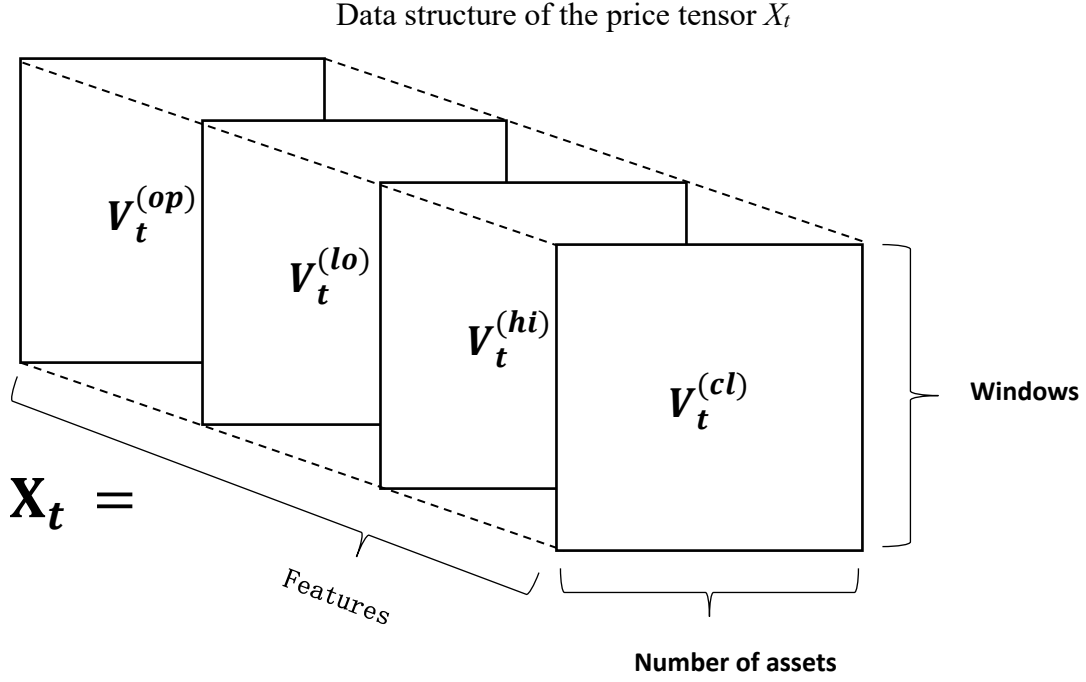$$\bar{R} = \frac{1}{t_f} \sum_{t=1}^{t_{f+1}} \gamma_t$$

In formula (1), $t_f$ denotes the $f^{\text{th}}$ trading period, $\gamma_t$ represents the payoff obtained during a single trading period at time $t$. It should be noted that this paper studies the Chinese stock market, which currently operates under the T+1 trading system, where a single stock asset can only be traded once a day in most cases. Therefore, we specify that each asset in a portfolio is only traded once a day, in addition, we also allow the position of the previous trading day to be opened once again after the position has been cleared for that day. Therefore, $\gamma_t$ denotes the daily logarithmic ROR of a portfolio in this paper.

### 2.2.    Definition of State Space and Price Tensor

We follow the definition of $S_t$ from Jiang et al.'s literature with some modifications to $X_t$

and $W_t$. Therefore, we set state $S_t = (X_t, W_t)$, where $X_t$ is a tensor consisting of price features, $W_t$ is an action vector consisting of the weight of each asset in a portfolio. We will discuss the structure of $S_t$ in the neural network design section. The structure of $X_t$ shows in figure 1.

FIGURE 1

Data structure of the price tensor $X_t$



Tensor $X_t$ is composed of four price features: close price ($V_t^{(cl)}$), high price ($V_t^{(hi)}$), low price ($V_t^{(lo)}$), open price ($V_t^{(op)}$). Since we set the length of observation windows equal to 50, namely n = 50, the data window is constructed at least from the $51^{st}$ trading period.

The formula for each price feature is defined as

$$V_t^{(cl)} = [v_{t-n+1} \oslash v_t | v_{t-n+2} \oslash v_t | \cdots | v_{t-1} \oslash v_t | 1]$$

$$V_t^{(hi)} = [v_{t-n+1}^{(hi)} \oslash v_t | v_{t-n+2}^{(hi)} \oslash v_t | \cdots | v_{t-1}^{(hi)} \oslash v_t | v_t^{(hi)} \oslash v_t]$$

(2)

$$V_t^{(lo)} = [v_{t-n+1}^{(lo)} \oslash v_t | v_{t-n+2}^{(lo)} \oslash v_t | \cdots | v_{t-1}^{(lo)} \oslash v_t | v_t^{(lo)} \oslash v_t]$$

$$V_t^{(op)} = [v_{t-n+1}^{(op)} \oslash v_t | v_{t-n+2}^{(op)} \oslash v_t | \cdots | v_{t-1}^{(op)} \oslash v_t | v_t^{(op)} \oslash v_t]$$

Among them, the lowercase letter $v_t$ denotes the close price vector of each asset in a

portfolio on trading day $t$, and the symbol $\oslash$ is element-wise division.

## 2.3. Definition of Action Space

We take the weights of each asset of portfolio as actions, so we set $Actions = Weights$. That is to say, action vector $A_t$ is weight vector $W_t$, where t represents the $t$[th] trading period. The $W_t$ can be written as

(3) $$W_t = \left(\omega_{0,t}, \omega_{1,t}, \omega_{2,t}, \cdots, \omega_{m,t}\right)$$

The first weight $\omega_{0,t}$ represents the proportion of cash value owned by investors to the total market value of a portfolio in the $t$[th] trading period, $\omega_{m,t}$ is the proportion of assets' market value invested by investors in CSI300 to total market value of the portfolio in the $t$[th] trading period.

Here is an explanation of the short mechanism. Taking the futures market as an example, the assets in the market can be shorted, and the transaction is margin trading. Either long position or short position, the money in the margin is occupied by the position, and this part of occupied money cannot be used for trading. That means, this occupied margin is actually the collateral asset, while the rest part of margin that is not occupied by the position can treated as a risk-free asset, which is cash in this article. Shorting in the stock market is similar, but instead of trading on margin, if you want to short stocks you have to take assets as collateral. Therefore, the main process of short selling in the stock market is to borrow the stock (since you do not own the stock, you have to borrow it, while the lender will withhold the collateral of yours), then sell the stock at the appropriate price. When the stock price fluctuates to a desired price, the same number of the stocks are purchased back, and all the stocks are finally returned to the debit side from which you borrow. Meanwhile you recover the collateral and getting the sell-buy spread among the stocks. Here we do not consider the dividends or bonuses obtained when holding the borrowed stocks, and we will discuss the holding cost of borrowed stocks in experiments section. In particular, we assume that the leverage ratio for short is 1:1, which means that the price of mortgage assets is the same as the real-time price of borrowed stocks. Negative weights mean short position and gain a positive return when the stock price drops, on the contrary, positive weights means long position and gain a

negative return when the stock price drops. Long position can only obtain positive returns when the stock price rises. Therefore, the value range of the action vector $W_t$ is [-1,1].

All the asset weights of a portfolio conform to the following relationship,
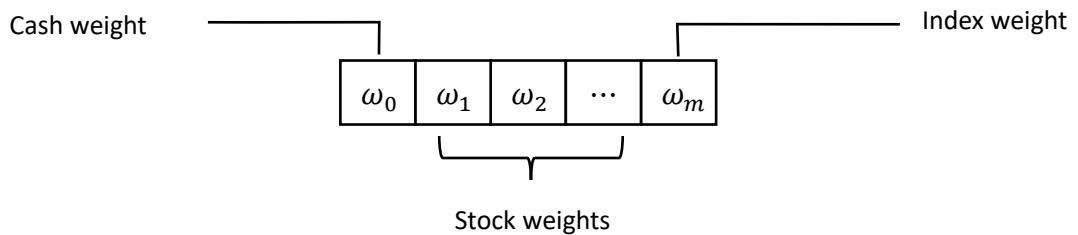
$$(4) \qquad \sum_{i=0}^{m} |\omega_i| = 1$$

That is, in an action vector, the sum of the respective absolute value of all action weights is 1, where $\omega_0$ represents cash's weight, since cash cannot be shorted, the value range of its weight is [0,1].

