

Deep Reinforcement Learning for Portfolio Management

Yue Ma*

Northern Illinois University, DeKalb, IL

Ziping Liu[†] and Charles D. McAllister[†]

Southeast Missouri State University, Cape Girardeau, MO

Abstract

This paper discusses how to build deep reinforcement learning (DRL) agents to determine the allocation of money for assets in a portfolio so that the maximum return can be gained. The policy gradient method from reinforcement learning and convolutional neural network/recurrent neural network/convolutional neural network concatenated with the recurrent neural network from deep learning are combined to build the agents. With the proposed models, three types of portfolios are tested: stocks portfolio which has a positive influence due to the Covid-19; stocks portfolio which has a negative influence due to the Covid-19; and stocks, cryptocurrency combined portfolio which are randomly selected. The performance of our DRL agents is compared with that of equal-weighted agent and all the money fully invested in one-stock agents. All of our DRL agents showed the best performance on the randomly selected portfolio, which has an overall stable increasing trend. In addition, the performance of a linear regression model is also tested with the random selected portfolio, and it shows a poor result compared to other agents.

Key Words: Deep reinforcement learning; financial portfolio management; machine learning; neural networks; stock.

1 Introduction

Deep reinforcement learning shows its advantages in various fields, such as gaming [17, 7, 5], transportation [18, 11, 8] and the medical field [1, 24, 20]. These successful applications have inspired more and more researchers to apply this framework to other fields. For example, use in the financial market can expect to achieve the same success.

Investment in the financial market plays an important role in a global economy, especially the stock market investment [3]. However, the financial market is not only large and volatile [16], but also can be heavily impacted by all sorts of unpredictable natural [9], social [6], and political factors [13]. Hence, the impacts of actions taken on the financial market will not be measured easily and quickly, which adds difficulties on the building of effective financial models. In deep reinforcement learning, the goal is to maximize the reward. Because of this special feature, deep reinforcement learning has gained a lot of attention in the field of gaming. Bearing the goal to

maximize the overall return in financial portfolio management, which is similar to gaming, deep reinforcement learning should be applicable to manage financial portfolios. Considering the impact of the volatility of Covid-19 on the financial market, three different types of portfolios were selected to test the hypothesis in this paper.

This paper investigates the applicability of deploying deep reinforcement learning agents in the design of portfolio management agents such that weights in the portfolio can be allocated automatically to maximize the overall return. Convolutional neural network and recurrent neural network are investigated in the construction of the policy network, and the performances between agents with different network structures are compared.

2 Literature Review

Financial portfolio management is about the allocation of financial assets to meet the financial goals of investors. Generally, most investors have the same goal: maximize the expected return. In many different fields, such as medical diagnosis, speech recognition, and market clustering analysis, machine learning methods have proved to be powerful analytical tools. The above successful cases have motivated researchers to design machine learning models for portfolio management.

In [4], Gah-Yi Ban et al. introduced performance-based regularization (PBR) and performance-based cross-validation for the portfolio optimization problem and investigated them in detail. Tianello [2] used regression analysis to predict how each type of asset performs on a risk adjusted basis, and then allocates funds to those assets proportionately. Compared to classic machine learning, deep reinforcement learning prefers to achieve a pursued long-term result and portfolio management has a similar goal to maximize the rewards. The aforementioned similar goal has driven many researchers to develop and implement different deep reinforcement learning methods to manage financial portfolios. Liang et al. in [23] used deep reinforcement learning algorithms with continuous action space for the asset allocation. The agents were tested with risk adjusted accumulative portfolio value as objective function and the inputs took the combinations of different features. In [22], Liang et al. implemented three state-of-art continuous reinforcement learning algorithms, Deep Deterministic Policy Gradient (DDPG), Proximal Policy Optimization (PPO) and Policy Gradient (PG) in portfolio management and proposed

*e-mail: yuema@niu.edu

[†]e-mail: zliu,cdmcallister@semo.edu

an adversarial training method and showed that it can greatly improve the training efficiency and significantly raise average daily return and Sharpe ratio in backtests. Zhang et al. proposed a cost-sensitive portfolio selection method with deep reinforcement learning in [21], where a novel two-stream portfolio policy network was devised to extract both price series patterns and asset correlations. The proposed method also includes a new cost-sensitive reward function to maximize the accumulated return as well as to constrain both costs via reinforcement learning. Kanwar in [10], applied model free reinforcement learning algorithm which is a computationally efficient algorithm that directly estimates the optimal policy or value function through algorithms such as policy iteration or value iteration, with emphasis on policy gradient and Actor Critic Methods to build the trading agent in order to maximize its overall return. Opposite to [14], Yu et al. used a model-based approach which allows agents to plan ahead a range of possible choices to build the deep reinforcement learning agents. Jiang et al. also presented a model-free convolutional neural network with historic prices on a set of financial assets as its input, which outputs portfolio weights of the set in [26]. In [25], Jiang et al. proposed portfolio vector machine (PVM) which is an online stochastic batch learning scheme. In the proposed scheme, PVM works together with the policy network via convolutional neural network (CNN), recurrent neural network (RNN) and the Long Short-Term Memory (LSTM) to re-balance the cryptocurrency portfolio every 30 minutes. Selim Amrouni et al. in [15] extended the work of Jiang et al., applied the PVM to trade the daily stock and cryptocurrency data with a daily rebalance.

This research was motivated by both Jiang et al. [25] and Amrouni et al. [15], and used their work as reference for the basic policy network architecture of a convolutional neural network but implemented with a recurrent neural network and a mixed neural network (RNN+CNN) to build agents and used various financial asset combinations for testing. Since Covid-19, the stock market has experienced the most volatile ups and downs. Meanwhile, there is less related research taking Covid-19 into consideration regarding deep reinforcement learning on portfolio management. In order to gain insights, this research selected assets with different combinations to build the financial portfolio on testing the performance of deep reinforcement learning agents. To compare the performance of deep reinforcement learning with classic machine learning approaches in portfolio management, linear regression is also tested for weights allocation.

3 Methodology

3.1 Deep Learning

Deep learning is one branch among the many fields of machine learning, and it is based on artificial neural networks. Convolutional neural networks are designed to work with grid-structured inputs, which have strong spatial dependencies in local regions of the grid. The most obvious example of grid-

structured data is a 2-dimensional image, while other forms of sequential data like text, time-series, and sequences can also be considered special cases of grid-structured data with various types of relationships among adjacent items. The core concept of convolutional neural network is to simplify complicated questions. It reduced the dimensionality of a large number of parameters into a small number of parameters before processing.

Certain data types such as time-series, text, and biological data contain sequential dependencies among the attributes. Taking time-series data set as example, the values on successive timestamps are closely related to one another. If one uses the values of these timestamps as independent features, then key information about the relationships among the values of these timestamps is lost. For example, the value of a time-series at time t is closely related to its values in the previous window. However, this information is lost when the values at individual timestamps are treated independently of one another. A recurrent neural network, also known as RNN, can solve this problem. RNN is a class of neural networks which allow previous outputs to be used as inputs while having hidden states. The main advantage of recurrent neural network is that it can model sequence data, which is also the reason why RNN models are mostly being used by processing natural language.