Here is an example of how to calculate the weights of shorted assets. Suppose that we borrowed some stocks with a market value of ¥300,000, namely we make a short position of these borrowed stocks, meanwhile, the market value of other assets in the portfolio is ¥700,000, at this moment, the weight of the borrowed stocks can be calculated as: $\omega_{borrowed} = \frac{-300,000}{700,000 + |-300,000|} = -0.3$.

We initialize the weight of a portfolio to $W_0 = (1,0,\cdots,0)^T$, which means that we only have cash asset in our portfolio and zero money invested in other assets at the beginning of a trade. We assume that the entire portfolio has 1 cash asset + m risk assets (risk assets refer to the constituents of CSI300 and the CSI300 stock index in this article). The structure of the weight vector is shown in figure 2.

FIGURE 2

The Structure of Weight Vector



We take the CSI300 financial futures into our portfolio. Since the constituents of CSI300 are high-quality stocks in Chinese market, the return of CSI300 is actually better than the

comprehensive performance of the market itself, but here we assume that the return of CSI300 can represent the overall return of the market. If you put all your money into CSI300, you are getting the average return of the market, and you cant get the excess return.

## 2.4. Arbitrage

Although we have pointed out some shortcomings of traditional econometric analysis in INTRODUCTION section, Economic theory can help us solve prediction policy problems (Mullainathan & Spiess 2017), namely, Although the ROR of each asset in the portfolio is not necessarily linear relationship with each other, Arbitrage Pricing Theory (APT) is still worthy to apply. We take CSI300 as the investment benchmark. Based on APT (Ross 2013), stocks' ROR has a linear relationship with each other because all stocks' ROR can be represented by some common factors. Make a long story short, if some stocks are overvalued and the price of the stocks are predicted to fall, the best way to get excess returns is to sell those overvalued stocks and have a long position of investment benchmark (CSI300) to reduce risk at the same time. Vice versa, if an asset is undervalued, go long this undervalued asset and short the CSI300 simultaneously. Therefore, here we introduce a very important mechanism in a portfolio, that is the weight of investment benchmark and the weight of other assets in a portfolio cannot be positive or negative at the same time. Therefore, the weight vector subject to the following restrictions,

(5) $$\sum_{i=1}^{m}|\omega_i| \neq |\sum_{i=1}^{m}\omega_i| \text{ and } \sum_{i=1}^{m}|\omega_i| - |\omega_m| \neq 0$$

As the formula (5) shows above, in a portfolio, the sum of the absolute values of the weight vector $\omega_i$ cannot be equal to the absolute value of its sum. We allow the agent to adopt a conservative strategy, that means all the funds will be allocated to CSI300 or both CSI300 and cash, and the weights of other risk assets are zero. We also allow the agent to adopt an aggressive strategy, all the funds will be invested in risk assets, or both risk assets and cash except investment benchmark, namely the weight of CSI300 is zero. If all the weights are positive or negative simultaneously while updating gradient, then we immediately reverse the weight of CSI300 $\omega_m$, i.e.,

(6)
$$\omega_m := -\omega_m$$

### 2.5. Reward Function

We take the daily logarithmic return of portfolio as the reward function. Defining the relative price vector $Y_t$ as

(7)
$$Y_t \triangleq P_t \oslash P_{t-1} = \left(1, p_{1,t} / p_{1,t-1}, \cdots, p_{i,t} / p_{i,t-1}\right)^T$$

The first element in $Y_t$ represents the relative price of cash, the last element represents the relative price of CSI300, and $P_t$ represents the closing price vector of each asset of a portfolio in trading period $t$. We assume that the growth rate of cash over the investment period is always zero (i.e., the value of a unit cash will neither increases nor depreciates over the investment period), therefore, the first element of the relative price vector $Y_t$ is always one.

If $\rho_t$ is used to represent the price of the portfolio in trading period $t$, ignoring transaction costs, the price of the portfolio can be expressed as

(8)
$$\rho_t = \rho_{t-1}\exp[(\ln Y_t) \cdot W_{t-1}]$$

Among formula (8), symbol $\cdot$ denotes the dot product of the vector, ln represents the natural logarithm calculator, and exp represents the natural exponential function.

The daily logarithmic ROR of a portfolio is

(9)
$$\gamma_t = \ln(\rho_t / \rho_{t-1})$$

In formula (1), we have the definition of $\bar{R}$, which is an average cumulation of $\gamma_t$. The purpose of training agent is to maximize $\bar{R}$ by using DRL.

### 2.6. Transaction cost

There are two kinds of transaction costs: fixed costs and variable costs. Variable costs include impact and opportunity costs which are difficult to quantify. So, we only consider the fixed costs in our paper. Inspired by Qian *et al.* (2007), here we analyze the process of calculating transaction costs. At the beginning of trading time $t$, the portfolio's weight vector is $W_{t-1}$. Due to price changes of assets in the market, the weight vector changes to $W_t'$ at the end of the trading time $t$. The formula for the weights' evolution is

$$(10) \qquad\qquad W_t' = (Y_t \odot W_{t-1}) \,/\, (Y_t \cdot |W_{t-1}|)$$

In formula (10), The symbol $\odot$ denotes Hadamard product, which means element-wise multiplication, and the symbol $|\ \ |$ represents the absolute value for each element in the vector. Since the objective function (reward function) will maximize the average returns of portfolio while updating gradient, the weight vector will evolve from $W_t'$ to $W_t$ before the start of trading time $t+1$, so that the whole rate of transaction cost of portfolio at this point is

$$(11) \qquad\qquad C_t = \mu_t \left( \sum_{i=1}^{m} \left| \omega_{i,t}' - \omega_{i,t} \right| \right)$$

$\mu_t$ is the single asset's rate of transaction cost in stock exchange, which include commission and other fixed costs in the transaction process. If the transaction cost is considered, the price of portfolio at this point is

$$(12) \qquad\qquad \rho_t = \rho_{t-1}(1 - C_t)\exp[(\ln Y_t) \cdot W_{t-1}]$$

Then we analyze the specific process of weight evolution. To make it easier to understand, let $portfolio\_\gamma_t = \gamma_t$. Assuming that at time $t-1$, the trader has allocated the weight of each asset in the portfolio, which is $W_{t-1}$, when coming to time t, the price of each asset will change accordingly (even if the price remains the same, we also believed that the asset price has changed, but the degree of change is zero), at this time the portfolio's ROR is,

$$(13) \qquad portfolio\_\gamma_t = \gamma_{0,t}\omega_{0,t-1} + \gamma_{1,t}\omega_{1,t-1} + \cdots + \gamma_{m,t}\omega_{m,t-1} = \Gamma_t \cdot W_{t-1}$$

Among formula (13), $\Gamma_t$ is the vector which represent each asset's ROR in portfolio at time $t$, it can be calculated through $\ln Y_t$ (note that at this timepoint, when $\Gamma_t$ has been calculated, agent have not reassigned weights to the assets), due to price changes, the weight of each asset in the portfolio will automatically change to $W_t'$ accordingly. And the calculation formula for the $W_t'$ is shown in formula (10). After that, the agent will predict the price movement of assets at time $t+1$ and reassign the weight $W_t$ to the assets before the end of time $t$. We can calculate the transaction cost through the weights' change from $W_t'$ to $W_t$ based on formula (11), and the agent will repeat this transaction process until the trading is closed. Actually, the $W_t$ is obtained by the neural network in DRL, which we will discuss later. Please note that all the weights must subject

to equation （4）.