3.2 Deep Reinforcement Learning

Learning can be divided into supervised learning, unsupervised learning and reinforcement learning. Supervised learning models receive features (X) and labels (y) as training data, using them to get predicted values (\hat{y}); while unsupervised learning models receive only features (X) and to do further analysis; but reinforcement learning can be different with these learning types. Reinforcement learning (RL) [12] is the study of how an agent can interact with its environment to learn a policy which maximizes expected cumulative rewards for a task. Environment in DRL represents the outside world the agent interacts with. Agent is the individual which makes decisions to decide what actions to take in the environment. State is the current situation of agent, what the agent observed in the environment. Action is what the agent can do in the environment, the input provides to the environment. Reward is a feedback signal from the environment. Policy is used by an agent as a strategy to decide what actions should be taken given state s . Policy had deterministic policy and stochastic policy.

3.3 Financial Portfolio Management

Markowitz portfolio theory, which is also called Modern portfolio theory (MPT), is a theory on how risk-averse investors can construct portfolios to maximize expected return based on a given level of market risk. The investor can decide the risk or the expected rule, but there is a rule to follow, which is: More risk will be taken, when more money is expected in return. How we compute the expected rate of return for an individual investment

as shown in the formula below:

$$E(R_{port}) = \sum_{i=1}^n w_i R_i$$

where:

w_i = the weight of an individual asset in the portfolio, or the percent of the portfolio in Asset i

R_i = the expected rate of return for Asset i

The expected rate of return for a portfolio of investments is simply the weighted average of the expected rates of return for the individual investments in the portfolio.

4 Research Methodology

4.1 Overall Architecture

The concept of using efficient portfolio management agents to maximize the expected return and to avoid risk follows Markowitz portfolio theory. In this paper, a model with deep reinforcement learning (DRL) agents is studied to determine the allocation of money for assets in a portfolio so that the maximum return can be gained. The policy gradient method from reinforcement learning and convolutional neural network/recurrent neural network/recurrent neural network concatenated with the convolutional neural network from deep learning are combined together to build the agents. Figure 1 shows the overall architecture of the deep reinforcement learning agent. Data should be pre-processed based on different choices of neural networks, and then the pre-processed data is passed to the neural networks. Softmax is applied in order to get the weight distribution of the assets. Input data, neural networks and the weight distribution formed the policy network, and reward is calculated based on the calculated weights. Reward is next fed into an objective function to update the policy network.

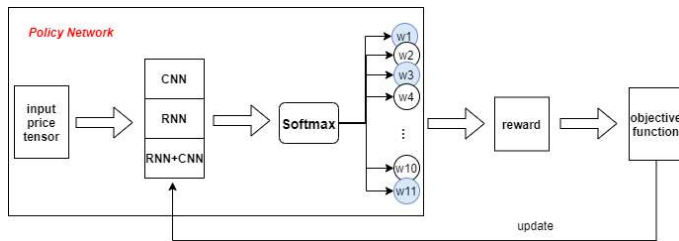


Figure 1: Overall deep reinforcement learning architecture

4.2 Data Collection

Considering the impact of the volatility of Covid-19 on the financial market, in this research, various combinations of stocks and cryptocurrency are selected to be tested, and all the data were collected from Yahoo! Finance [19].

4.2.1 Negatively Influenced Portfolio. The stocks which are heavily impacted by the Covid-19 are selected as the negatively influenced portfolio. These stocks are from the following companies: American Airlines Group

(AAL), Delta Airlines (DAL), Southeast Airlines (LUV), United Airlines Holdings (UAL), Spirit Airlines (SAVE), Hilton Worldwide Holdings (HLT), Marriott International (MAR), Choice Hotels International (CHH), InterContinental Hotels Group PLC (IHG), and International Consolidated Airlines Group (ICAGY). Figure 2 shows the daily closing price evolution for the negatively influenced portfolio from January 25, 2016 to January 25, 2021. X-axis represents the days, and y-axis represents close price (US dollars).

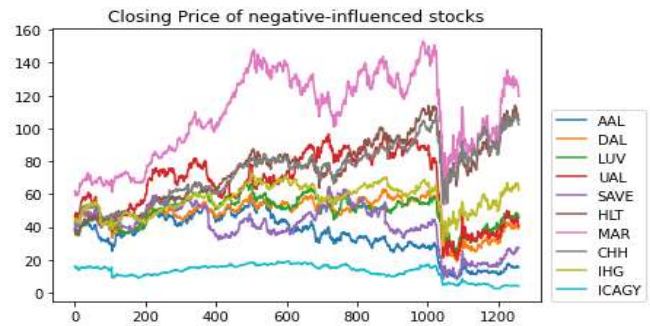


Figure 2: Negatively influenced stocks

4.2.2 Positively Influenced Portfolio. The stocks which have a positive influence because of the Covid-19 are selected as the positive influenced portfolio. These stocks are from the following companies: Alibaba Group Holding Limited (BABA), Walmart (WMT), eBay (EBAY), Best Buy (BBY), Allied Healthcare Products (AHPI), Amedisys (AMED), Viridian Therapeutics (MGEN), Meridian Bioscience (VIVO), NanoViricides (NNVC), Netflix (NFLX). Figure 3 shows the daily closing price evolution for the positive influenced portfolio from January 25, 2016 to January 25, 2021. X-axis represents the days, and y-axis represents close price (US dollars).

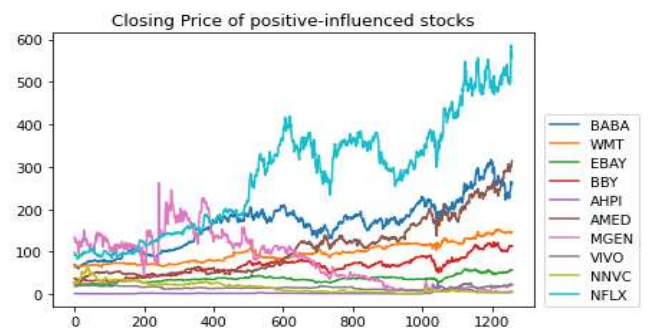


Figure 3: Positively influenced stocks

4.2.3 Randomly Selected Portfolio. To consider neither the positive impact nor the negative impact of Covid-19, seven random US stocks and three cryptocurrencies are selected from the following companies: American Airlines Group (AAL), Facebook (FB), Bank of America (BAC), Best Buy (BBY), Alibaba Group (BABA), Axon Enterprise (AAXN), BlackRock (BLK), and cryptocurrency of Ethereum (ETH), Ripple (XRP),

Litecoin (LTC). Figure 4 shows the daily closing price evolution for the positive influenced portfolio from January 25, 2016 to January 25, 2021. X-axis represents the days, and y-axis represents close price (US dollars).

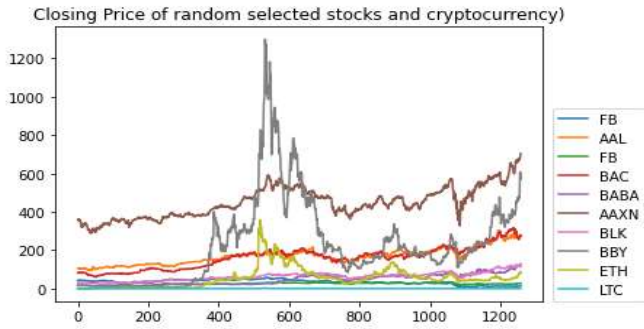


Figure 4: Randomly selected stocks and cryptocurrency

4.3 Data Preprocessing

For both stocks and cryptocurrencies, data between January 25, 2016 to January 25, 2021 were downloaded and they span totally five years, 1258 days for stocks and 1828 days for cryptocurrencies. The reason why only five years data were selected is that cryptocurrency does not have historical data before 2009 when the first cryptocurrency was adopted as bitcoin. Since then, more and more newly generated cryptocurrencies have come out. To make more use of stable cryptocurrency data in the model, our research restricted the data selection to the recent five years.

There was a problem when data for the random selected portfolio was processed, since the length of two types of investment data do not have the same length and has a gap of 569 trading days. In this research, for the days which cryptocurrency trading occurred while the stock trading did not occur, the price information for cryptocurrency was removed.

The data is divided as below:

- Total days: 1258
- Training days: 60% of total days
- Validation days: 20% of total days
- Test days: 20% of total days

Normal indicators for stock include "Date," "Open," "High," "Low," "Close," "Adj Close" and "Volume," but for this research, only "Open," "High," "Low" and "Close" were used. Strictly speaking, the market for cryptocurrency never closes, but there is still data available as "Open price" and "Close price" provided for cryptocurrency. "Open" generally refers to the price at 12:01 AM UTC of any given day and "close" generally refers to the price at 11:59 PM UTC of any given day. Hence, stock data and cryptocurrency data can be mixed together.

Figure 5 shows the data set weight distribution for this research, testing data set contains 253 days.

4.3.1 CNN Price Tensor Pre-Processing. The input data are arranged into 3D in shape of (4,10,1258) as shown in Figure

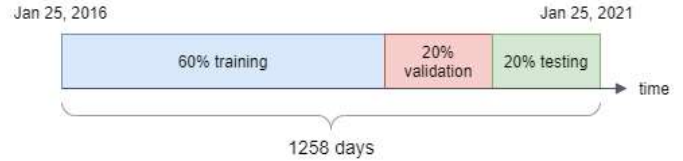


Figure 5: Dataset weight distribution

6. The first dimension 4 represents the number of features, and they are:

- Close(t-1)/Open(t-1)
- High(t-1)/Open(t-1)
- Low(t-1)/Open(t-1)
- Open(t)/Open(t-1)

Here the prices are normalized to the same scale. The second value for the dimension is 10, where it represents ten asset items (stocks+cryptocurrency) in the portfolio. The third value for the dimension is 1258, which represents the time steps (days). 1258 was sliced to 50 days each before passing to the network.

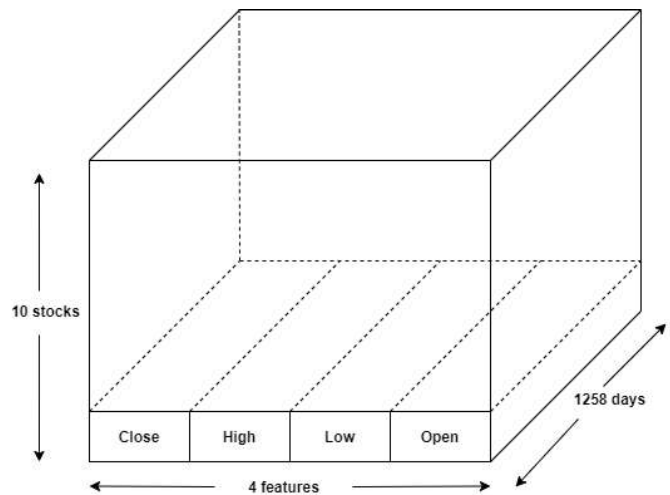


Figure 6: CNN data pre-processing

4.3.2 RNN and RNN+CNN Price Tensor Pre-Processing.

For the architecture of recurrent neural network, and the architecture for the mixed (RNN+CNN) neural network, these two networks have the same input shape. There is only one feature (the mean value for the open price, close price, high price and the low price) used. The mean prices are passed to MinMaxScaler which helps to make all the value in the same scale.

$$MeanPrice = \frac{OpenPrice + ClosePrice + HighPrice + LowPrice}{4}$$

After data scaling, data was grouped to 50 days each. In the previous section's CNN policy network, the length of tensor was set to be 50 days. In order to do a comparative analysis with the CNN policy network, the tensor should be processed as the same

length for batch processing. After transposing and dimension expanding, pre-processed input data has the shape (25,1,10,50), which is illustrated in Figure 7. The diagram shows that there are 25 batches of (1,10,50), for each step, one batch (1,10,50) from (25,1,10,50) was passed to the policy network to calculate the weight distribution. 1 means there is only one feature–mean value, 10 represents 10 stocks, and 50 is the length of the tensor(50 days).

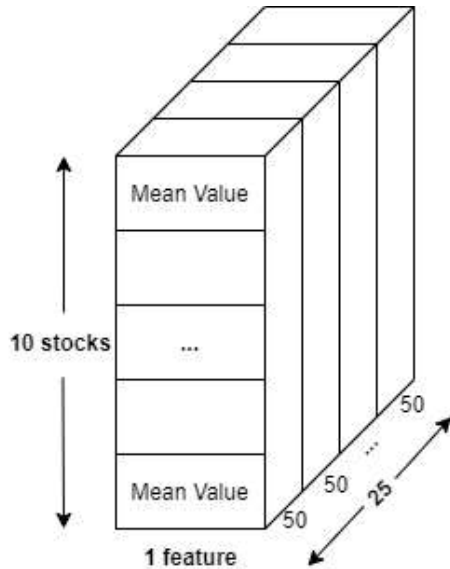


Figure 7: RNN and RNN+CNN data pre-processing

4.3.3 Linear Regression Price Tensor Pre-Processing. Data processing for linear regression is similar as for RNN. However, in linear regression, it needs X (features) and y (target values). In this research, the initial investment and mean value (open price, close price, high price and low price) of 10 stocks are chosen for X , which has 11 features in total. Initial investment is set to 10000 for every data time, and y is the portfolio value which is randomly generated between values of (10000,13000). After setting the portfolio value for each data time, the initial investment of each day will be adjusted to the previous portfolio value. Figure 8 shows the shape of pre-processing data for linear regression model.

Money	Stock1	Stock2	...	Stock10	Portfolio Value
10000	12.3	25.6		28.8	12500
12500	12.9	22.6		15.8	12907
12907	17.8	11.1		15.9	11008
...