### 2.7. Portfolio return after leveraging assets

If leveraging assets in portfolio, the formula for calculating the portfolio's ROR is

$$portfolio\_r_t = \lambda_{0,t-1}\gamma_{0,t}\omega_{0,t-1} + \lambda_{1,t-1}\gamma_{1,t}\omega_{1,t-1} + \cdots + \lambda_{m,t-1}\gamma_{m,t}\omega_{m,t-1}$$

(14)

$$\text{i.e., } portfolio\_r_t = \Lambda_{t-1} \cdot \Gamma_t \cdot W_{t-1} \text{ and } \sum_{i=0}^{m}|\omega_{i,t}| = 1$$

where $\lambda_{0,t-1}, \lambda_{1,t-1} \cdots \lambda_{m,t-1}$ are the leverage ratios of each asset in portfolio at trading period *t-1*, which are treated as the leverage vector $\Lambda_{t-1}$. To simplify the analysis process and for fair comparison in the back-test session, all assets in this paper are unleveraged and the leverage ratio of each asset is 1:1, i.e., the leverage vector $\Lambda_{t-1}$ is constant 1.

## 3. Deep Reinforcement Learning Algorithm

### 3.1. DDPG algorithm

Once we have set up the environment, the reward function and the action vector, the next step is to choose a strategy for the agent to explore the environment and obtain the rewards. This strategy for exploring the environment and gaining rewards is a reinforcement learning algorithm. Deep reinforcement learning algorithm is a reinforcement learning algorithm that incorporates deep neural networks.

In our paper, we take Deep Deterministic Policy Gradient (DDPG) as the algorithm of reinforcement learning. DDPG has relatively good properties in various algorithms for reinforcement learning. There are two points to note when applying the DDPG algorithm. First, all reinforcement learning algorithms have the two characteristics of Explore and Exploit, that means, reinforcement learning algorithms are about how to make the agent explore more rewards while exploiting the environment, and the DDPG algorithm is no exception. Since the policy function in DDPG algorithm is deterministic, it does not explore the environment through probabilistic algorithms like PG algorithm, so in order to balance exploration and exploitation, the policy function in the DDPG algorithm will explore data with a certain probability. Therefore, the

policy function of DDPG is added noise from a normal distribution *N(0.05, 0.25)* in this paper. Second, DDPG is an off-policy algorithm that learns the Q value function and the policy $\pi$ in a continuous action space and observation space. In DDPG algorithm, the Q value function is obtained from the Critic network and the policy $\pi$ is obtained from the Actor network, if we require better fitting accuracy for Q values and policy $\pi$, we need to change the structure of Critic and Actor networks. These two points are what we need to pay attention when applying DDPG algorithm. We will cover the design of Critic and Actor networks in detail in the next section.

More details about DDPG algorithm can be found in the related literatures (Silver *et al.* 2014; Lillicrap *et al.* 2015). We recommend reading these two papers after understanding the policy iteration method, PG algorithm, Temporal-Difference algorithm and Deep Q-Networks algorithm, or directly reading the code of DDPG algorithm, which is relatively easier to understand.

3.2.    Deep Neural Network Design

Deep Neural Network solves the previous problem of exploding gradient and vanishing gradient when updating gradients in Artificial Neural Network, it can overlay more network layers and thus fit the objective function more accurately. Let's suppose this situation, if a DRL algorithm is pretty good but the performance of the neural network adopted by the DRL algorithm is mediocre, it is unquestionable that the agent cannot obtain the optimal strategy. We describe this situation as "the neural network is lack of expression."
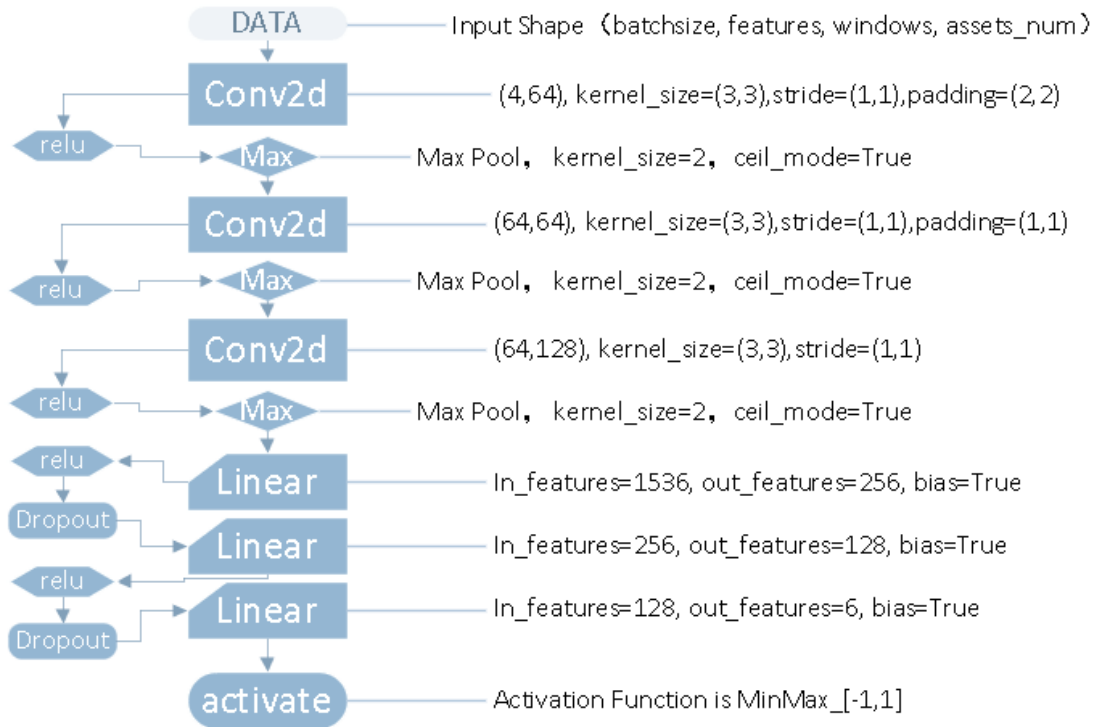
Since the tensor $X_t$ in the state $S_t$ has the same structure as the image data, it is natural that we thought of designing our deep neural network by using the network structure for processing image data. It is meaningful by applying image processing network in financial data. For instance, a 2-dimensional time series data can be treated as an image data, in which the information of times series data looks like pixel images (Sezer & Ozbayoglu 2018), in other words, this image data is the data which include a piece of financial information. In our experiments, the price tensor $X_t$ is a 3-dimensional data, which can be viewed as a group of 2-dimensional data with a lot of financial information. We chose VGG net (Simonyan & Zisserman 2014) as a reference for network design.

The DDPG algorithm has two neural network structures which are Actor and Critic. Actor is used to explore the policy function, and Critic is used to represent the value function, we will introduce the network design of these two structures in the following.

In our paper, Actor is used to explore the optimal action vector, i.e., the optimal weight vector. Since the purpose of training the neural network is to maximize the $\bar{R}$, and $\bar{R}$ is calculated from the closing price of the assets which is also an element of the tensor $X_t$, and the weight vector $W_t$ is also used to calculate $\bar{R}$, therefore, maximizing the value of $\bar{R}$ will definitely change the value of $W_t$ when updating gradients. Based on the analysis above, the structure of Actor is,

FIGURE 3

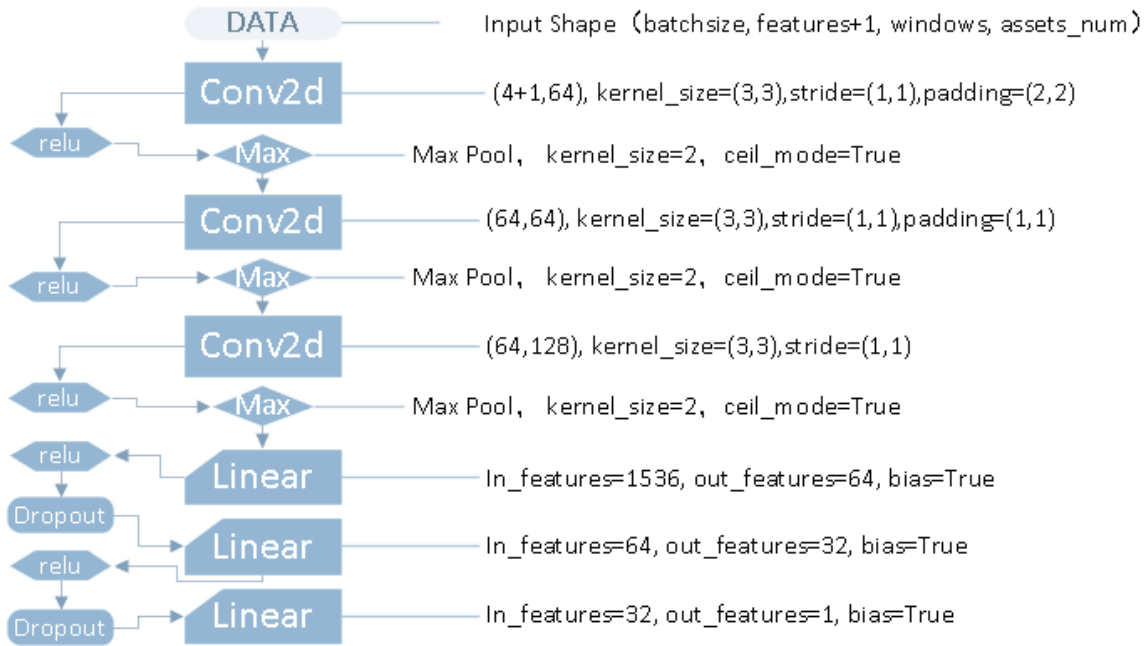The Neural Network Structure of Actor



We do not take into account cash in the input data because the Actor neural network will directly output the optimized weight of cash. In particular, here we explain how to design the activation function used for outputting weight vector. We have tried several activation functions, and finally we decided to use the Min-Max normalized activation function to obtain the weight vector. First, it is necessary to Min-Max normalize the output of the Actor network to limit the

value range of each element of action vector to [0,1], this normalized action vector can be set to Action$_{min\_max[0,1]}$, then we use linear transformation to change its value range to [-1,1], that is, Action$_{min\_max[-1,1]}$ = 2(Action$_{min\_max[0,1]}$ − 0.5). Finally, Action$_{min\_max[-1,1]}$ is scaled (divided by the sum of the absolute value of each asset's weight) so that the weights of a portfolio subject to Equation（4）.

Critic represents the action-value function, which is used to evaluate the quality of the action $A_t$ in the state $S_t$. Therefore, both tensor $X_t$ which represents state and tensor $W_t$ which represents action, are the elements of Critic's input. The structure of Critic is,
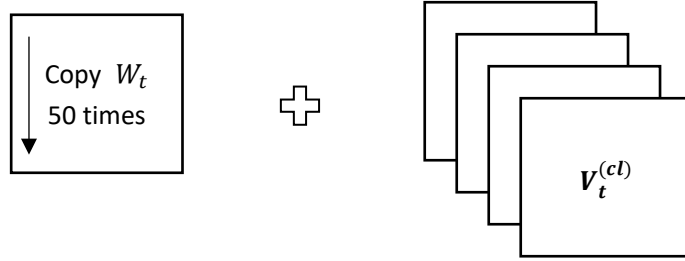
FIGURE 4

The Neural Network Structure of Critic



The network structure of Critic is similar to the structure of Actor, with some simplification of the parameters, and finally the action-value is obtained directly through the linear layer not the activation function, this is different from Actor. It is important to note that the way to combine $W_t$ and $X_t$ as input data of Critic, here we directly replicate $W_t$ 50 times in the longitudinal direction and insert it into the tensor $X_t$ as a feature layer.

FIGURE 5

Critic's Input Data Structure



Look at Figure 5, it is equivalent to add a new channel to the tensor $X_t$, i.e., add a new feature layer to $X_t$, and this is the structure of $S_t$ that we mentioned in the Definition of State Space. In this paper, a portfolio includes 5 risky assets (4 stocks + 1 stock index), if the number of assets is changed, the network parameters of Actor and Critic need to be modified accordingly.

## 4. Experiments

### 4.1. Data selection, pre-processing, and hypotheses

In order to randomly select a stock, we first find the CSI300 constituent section in the brokerage client, in which each stock is numbered from 1 to 300, and each number uniquely corresponds to a stock. Then we use the python code: "numpy.random.randint(300,size=4)" to select four random numbers. This numbers have a one-to-one correspondence with the stock's section number in CSI300's constituent section. To ensure that the number of training sample is as large as possible, we only select stocks of which the listing date is before December 31, 2010. If a randomly selected stock does not meet this condition, we will continue to search for stocks by adding +1 to the stocks sec tion number in the client until we find a stock that meets the listing date requirement. Repeat this selection process four times, and each selection constitutes a randomly selected portfolio. Each portfolio includes the CSI300 and 4 randomly selected constituents of CSI300, and the stocks in each portfolio are all not repeating.

Since the stocks are the constituents of CSI300, the stocks can be shorted. Our data comes from the Wind database, and all missing values of asset are filled with zero. We also assumed that all risk assets are liquid enough, each trading is executed immediately, and the trade does not cause

a market impact.

4.2.    Hyperparameter setting and optimization ideas for DRL model

The main hyperparameter settings in this paper are: replay buffer = 600, batch size = 64, optimizer for both Critic and Actor networks is Adam, learning rate is $5 \times 10\text{-}4$ and $4 \times 10\text{-}5$ respectively. In addition, we refer to wassname (2019)'s open source code idea on github to make agent randomly sample 252 consecutive trading days from the entire training set for each episode, i.e., steps = 252 for each training episode. Please note that this is different from the commonly used reinforcement learning ideas. Normally, the entire training set is used as a training trajectory. However, due to the vagaries of the capital market and the overpowering ability of DRL to fit a certain nonlinear function when exploring environment, it is easy to overfit if we train the entire training set as one trajectory. Therefore, in order to avoid overfitting, agent randomly sample 252 consecutive trading days' data from the entire training set for training at each episode. The total training steps is 300000, that means, after 1191 training episodes, the entire training process will end($300000 \div 252 \approx 1191$).

In our paper, we did not treat DRL as a black box, here are two reasons to support this opinion. First, the principle of neural network is to imitate the neurons of human brain, that means deep neural networks can be recognized as a simulation of human brain. Second, reinforcement learning is a well-established theoretical system for optimal control problem (Sutton & Barto 2018), in which trial-and-error search and delayed reward are the most important characters, these features are the same as human's, which are to imitate the human brains reaction to things. From these two reasons, it can be seen that DRL is to imitate the way of thinking of the human brain towards things. Lets imagine a situation where if a gamer beats his rivals in Starcraft computer game, can we assume that this gamers brain is a black box? Obviously not. Similarly, if a human gamer is defeated in StarCraft by AI using DRL, the AI cannot be considered as a black box too. That's because DRL is a simulation of the human brain, which imitating the way that human process things, and there is a complete theoretical system to support DRL. Therefore, if the designed DRL model lacks explanatory power or has deficient performance, we should look for

improvements from the theoretical system of DRL, rather than treating it as a black box.

4.3. Performance Metrics

Our purpose is to verify whether a randomly selected portfolio's ROR can outperform the market's average ROR by using DRL, therefore, we chose the following metrics,

1) Simple-daily-return: Average daily simple return,

2) Log-daily-return: Average daily logarithmic return,

3) Simple-annual-sharpe: Annualized Sharpe ratio using simple return,

4) Log-annual-sharpe: Annualized Sharpe ratio using logarithmic return,

5) Simple-annual-sortino: Annualized Sortino ratio using simple return,

6) Log-annual-sortino: Annualized Sortino ratio using logarithmic return,

7) MDD: Maximum drawdown is the maximum observed loss from a peak to a trough of an asset being traded.

4.4. AI Verification

4.4.1. Transaction cost setting

With respect to transaction cost, we only consider the stamp duty, commission fees and cost of borrowed stocks. We set the rate of stamp tax at 2‰ when selling or buying, and set the commission rate at 2‰ when trading. Besides, the annual transaction cost of holding borrowed stocks is 10.6% in Chinese stock market, $10.6\% \div 365 = 0.29‰$, namely the daily transaction cost of shorting an asset is 0.29‰. which is approximately equal to 3‰. Therefore, the whole transaction cost is 2‰+2‰+3‰=2.5‰ for every single asset being traded, i.e., $\mu_t = 0.0025$. Actually, not all assets of a portfolio can be shorted simultaneously because of arbitrage mechanism. And in reality the rate of stamp tax is 1‰ from the year of 2008, which only can be charged on the sale of assets. Consequently, such transaction cost of $\mu_t = 0.0025$ is actually at a very high level in Chinese stock market.