Figure 8: Linear regression data pre-processing

4.4 Reinforcement Learning Framework

Markov Decision Process (MDP) provided the framework for describing the reinforcement learning problem. The state of the agent is the input matrix X_t and previous portfolio weights at time $t-1$.

$$state = S_t = (X_t, w_{t-1})$$

The action of the agent is the weight vector of portfolio at time t .

$$action = w_t = (w_1, \dots, w_n)$$

The reward function is defined such as it is the agent’s return minus a baseline return minus a term proportional to the maximum of the weight (this term is set up to make the agent avoid investing fully in one stock), and baseline is an equal weighted agent which divides the investment into the same proportion.

$$r_t = \sum w'_i y_i - \frac{1}{m} \sum y_i - \alpha max(w_1, \dots, w_m)$$

where:

y_i is the return of single stock or cryptocurrency.

$\sum w'_i y_i$ is the sum of each stock or cryptocurrency’s return. The formula for calculating the expected rates of return was used as reference.

m is number of items(stocks and cryptocurrency) in the portfolio.

α is the parameter of regularization, which can avoid investing fully in one stock.

The input price tensor is the state for the deep reinforcement learning model, and action is the weight distribution which was calculated by the policy network. Reward is the reward function in the overall architecture which helped to build the objective function. The three components work together to train the model to pursue maximum return.

4.5 Network Structure

The convolutional neural network structure was designed by Amrouni et al [15]. Motivated by the work in [15], this research also studied recurrent network, and a mixed network with RNN+CNN. The network structure design for the recurrent neural network and the mixed network is a new attempt for the training process which will be introduced in the next section.

4.5.1 Convolutional Neural Network. The detailed parameters for the convolutional neural network model are listed in Table 1. The price tensor starts in a shape of $[None, nb_feature, m, n]$, where $nb_features$ represents number of features, m represents the number of stocks and cryptocurrency, n is the length of tensor. It is then transposed to $[None, n, m, nb_features]$ where it is used as the input for the convolutional neural network.

The policy network which uses convolutional neural network comprises input layer and three convolution layers where conv1 extracted and integrated effective features to 2 kernels, conv2 kept on the process and output 20 kernels, and in the end the

conv3 will output the last feature map, which concatenated with the cash bias. After the concatenation with the cash bias, the tensor got squeezed, and in the end, softmax function is used to get the probability distribution which is also the weight distribution of the portfolio. The convolutional neural network model is shown in Figure 9. Batch size was set to 50 in this research, which is a common number picked by most researchers and it can also be changed to another value.

The output of the network is then used to calculate the adjusted reward. At first it needs to calculate the instantaneous reward as shown in the formula below:

$$InstantaneousReward = \frac{PortfolioValue - PrevPortfolioValue}{PrevPortfolioValue}$$

In the formula, the trading cost is also taken into consideration when the portfolio value is calculated.

Considering the formulation of the objective function, since there is no maximization function for the reward, minimizing the negative of the original value is equivalent to maximizing the original value; therefore, minimize (-adjusted reward) was chosen to maximize the reward over the batch.

Table 1: Parameters for convolutional neural network model

Input	(?,4,10,50)	
Conv1	(?,50,10,2)	activation=tf.nn.relu, filters=2, strides=(1, 1), kernel_size=(1,3), padding='same'
Conv2	(?,1,10,20)	activation=tf.nn.relu, filters=20, strides=(50, 1), kernel_size=(1, 50), padding='same'
Conv3	(?,1,10,1)	activation=tf.nn.relu, filters=1, strides=(21, 1), kernel_size=(1, 1), padding='same'
Tensor4	(?,1,11,1)	Concatenate conv3 with cash_bias
Squeezed_tensor	(?,11)	
Action	(?,11)	tf.nn.softmax

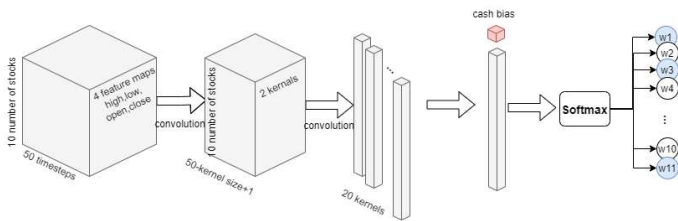


Figure 9: Convolutional neural network model

4.5.2 Recurrent Neural Network. The detailed parameters for the recurrent neural network model are listed in Table 2. The tensor which is the input to the network has the shape of (None, m, n), where m is the number of stocks and n is the length (time step) of the tensor. It has two recurrent layers. For the first recurrent layer, the number of recurrent units was set to 20, and it came from the inspiration from figure 4.9's CNN of filters=20. For the second recurrent layer, the number of the recurrent units

was set to 1. The output dimension after the two recurrent layers was expanded to 4D. And after that it can be concatenated with a 4D-shape cash-bias. The following steps such as to squeeze the tensor, to use softmax to get the probability distribution, and to calculate the reward are the same as in the convolutional neural network. The recurrent neural network model is shown in Figure 10.

Table 2: Parameters for recurrent neural network model

Input	(?,10,50)	
Rnn1	(?,10,20)	num_units=20, activation=tf.nn.relu
Rnn2	(?,10,1)	num_units=1, activation=tf.nn.relu
Outputs(expand dimension)	(?,1,10,1)	
Tensor4	(?,1,11,1)	Concatenate Outputs with cash_bias
Squeezed_tensor	(?,11)	
Action	(?,11)	tf.nn.softmax

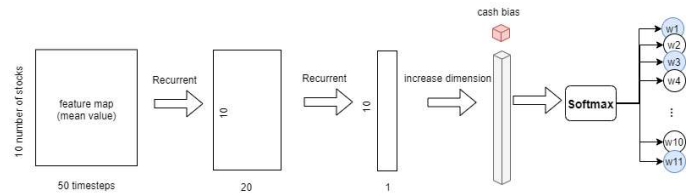


Figure 10: Recurrent neural network model

4.5.3 Mixed Neural Network. The mixed neural network model combines both the convolution layers and recurrent layers when processing the price tensor. Table 3 shows the details of the parameters. The input tensor has the same shape as the one of the recurrent neural networks and it is (None,10,50). First, it is passed to the recurrent layer, which has 20 recurrent units. The output of the recurrent layer's dimension is then expanded in order to be processed by the convolutional neural network. For the convolution layer, there is one filter. After the dimension expansion, all the following steps are the same as the convolution neural network's. The mixed neural network structure is as the Figure 11 below.

Table 3: Parameters for CNN+RNN model

Input	(?,10,50)	
Rnn1	(?,10,20)	num_units=20, activation=tf.nn.relu
Outputs(expand dimension)	(?,1,10,20)	
Conv3	(?,1,10,1)	activation=tf.nn.relu, filters=1, strides=(21, 1), kernel_size=(1, 1), padding='same'
Tensor4	(?,1,11,1)	Concatenate conv3 with cash_bias
Squeezed_tensor	(?,11)	
Action	(?,11)	tf.nn.softmax

4.5.4 Linear Regression Model. The training data is passed to the standard scaler first to be processed into the same scale.