4.4.2. Training and Back-testing Periods

In our experiments, we randomly select 4 stocks which are constituents of CSI300 to construct a portfolio with CSI300, repeat the randomly selected process 4 times, and then observe the results.

Table 1. Training and Back-testing Periods

| Group | Training periods | Testing periods |
| --- | --- | --- |
| Stochastic Portfolio_1 | 2007/11/06 to 2020/2/05 | 2020/4/13 to 2021/4/26 |
| Stochastic Portfolio_2 | 2007/11/06 to 2020/2/05 | 2020/5/20 to 2021/6/02 |
| Stochastic Portfolio_3 | 2010/9/29 to 2020/6/09 | 2020/7/2 to 2021/7/14 |
| Stochastic Portfolio_4 | 2010/9/29 to 2020/6/09 | 2020/7/2 to 2021/7/14 |

The testing data is all out-of-sample, none of which is included in the training data. That means, the agent only knows the testing data when back-testing, and until then the agent is completely unaware of what future price movements will look like.
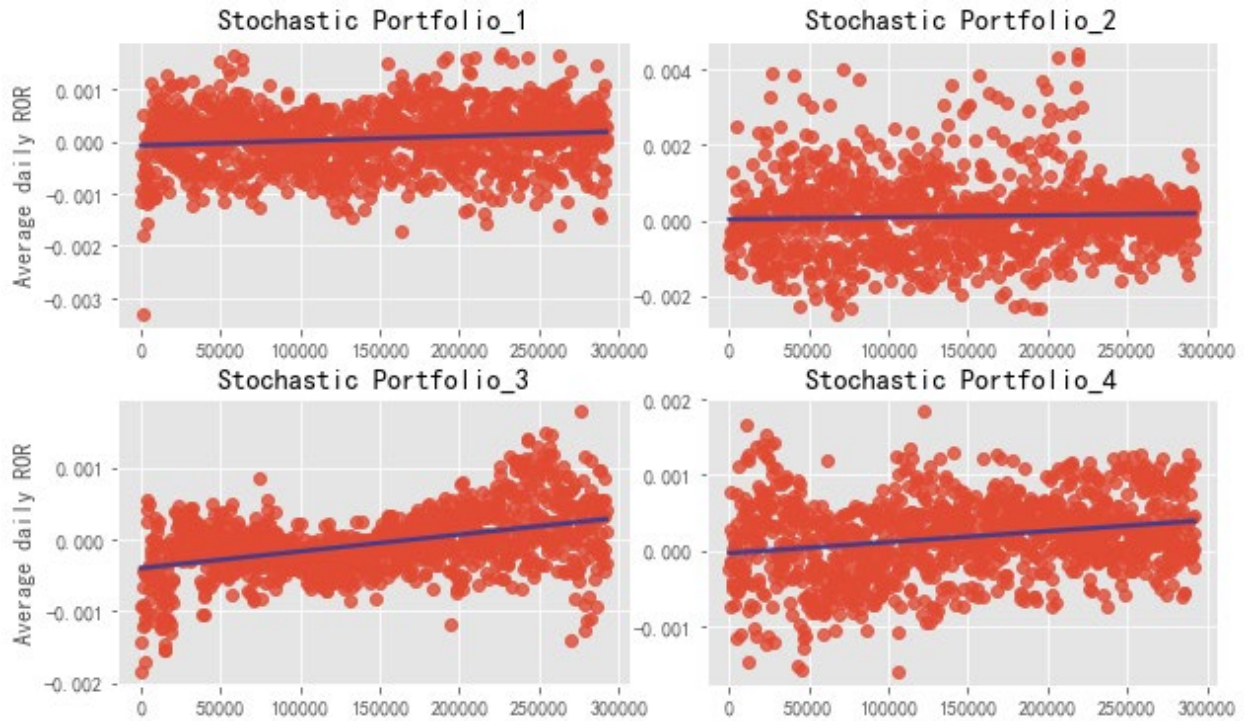
As can be seen in Table 3-1, the testing periods is about 1 years, which was chosen in this way for two reasons: First, unlike long-term trading strategies that investors trading assets every two weeks or one month and choose 3 to 5 years for back-testing, our trading frequency belongs to medium frequency, there are 252 trading decisions need to be made during 1 year back-testing period. So, such a back-testing period is long enough to meet the demand of investor's observation. Second, in the hyperparameter setting section we mentioned that at each training episode agent will randomly sample 252 consecutive trading days data from the training set, so choosing 1-year back-testing period also meets our sampling rule. Therefore, it is very reasonable for us to set a back-testing period of 1 year.

4.4.3. Training performance indicator

As we mentioned in the introduction section, the data-driven approach is different from the traditional econometric approach, it does not set an equation and estimate the parameters of the equation from the overall sample, so the traditional econometric significance testing is not applicable to DRL at all. Besides, some performance metrics in Tensorboard are commonly used in deep learning to discriminate the training effect of the model, but deep learning belongs to

supervised learning, which is completely different from reinforcement learning, so it is not suitable to discriminate the training effect of DRL by the criteria of deep learning either. At present, there are no universally accepted metrics to discriminate the training performance of DRL in finance. In order to get appropriate training performance metrics for DRL, we draw on the idea of open source code on github by wassname (2019), if the training performance gets better in general over time, the agent is considered to be getting smarter as it learns, and the training performance meets the requirements.

FIGURE 6
Training results regression line



In Figure 6, the horizontal coordinate is the number of training steps and the vertical coordinate is the average daily ROR. If the slope of the regression line is positive, the agent is considered to be smarter as it learns. It is intuitive from figure 6 that the slope of each regression line is positive, so the training performance meets the requirements.

### 4.4.4. Back-testing results
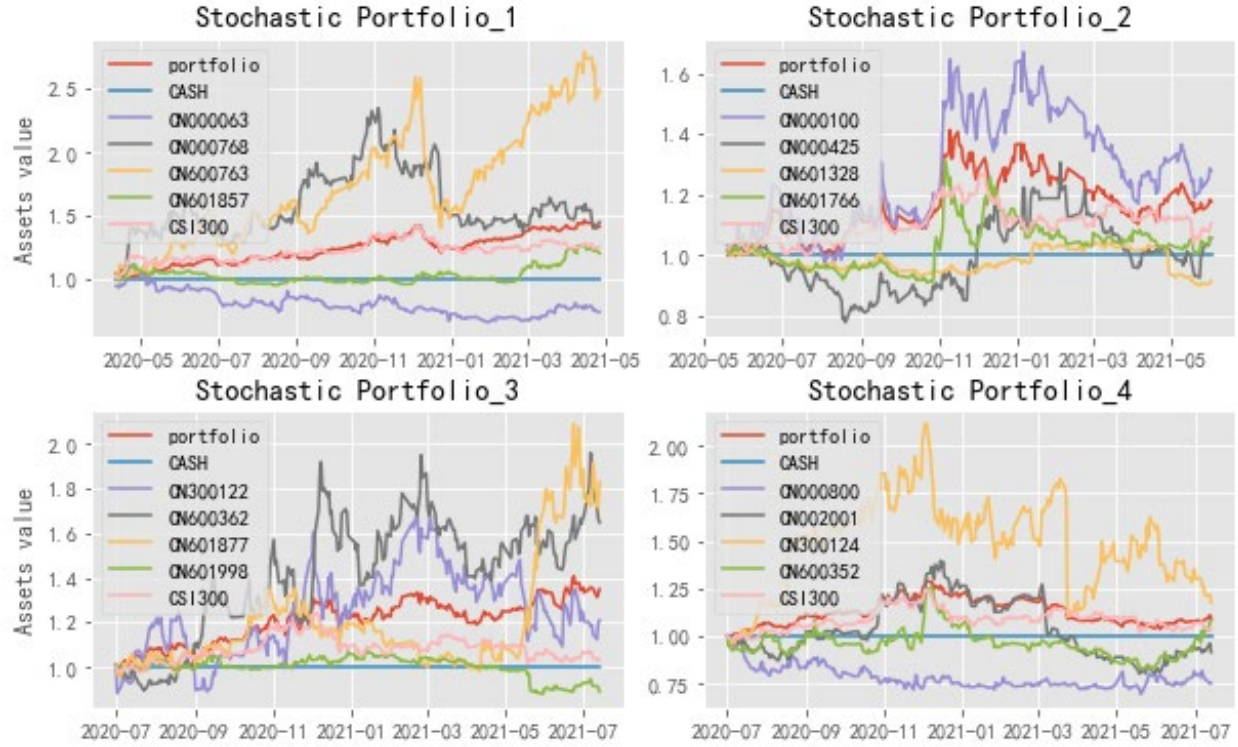
FIGURE 7
Assets value in Back-testing



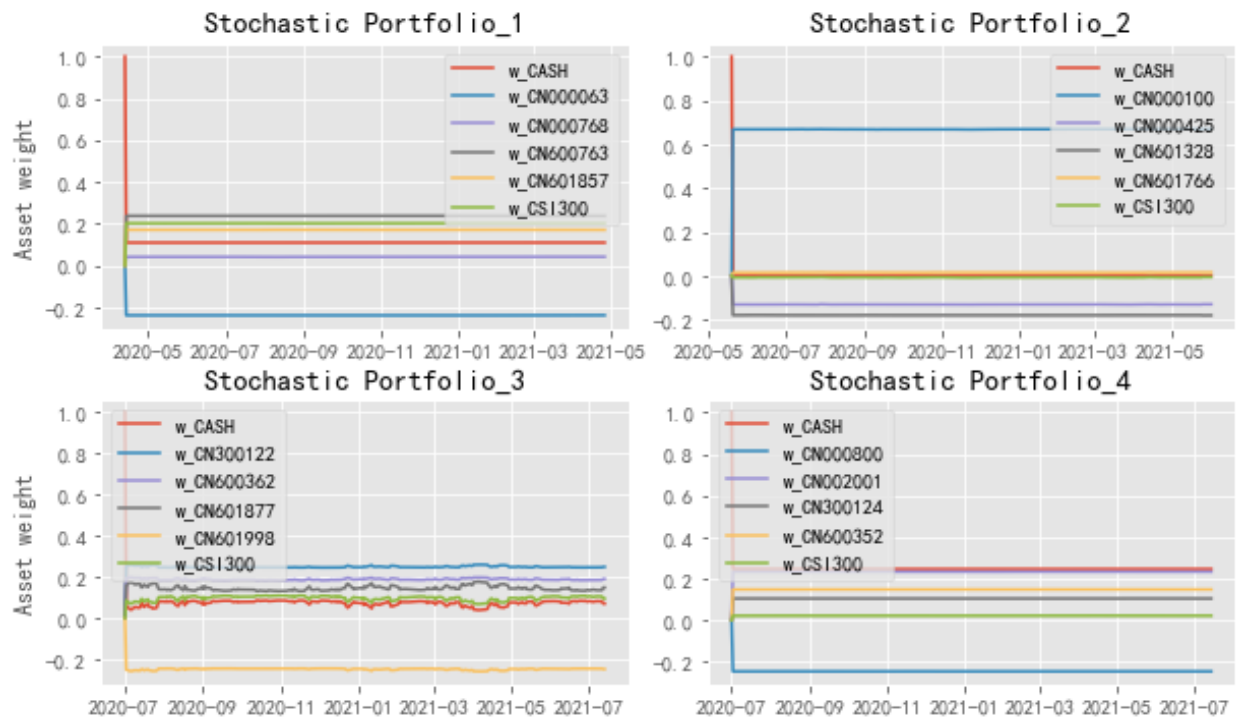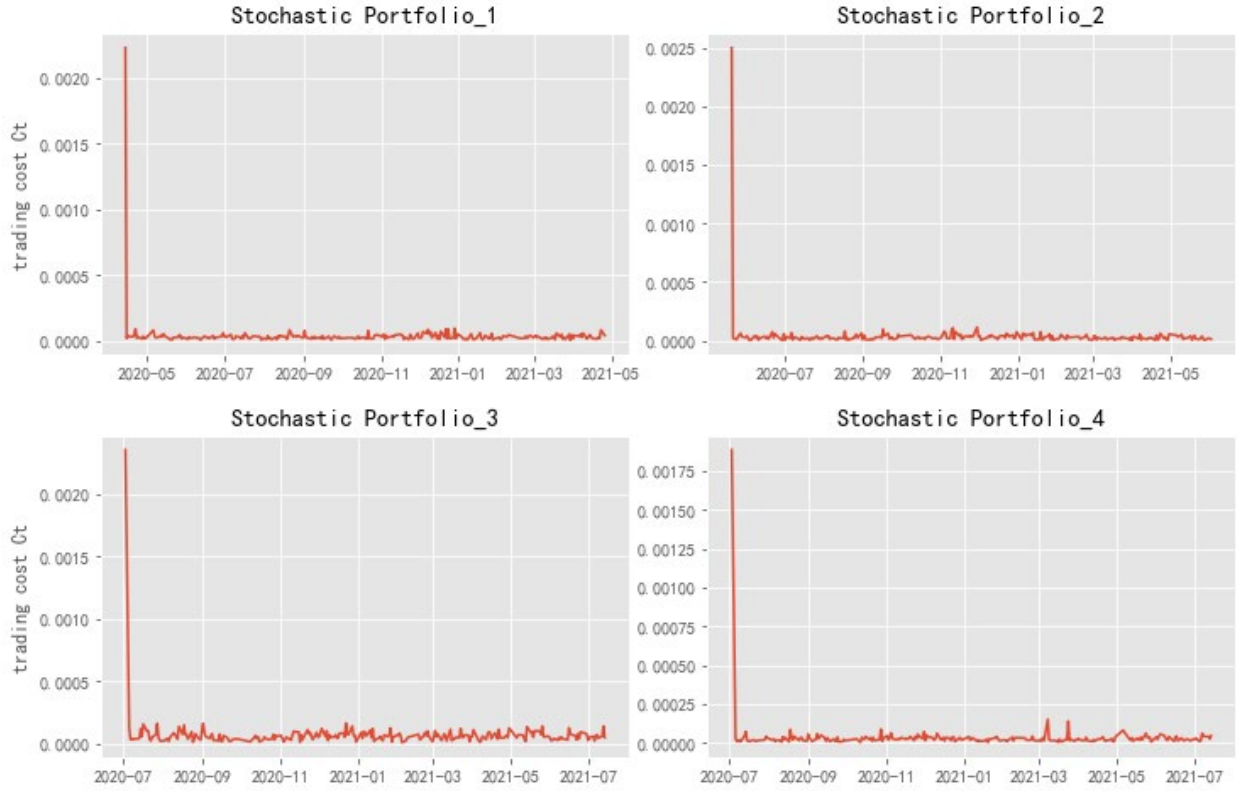FIGURE 8
Assets' weight in Back-testing

FIGURE 9
Transaction cost in Back-testing



In Figure 7, all portfolios outperformed the CSI 300 in value at the end of the back-testing period. Actually, each portfolio outperformed the CSI 300 on all metrics.

In Figure 8, almost all asset weights of all portfolios remain at a fixed value throughout the back-testing period, with some assets fluctuating slightly, such as the 600362 stock in stochastic portfolio 3, whose weight fluctuates slightly around the value of 0.185, but overall remains at a fixed value. This is not deliberate, but a result of our models decis ion making.

In Figure 9, the transaction cost rate $C_t$ for each portfolio is 0.0025 on the first day of trading, after which the agent maintains the transaction cost at a very low level. Combined with Figure 8 we can see that the agent is trading every day, and almost maintaining the assets weights at a fixed value. The transaction cost rate after opening a position is actually not high.