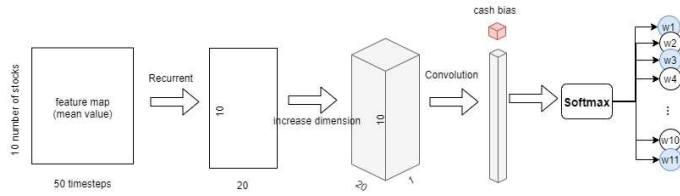


Figure 11: Mixed neural network model

Next, the processed data passes into the linear regression model in order to get the coefficient list which is also the weight list. The accumulated value for each value in the portfolio weight list should be 1, so the proportional value is calculated. The linear regression model is shown in Figure 12.

$$sum_coef = |coef_1| + |coef_2| + \dots + |coef_{11}|$$

$$weight_list = \left[\frac{|coef_1|}{sum_coef}, \frac{|coef_2|}{sum_coef}, \dots, \frac{|coef_{11}|}{sum_coef} \right]$$

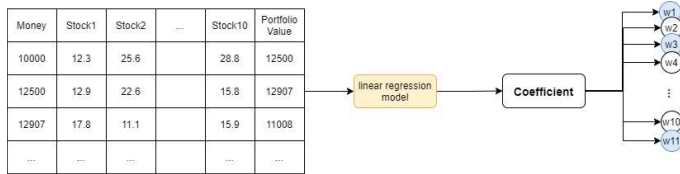


Figure 12: Linear regression model

4.6 Network Training Process

Portfolio Vector Memory was introduced by Jiang et al. in [9], and was implemented by Selim Amrouni in [10] for convolutional neural network model training. Figure 13 shows the CNN model training process. For each episode (total of 2 episodes in this paper), the PVM was initialized with training parameters. For each batch (total of 10 batches in this paper, 1 batch has 50 days), the starting point of the batch, initial portfolio value and initial weights were drawn from the PVM at time t-1. Based on the starting point, a price tensor which has the length of 50 was passed to the policy network to get the weight distribution. The weight distribution which is also the action wrote to the PVM at time t and made the agents step to the next day. The whole process recursively performed until all the batches and episodes ended.

This paper also implemented the PVM with the recurrent neural network model and the mixed RNN+CNN model. Figure 14 shows the process. The main difference between them is that for RNN and the mixed model, the price tensor retrieval and days traversal are separated. The starting point which was fed from the PVM was generated from the RNN pre-processed data, instead of generating the price tensor from CNN pre-processed data. CNN pre-processed data was here to help traverse the days. One price tensor which has length of 50 was retrieved and sent to the policy network to get the weight distribution. The remaining steps were similar as the CNN's.

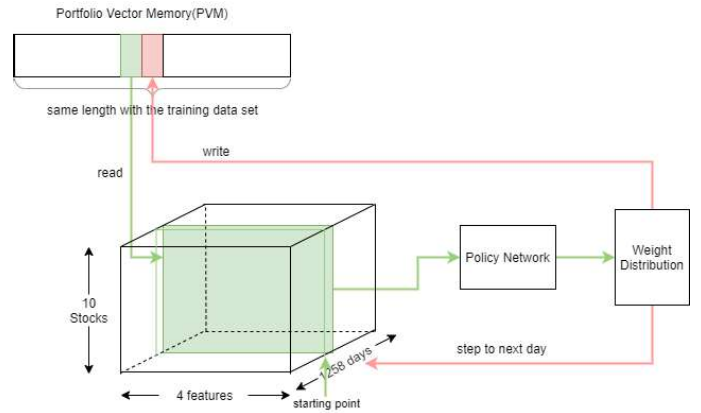


Figure 13: CNN training process

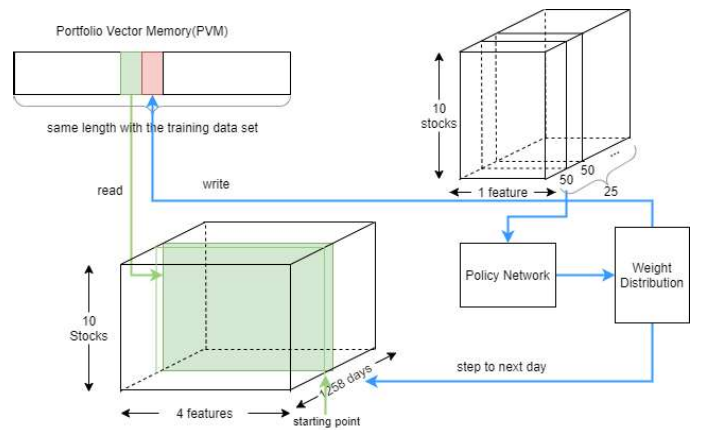


Figure 14: RNN and RNN+CNN training process

5 Results

5.1 Overall Results

Three different types of portfolio and three different networks were tested. For linear regression, the random selected portfolio was tested. The portfolio value evolution, min and max value, which represents minimum and maximum portfolio value, and the weight evolution during testing are presented in the following subsections. For testing, result figures are presented in the next subsections, except three deep reinforcement learning agents, where are other financial portfolio management agents being tested. Equal-weighted agent is an agent used for all the financial assets. Initial investment is 10000 dollars which acts as a baseline to show a clear performance comparison. For the other "full stock AAL, DAL, LUV..." means the money fully invest on the single stock.

5.1.1 Negatively Influenced Portfolio. Figure 15 shows the negatively influenced portfolio value evolution for 253 days. For both three types of deep reinforcement learning agents, there exist certain periods in which portfolio value is greater than the initial investment. But overall, they do not show an increasing trend. Among the three models, RNN model has

the best performance(most return) at the end of the testing day, however it still loses 45.58% off the initial investment. Table 4 shows the maximum portfolio value and the minimum portfolio value for testing data set, CNN reached 10749 dollars which is the maximum portfolio value compared to the other two deep reinforcement learning agents. Tables 5, 6 and 7 shows the weight evolution every 50 days during testing for these three deep reinforcement learning agents.

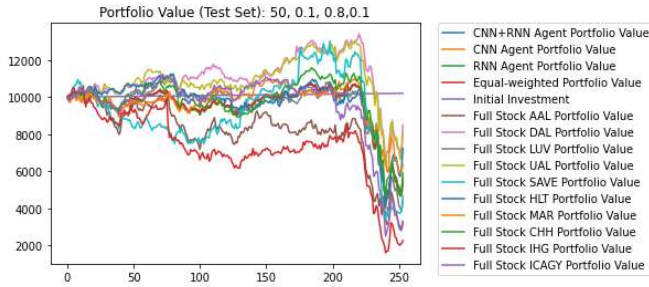


Figure 15: Negatively influenced portfolio

Table 4: Negatively influenced portfolio

	RNN	CNN	RNN+CNN
MAX	10744.384170389269	10749.582654440033	10746.332953411205
MIN	4296.119453575351	4257.566520880737	4282.784544420906

Table 5: Weight evolution during testing of CNN network

	Mo ney	AAL	DAL	LUV	UAL	SAV E	HLT	MA R	CHH	IHG	ICA GY	Portf olioV alue
St art	1	0	0	0	0	0	0	0	0	0	0	1000 0
50 da ys	0.04 202 312	0.09 757 042	0.09 566 913	0.09 583 758	0.09 606 861	0.09 529 205	0.09 550 374	0.09 521 315	0.09 623 07	0.09 488 089	0.09 571 061	9869
10 0d ay s	0.04 302 808	0.09 447 327	0.09 794 327	0.09 643 734	0.09 649 739	0.09 397 549	0.09 667 79	0.09 640 097	0.09 637 725	0.09 359 015	0.09 459 889	9043
15 0d ay s	0.04 197 768	0.09 555 684	0.09 728 851	0.09 613 146	0.09 501 63	0.09 509 232	0.09 558 895	0.09 559 92	0.09 481 207	0.09 690 978	0.09 602 688	1011 1
20 0d ay s	0.04 253 444	0.09 510 022	0.09 617 578	0.09 571 166	0.09 539 249	0.09 750 834	0.09 640 942	0.09 539 837	0.09 467 928	0.09 564 289	0.09 544 712	1047 1
25 0d ay s	0.04 445 002	0.09 371 601	0.10 153 517	0.09 098 017	0.09 698 123	0.09 955 532	0.09 356 5	0.09 537 336	0.10 010 976	0.09 723 239	0.08 650 158	4815

5.1.2 Positively Influenced Portfolio. Figure 16 shows the positively influenced portfolio value evolution for 253 days. Due to the explosive increasing trend for the full stock on BABA portfolio value, the graph does not show a clear trend for the

Table 6: Weight evolution during testing of RNN network

	Mo ney	AAL	DAL	LUV	UAL	SAV E	HLT	MA R	CHH	IHG	ICA GY	Portf olioV alue
St art	1	0	0	0	0	0	0	0	0	0	0	1000 0
50 da ys	0.05 115 532	0.09 664 03	0.09 475 713	0.09 492 398	0.09 515 281	0.09 438 365	0.09 459 332	0.09 430 55	0.09 531 335	0.09 397 641	0.09 479 822	9870
10 0d ay s	0.05 236 406	0.09 355 161	0.09 698 776	0.09 549 653	0.09 555 599	0.09 305 869	0.09 573 473	0.09 546 051	0.09 543 702	0.09 267 71	0.09 367 601	9052
15 0d ay s	0.05 109 973	0.09 464 697	0.09 636 215	0.09 521 612	0.09 411 158	0.09 418 687	0.09 467 877	0.09 468 893	0.09 390 93	0.09 598 703	0.09 511 254	1011 1
20 0d ay s	0.05 176 996	0.09 418 29	0.09 524 809	0.09 478 844	0.09 447 235	0.09 656 779	0.09 547 947	0.09 447 817	0.09 376 602	0.09 472 033	0.09 452 645	1046 9
25 0d ay s	0.05 408 524	0.09 277 103	0.10 051 135	0.09 006 277	0.09 600 332	0.09 855 146	0.09 262 155	0.09 441 167	0.09 910 031	0.09 625 195	0.08 562 935	4854

Table 7: Weight evolution during testing of RNN+CNN network

	Mo ney	AAL	DAL	LUV	UAL	SAV E	HLT	MA R	CHH	IHG	ICA GY	Portf olioV alue
St art	1	0	0	0	0	0	0	0	0	0	0	1000 0
50 da ys	0.04 798 819	0.09 696 288	0.09 507 342	0.09 524 083	0.09 547 042	0.09 469 869	0.09 490 906	0.09 462 028	0.09 563 149	0.09 429 009	0.09 511 465	9870
10 0d ay s	0.04 912 673	0.09 387 12	0.09 731 909	0.09 582 276	0.09 588 243	0.09 337 659	0.09 606 178	0.09 578 662	0.09 576 306	0.09 299 371	0.09 399 602	9049
15 0d ay s	0.04 793 61	0.09 496 253	0.09 668 343	0.09 553 357	0.09 442 535	0.09 450 089	0.09 499 443	0.09 500 462	0.09 422 239	0.09 630 705	0.09 542 965	1011 1
20 0d ay s	0.04 856 724	0.09 450 101	0.09 556 98	0.09 510 86	0.09 479 144	0.09 689 396	0.09 580 196	0.09 479 728	0.09 408 273	0.09 504 026	0.09 484 572	1047 0
25 0d ay s	0.05 074 487	0.09 309 864	0.10 086 629	0.09 038 082	0.09 634 235	0.09 889 948	0.09 294 863	0.09 474 508	0.09 945 027	0.09 659 185	0.08 593 173	4841

other agents. Actually, for the three types of deep reinforcement learning agents, during more than half of the testing days, the portfolio value seldom exceeds the initial investment, but at the end of the testing days, all agents show an increasing trend. For CNN model which has the best performance (most return) at end of the testing days, it gains 110.87% on the initial investment. Table 8 shows the maximum portfolio value and the minimum portfolio value for the testing data set, CNN reached 22784.7 dollars which is the maximum portfolio value compared to the other two deep reinforcement learning agents. Tables 9, 10 and

11 show the weight evolution during testing for these three deep reinforcement learning agents.

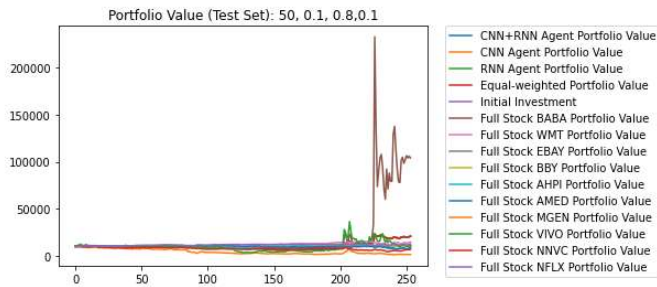


Figure 16: Positively influenced portfolio

Table 8: Positively influenced portfolio

	RNN	CNN	RNN+CNN
MAX	22696.510164433843	22784.707536769543	22604.74045890975
MIN	7621.239917016383	7607.848561432198	7635.092109920798

Table 9: Weight evolution during testing of CNN network

	Mo ney	BAB A	WM T	EBA Y	BBY	AHP I	AM ED	MG EN	VIV O	NN VC	NFL X	Portf olioV alue
St art	1	0	0	0	0	0	0	0	0	0	0	1000 0
50 da ys	0.04 045 456	0.09 421 504	0.09 656 75	0.09 705 374	0.09 904 441	0.09 744 479	0.09 691 681	0.09 784 962	0.08 812 246	0.09 594 888	0.09 638 219	9509
10 Od ays	0.04 151 763	0.09 493 009	0.09 508 458	0.09 717 265	0.09 681 507	0.09 684 529	0.09 672 782	0.09 159 769	0.09 385 134	0.09 725 982	0.09 819 801	8612
15 Od ays	0.04 063 14	0.09 541 863	0.09 672 963	0.09 640 315	0.09 806 951	0.09 753 956	0.09 593 144	0.08 993 582	0.09 915 782	0.09 297 967	0.09 720 337	8168
20 Od ays	0.04 009 868	0.09 856 441	0.09 598 776	0.09 656 359	0.09 593 668	0.09 655 176	0.09 554 054	0.08 520 287	0.10 651 74	0.09 353 385	0.09 550 245	9840
25 Od ays	0.04 078 758	0.10 094 165	0.09 346 066	0.09 511 324	0.09 493 542	0.09 390 534	0.09 331 918	0.09 932 903	0.08 692 816	0.10 371 155	0.09 756 818	1994 8