Table 2. Back-testing results

| | Stochastic Portfolio_1 | CSI300 |
|---|---|---|
| **Simple-daily-return** | 0.001409 | 0.000189 |
| **Log-daily-return** | 0.001358 | 0.000025 |
| **Simple-annual-sharpe** | 2.253346 | 0.165819 |
| **Log-annual-sharpe** | 2.169531 | 0.022390 |
| **Simple-annual-sortino** | 2.170365 | 0.118142 |
| **Log-annual-sortino** | 3.083598 | 0.268245 |
| **MDD** | 15.06% | 9.92% |
| | **Stochastic Portfolio_2** | **CSI300** |
| **Simple-daily-return** | 0.000837 | 0.000093 |
| **Log-daily-return** | 0.000643 | -0.000049 |
| **Simple-annual-sharpe** | 0.674304 | 0.087690 |
| **Log-annual-sharpe** | 0.518942 | -0.046470 |
| **Simple-annual-sortino** | 1.146697 | 0.147968 |
| **Log-annual-sortino** | 0.859027 | -0.077108 |
| **MDD** | 21.95% | 6.49% |
| | **Stochastic Portfolio_3** | **CSI300** |
| **Simple-daily-return** | 0.001295 | 0.000096 |
| **Log-daily-return** | 0.001179 | -0.000049 |
| **Simple-annual-sharpe** | 1.353393 | 0.089181 |
| **Log-annual-sharpe** | 1.233873 | -0.045915 |
| **Simple-annual-sortino** | 2.261293 | 0.155939 |
| **Log-annual-sortino** | 2.026551 | -0.078898 |
| **MDD** | 10.70% | 6.49% |
| | **Stochastic Portfolio_4** | **CSI300** |
| **Simple-daily-return** | 0.000402 | 0.000096 |
| **Log-daily-return** | 0.000345 | -0.000049 |
| **Simple-annual-sharpe** | 0.599214 | 0.089181 |
| **Log-annual-sharpe** | 0.513162 | -0.045915 |
| **Simple-annual-sortino** | 0.776685 | 0.155939 |
| **Log-annual-sortino** | 0.654805 | -0.078898 |
| **MDD** | 19.89% | 6.49% |

In Table 2, we can see that all portfolios outperform the CSI 300on all the metrics except MDD. Although all the MDD are relatively high, the Sharpe-ratio and Sortino-ratio of portfolios are better than CSI300 throughout the back-testing periods, so the anti-risk ability of the optimized portfolio by DRL are significantly better thanCSI300. Therefore, our model has the ability to optimize investment decisions, and the optimized agent can obtain excess returns in the stock

market.

4.5.    Strategies Performance Comparison

In this section we compare the strategy generated by our DRL model in this paper with the strategy based on multi-factor model. According to Liu *et al.* (2019), the value factor can explain the anomalies of Chinese stock market and is one of the important factors, and the Earnings-to-Price ratio factor which belongs to value factor can be simply replaced by the turnover factor, where investors will long stocks with low turnover and short stocks with high turnover.

We select stocks based on Liu *et al.*'s study to form portfolios, using the two factors of Earnings-to-Price ratio (which is the inverse of the P/E ratio) and turnover as the basis for stock selection, and the selection pool is the constituents of CSI300. We first multiply the turnover factor by -1, then give the same weight to both factors for summing, that means (-1) × Earnings-to-Price and turnover will be multiplied with 0.5 respectively and summed together. The summed results are then sorted from largest to smallest, in which the larger values indicating that related stocks are more valuable for long positions and vice versa for short positions. Unlike the convolution operations in deep neural networks, which require strong computing power when training, the multi-factor model consumes almost no computing power, so a large number of stocks can be selected to form a portfolio. Therefore, every trading day we select the top 20 stocks from sorted pool to go long and the last 20 stocks to go short, for this total of 40 stocks to form a portfolio each day with equally absolute values of weights for each stock.

Specifically, the trading strategy of multi-factor model is to first select 40 stocks based on the previous trading days closing data, in which the stocks selection scheme is described above, and the weights of each stock are assigned by trader at the opening of the trading day. Take a note, the asset weights in the portfolio must comply to the requirements of Equation （4）, so that the weight of the asset being long is 1/40 and the weight of the asset being short is -1/40. We do not consider transaction cost for multi-factor model in the whole trading process, and there is no leverage for assets. And the portfolio's ROR is calculated using Equation （13）, where the daily ROR of each asset in portfolio is calculated using last trading day's close prices.
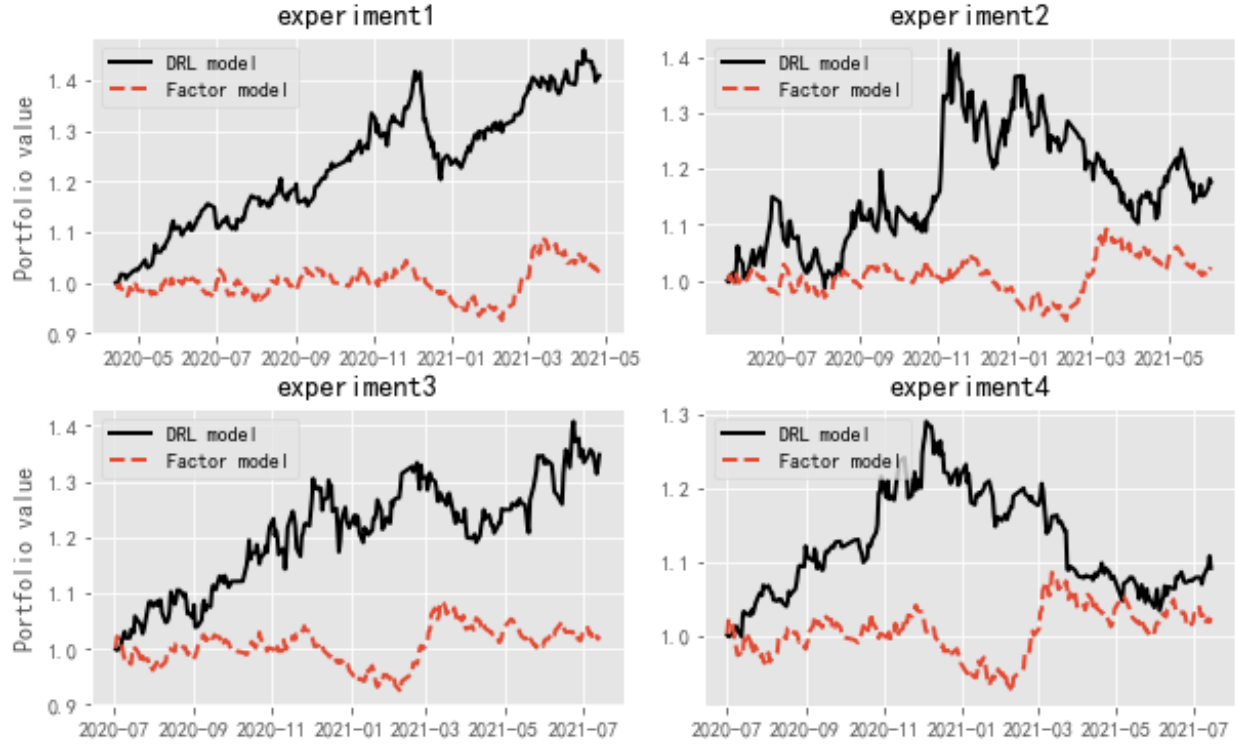
FIGURE 10
Strategies Performance Comparison



Table 3. Performance comparison in detail

| Group | Strategy | Log-daily-return | Log-annual-sharpe | Log-annual-sortino | MDD |
|---|---|---|---|---|---|
| Experiment 1 | DRL | 0.001358 | 2.169531 | 2.935230 | 15.06% |
| | Multi-factor | 0.000080 | 0.152778 | 0.356408 | 11.31% |
| Experiment 2 | DRL | 0.000643 | 0.518942 | 0.859027 | 21.95% |
| | Multi-factor | 0.000092 | 0.178042 | 0.414767 | 11.31% |
| Experiment 3 | DRL | 0.001179 | 1.233873 | 2.026551 | 10.70% |
| | Multi-factor | 0.000089 | 0.174027 | 0.394749 | 11.31% |
| Experiment 4 | DRL | 0.000345 | 0.513162 | 0.654805 | 19.89% |
| | Multi-factor | 0.000089 | 0.174027 | 0.394749 | 11.31% |

From the table 3, we can see that the strategy of DRL outperforms the strategy of multi-factor model on all the metrics except MDD. Although three portfolios' MDD are higher than the multi-factor model, the figure 10 shows the DRL model leads across the board in the assets' value

comparison for the entire back-testing periods. Besides, we do not consider the transaction cost when applying multi-factor model. Therefore, we can fully believe that the DRL model has application value.