5.1.3 Randomly Selected Portfolio. Figure 17 shows the randomly selected portfolio value evolution for 100 days. For the three types of deep reinforcement learning agents, they all show a stable increasing trend. The deep reinforcement learning agents have the best performance on this type of portfolio compared to the other two portfolios. The RNN model, which has the best performance (most return) gained 47.72% initial investment for 253 days. Table 12 shows the maximum portfolio value and the minimum portfolio value during testing, RNN

Table 10: Weight evolution during testing of RNN network

	Mo ney	BAB A	WM T	EBA Y	BBY	AHP I	AM ED	MG EN	VIV O	NN VC	NFL X	Portf olioV alue
St art	1	0	0	0	0	0	0	0	0	0	0	1000 0
50 da ys	0.04 617 789	0.09 365 308	0.09 599 151	0.09 647 485	0.09 845 364	0.09 686 356	0.09 633 874	0.09 726 599	0.08 759 685	0.09 537 658	0.09 580 73	9512
10 Od ays	0.04 738 496	0.09 434 898	0.09 450 253	0.09 657 781	0.09 622 242	0.09 625 245	0.09 613 57	0.09 103 698	0.09 327 683	0.09 666 445	0.09 759 69	8621
15 Od ays	0.04 638 161	0.09 484 672	0.09 614 986	0.09 582 533	0.09 748 171	0.09 695 493	0.09 535 646	0.08 939 677	0.09 856 349	0.09 242 237	0.09 662 076	8179
20 Od ays	0.04 577 628	0.09 798 143	0.09 542 002	0.09 599 244	0.09 536 923	0.09 598 068	0.09 497 544	0.08 469 892	0.10 588 738	0.09 298 062	0.09 493 757	9843
25 Od ays	0.04 656 188	0.10 033 4	0.09 289 805	0.09 454 068	0.09 436 393	0.09 334 004	0.09 275 742	0.09 873 109	0.08 640 487	0.10 308 722	0.09 698 083	1988 8

Table 11: Weight evolution during testing of CNN+RNN network

	Mo ney	BAB A	WM T	EBA Y	BBY	AHP I	AM ED	MG EN	VIV O	NN VC	NFL X	Portf olioV alue
St art	1	0	0	0	0	0	0	0	0	0	0	1000 0
50 da ys	0.05 211 806	0.09 306 983	0.09 539 37	0.09 587 403	0.09 784 05	0.09 626 632	0.09 573 876	0.09 666 024	0.08 705 131	0.09 478 26	0.09 521 064	9515
10 Od ays	0.05 347 262	0.09 374 604	0.09 389 861	0.09 596 063	0.09 560 751	0.09 563 735	0.09 552 135	0.09 045 521	0.09 268 075	0.09 604 672	0.09 697 32	8630
15 Od ays	0.05 234 932	0.09 425 317	0.09 554 816	0.09 522 566	0.09 687 167	0.09 634 819	0.09 475 972	0.08 883 733	0.09 794 668	0.09 184 4	0.09 601 611	8190
20 Od ays	0.05 166 957	0.09 737 629	0.09 483 07	0.09 539 959	0.09 478 023	0.09 538 79	0.09 438 887	0.08 417 582	0.10 523 341	0.09 240 637	0.09 435 123	9846
25 Od ays	0.05 255 428	0.09 970 339	0.09 231 418	0.09 394 649	0.09 377 084	0.09 275 34	0.09 217 443	0.09 811 056	0.08 586 181	0.10 243 931	0.09 637 131	1982 5

reached 14772 dollars which is the maximum portfolio value compared to the other two deep reinforcement learning agents. Tables 13, 14 and 15 show the weight evolution during testing for these three deep reinforcement learning agents.

5.1.4 Linear Regression on Randomly Selected Portfolio. Since randomly selected portfolio shows a stable increasing performance, we tested linear regression using the randomly selected portfolio. The linear regression agent's result is compared with the equal-weighted portfolio value and the initial investment in Figure 18. From the comparison, it can be seen

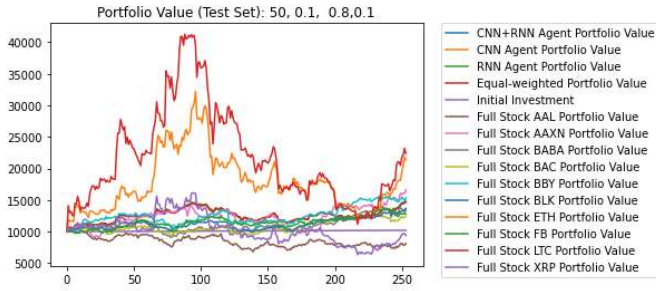


Figure 17: Randomly selected portfolio

Table 12: Randomly selected portfolio

	RNN	CNN	RNN+CNN
MAX	14772.48989557328	14679.754959892462	14764.415960353592
MIN	10000	10000	10000

that linear regression model shows a bad performance which has a clear downward trend. The weight evolution during testing is shown in Figure 19, the weight adjustment during testing is volatile for the linear regression agent.

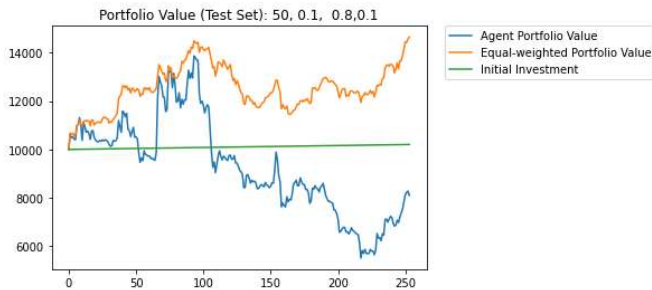


Figure 18: Linear regression on random selected portfolio

5.2 Portfolio Comparison Results

In this section, the performance results from the three different types of portfolios are discussed. Figure 20 shows the CNN model results for the three types of portfolio. Figure 21 shows the RNN model's results, and Figure 22 shows the RNN+CNN model's results. Among the three models, the negatively selected portfolio value presents a decreasing trend, while the positively selected portfolio value presents an explosive increasing trend around day 200, and the randomly selected portfolio presents an overall stable state with positive returns.