## 5.  Conclusion

The allocation of asset weights in portfolio management is a direct reflection of investment decisions, and this paper introduces a data-driven approach to optimize investment decisions using AI techniques. We redesigned the formula for calculating the portfolio asset weights that matches the actual trading process, in which the assets can be shorted and leveraged. And we make several innovations based on previous studies. Experimental results show that the data-driven AI model can optimize investment decisions and beat the stock market. Even with high transaction cost and randomly selected stocks, the stochastic portfolios can still get excess returns.

We found that the optimized agent will almost maintain the asset weights at certain fixed values throughout the trading periods once the agent starts trading. In our previous experiments, we also found that in a few cases, the optimized agent will only change some asset weights once or twice during the whole trading periods, and then maintain the changed asset weights until the end of the back-testing, indicating that the optimized agent will not change the asset weights frequently. From the experiments we can see, the optimized agent trading assets throughout the transaction periods. The transaction cost is maximum only on the first day of opening a position, after which the optimized agent maintains the asset weights at a certain constant value and make the transaction cost at a very low level. In back-testing periods, the optimized stochastic portfolios have been able to outperform the CSI300 on most metrics, especially the Sharpe ratio. In previous experiments, we have even seen a Sharpe ratio of around 2.8, and the optimized stochastic portfolios have significantly outperformed the CSI300 in terms of asset value at the end of back-testing, but it is a relatively rare occurrence.

We also explain that DRL cannot be considered as a black box, because there is a complete theory to support it, and if the DRL model has deficiencies or the experimental results are not

satisfactory, we have to look for improvements from the theory of DRL. In our model we add an arbitrage mechanism, which is actually a prior rule that instructs agent to follow human experience to find the optimal investment decision. To get better results in back-testing, more prior rules can be added in the DRL modeling. Although we utilize the network structure for processing images in our network design, financial data is different from image data, there have not been a specialized deep learning network for processing 3-dimension financial data, so our model still has much room for improvement.

In our opinion, most algorithmic trading models can be implemented with DRL, only unlike that setting up relational equations for each variable used in traditional algorithmic trading models, the DRL model is set up with reward function, action vector and the environment in which the agent interacts, and then finally the data is used with the DRL model to tell us the optimal path to obtain the return.

## REFERENCES

Bartram, S.M., Branke, J., Motahari, M., 2020. Artificial Intelligence in Asset Management. CEPR Discussion Papers

Bengio, Y., Simard, P., Frasconi, P., 1994. Learning long-term dependencies with gradient descent is difficult. IEEE transactions on neural networks 5, 157-166

Chan, ErnestP, 2009. Quantitative trading : how to build your own algorithmic trading business. Quantitative trading : how to build your own algorithmic trading business.

Chen, B., Pearl, J., 2013. Regression and causation: a critical examination of six econometrics textbooks. Real-World Economics Review, Issue, 2-20

Chincarini, L. B., and D. Kim. *量化股票组合管理：积极型投资组合构建和管理的方法*: 机械工业出版社 (2018).

Colquhoun, D., 2014. An investigation of the false discovery rate and the misinterpretation of p-values. Royal Society open science 1, 140216

Dempster, M.A., Leemans, V., 2006. An automated FX trading system using adaptive reinforcement learning. Expert Systems with Applications 30, 543-552

Deng, Y., Bao, F., Kong, Y., Ren, Z., Dai, Q., 2016. Deep direct reinforcement learning for financial signal representation and trading. IEEE transactions on neural networks and learning systems 28, 653-664

Gu, S., Kelly, B., Xiu, D., 2020. Empirical asset pricing via machine learning. The Review of Financial Studies 33, 2223-2273

Gu, S., Kelly, B., Xiu, D., 2021. Autoencoder asset pricing models. Journal of Econometrics 222, 429-450

Guo, X., Lai, T.L., Shek, H., Wong, S.P.-S., 2017. Quantitative trading: algorithms, analytics, data, models, optimization. CRC Press.

Harvey, C.R., Liu, Y., Zhu, H., 2016. … and the cross-section of expected returns. The Review of Financial Studies 29,

5-68

Heaton, J.B., Polson, N.G., Witte, J.H., 2017. Deep learning for finance: deep portfolios. Applied Stochastic Models in Business and Industry 33, 3-12

Hou, K., Xue, C., Zhang, L., 2020. Replicating anomalies. The Review of Financial Studies 33, 2019-2133

Jacobs, B.I., Levy, K.N., Markowitz, H.M., 2005. Portfolio optimization with factors, scenarios, and realistic short positions. Operations Research 53, 586-599

Jegadeesh, N., Titman, S., 2001. Profitability of momentum strategies: An evaluation of alternative explanations. The Journal of finance 56, 699-720

Jiang, Z., Xu, D., Liang, J., 2017. A deep reinforcement learning framework for the financial portfolio management problem. arXiv preprint arXiv:1706.10059

Kim, J.H., Kim, W.C., Fabozzi, F.J., 2016. Portfolio selection with conservative short-selling. Finance Research Letters 18, 363-369

Laloux, L., Cizeau, P., Bouchaud, J.-P., Potters, M., 1999. Noise dressing of financial correlation matrices. Physical review letters 83, 1467

Liang, Z., Chen, H., Zhu, J., Jiang, K., Li, Y., 2018. Adversarial deep reinforcement learning in portfolio management. arXiv preprint arXiv:1808.09940

Lillicrap, T.P., Hunt, J.J., Pritzel, A., Heess, N., Erez, T., Tassa, Y., Silver, D., Wierstra, D., 2015. Continuous control with deep reinforcement learning. Computer Science

Liu, J., Stambaugh, R.F., Yuan, Y., 2019. Size and value in China. Journal of Financial Economics 134, 48-69

Mnih, V., Kavukcuoglu, K., Silver, D., Rusu, A.A., Veness, J., Bellemare, M.G., Graves, A., Riedmiller, M., Fidjeland, A.K., Ostrovski, G., 2015. Human-level control through deep reinforcement learning. nature 518, 529-533

Moody, J.E., Saffell, M., 1999. Reinforcement learning for trading. Advances in Neural Information Processing Systems, 917-923

Mullainathan, S., Spiess, J., 2017. Machine learning: an applied econometric approach. Journal of Economic Perspectives 31, 87-106

Qian, E.E., Hua, R.H., Sorensen, E.H., 2007. Quantitative equity portfolio management: modern techniques and applications. Chapman and Hall/CRC.

Ross, S.A., 2013. The arbitrage theory of capital asset pricing. In: Handbook of the fundamentals of financial decision making: Part I. World Scientific, pp. 11-30.

Sezer, O.B., Ozbayoglu, A.M., 2018. Algorithmic financial trading with deep convolutional neural networks: Time series to image conversion approach. Applied Soft Computing 70, 525-538

Silver, D., Lever, G., Heess, N., Degris, T., Wierstra, D., Riedmiller, M., 2014. Deterministic policy gradient algorithms. In: International conference on machine learning, pp. 387-395. PMLR

Simonyan, K., Zisserman, A., 2014. Very Deep Convolutional Networks for Large-Scale Image Recognition. Computer Science

Song, Q., Liu, A., Yang, S.Y., 2017. Stock portfolio selection using learning-to-rank algorithms with news sentiment. Neurocomputing 264, 20-28

Sutton, R.S., Barto, A.G., 2018. Reinforcement learning: An introduction. MIT press.

wassname, 2019. rl-portfolio-management. URL https://github.com/wassname/rl-portfolio-management

Zhang, Z., Zohren, S., Roberts, S., 2020. Deep reinforcement learning for trading. The Journal of Financial Data Science 2, 25-40