5.3 Activation Functions

The performances of activation functions were compared using random selected portfolio. Three different activation functions were compared: sigmoid, ReLU and tanh. X axis shows the 253 days, and the y axis is the portfolio value. Figure

Table 13: Weight evolution during testing of CNN network

	Mo ney	AAL	AAX N	BAB A	BAC	BBY	BLK	ETH	FB	LTC	XRP	Portf olioV alue
St art	1	0	0	0	0	0	0	0	0	0	0	1000 0
50 da ys	0.07 450 122	0.09 033 511	0.09 407 107	0.09 289 475	0.09 233 032	0.09 176 801	0.09 219 527	0.09 322 398	0.09 446 668	0.09 161 997	0.09 259 363	1242 9
10 Od ay s	0.07 540 939	0.09 131 607	0.09 225 253	0.09 355 897	0.09 278 396	0.09 257 833	0.09 294 329	0.09 269 756	0.09 382 129	0.09 080 885	0.09 182 977	1415 7
15 Od ay s	0.07 423 953	0.09 346 664	0.09 401 267	0.09 278 13	0.09 427 755	0.09 246 214	0.09 293 781	0.09 294 169	0.09 060 269	0.09 087 423	0.09 140 374	1233 1
20 Od ay s	0.07 616 618	0.09 397 48	0.09 440 434	0.09 629 674	0.09 472 538	0.09 163 484	0.09 343 745	0.08 661 038	0.09 466 618	0.08 670 246	0.09 138 125	1239 7
25 Od ay s	0.07 331 774	0.09 164 177	0.08 960 15	0.08 934 477	0.09 293 649	0.09 313 318	0.09 140 766	0.09 780 822	0.09 983 148	0.09 658 493	0.09 439 227	1447 5

Table 14: Weight evolution during testing of RNN network

	Mo ney	AAL	AAX N	BAB A	BAC	BBY	BLK	ETH	FB	LTC	XRP	Portf olioV alue
St art	1	0	0	0	0	0	0	0	0	0	0	1000 0
50 da ys	0.05 744 246	0.09 200 017	0.09 580 498	0.09 460 698	0.09 403 215	0.09 345 947	0.09 389 461	0.09 494 228	0.09 620 789	0.09 330 87	0.09 430 031	1247 7
10 Od ay s	0.05 815 7	0.09 301 998	0.09 397 392	0.09 530 473	0.09 451 527	0.09 430 579	0.09 467 756	0.09 442 725	0.09 557 194	0.09 250 329	0.09 354 326	1424 3
15 Od ay s	0.05 723 598	0.09 518 336	0.09 573 942	0.09 448 543	0.09 600 916	0.09 416 04	0.09 464 481	0.09 464 876	0.09 226 68	0.09 254 333	0.09 308 257	1237 2
20 Od ay s	0.05 874 964	0.09 574 646	0.09 618 409	0.09 811 217	0.09 651 118	0.09 336 238	0.09 519 897	0.08 824 32	0.09 645 087	0.08 833 702	0.09 310 402	1243 9
25 Od ay s	0.05 651 004	0.09 330 393	0.09 122 665	0.09 096 526	0.09 462 212	0.09 482 238	0.09 306 556	0.09 958 222	0.09 146 08	0.09 833 674	0.09 610 431	1456 3

23 shows the RNN activation results. Figure 24 shows the combined model activation result, and Figure 25 shows the CNN activation results.

5.4 Optimizer

The performances of various optimizers are compared using random selected portfolio. Four types of optimizers are compared: Adam, RMSProp, Gradient Descent and Adagrad. X axis shows the 253 days, and y axis is the portfolio value. For all the models, Adam optimizer and Gradient Descent optimizer

Table 15: Weight evolution during testing of RNN+CNN network

	Money	AAL	AAXN	BABA	BAC	BBY	BLK	ETH	FB	LTC	XRP	Portfolio Value
Start	1	0	0	0	0	0	0	0	0	0	0	10000
50 days	0.05	0.09	0.09	0.09	0.09	0.09	0.09	0.09	0.09	0.09	0.09	1247
10 days	896	185	565	445	388	330	374	478	605	315	414	2
15 days	482	157	025	417	027	852	295	894	25	8	801	1423
20 days	0.05	0.09	0.09	0.09	0.09	0.09	0.09	0.09	0.09	0.09	0.09	1236
25 days	969	286	382	514	436	415	452	427	541	235	339	9
30 days	696	789	027	89	073	16	276	286	568	205	031	1236
35 days	0.05	0.09	0.09	0.09	0.09	0.09	0.09	0.09	0.09	0.09	0.09	1236
40 days	875	503	558	433	585	400	449	449	211	239	293	9
45 days	333	016	533	335	463	885	248	642	83	438	275	1236
50 days	0.06	0.09	0.09	0.09	0.09	0.09	0.09	0.08	0.09	0.08	0.09	1243
55 days	030	558	602	795	635	320	504	809	629	819	295	5
60 days	452	829	52	01	175	815	171	743	154	109	021	1243
65 days	0.05	0.09	0.09	0.09	0.09	0.09	0.09	0.09	0.09	0.09	0.09	1455
70 days	800	315	108	082	447	467	291	942	131	818	595	5
75 days	961	563	165	068	173	167	765	394	543	044	156	1455

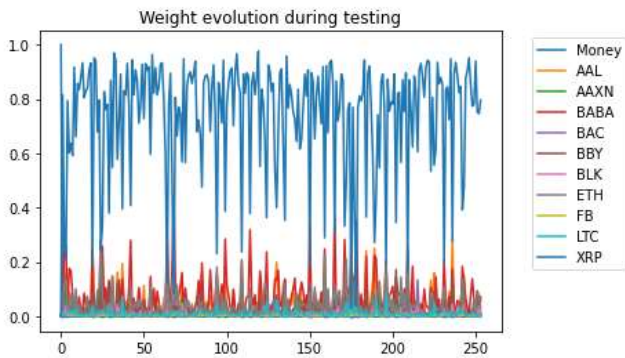


Figure 19: Weight evolution of linear regression

showed the best result. Figure 26 shows the RNN activation results, Figure 27 shows the combined model activation result, and Figure 28 shows the CNN activation results.

5.5 Learning Rate

The performances of different learning rates are also compared. Four different learning rates with 0.1, 0.01, 0.001 and 0.0001 were calculated. Figure 29 shows the RNN learning rate results. Figure 30 shows the combined model learning rate result, and Figure 31 shows the CNN learning rate results.

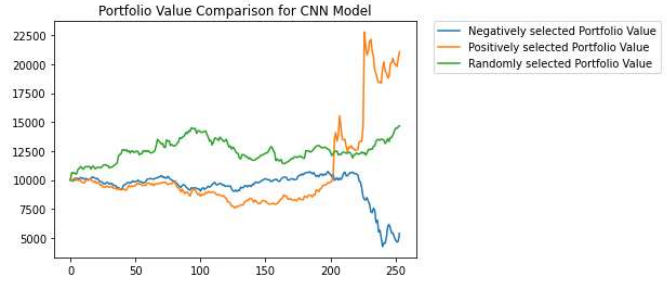


Figure 20: CNN model results

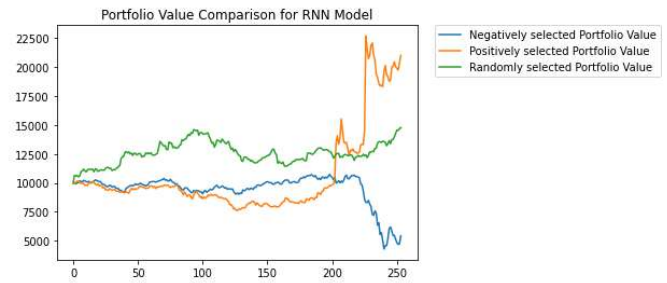


Figure 21: RNN model results

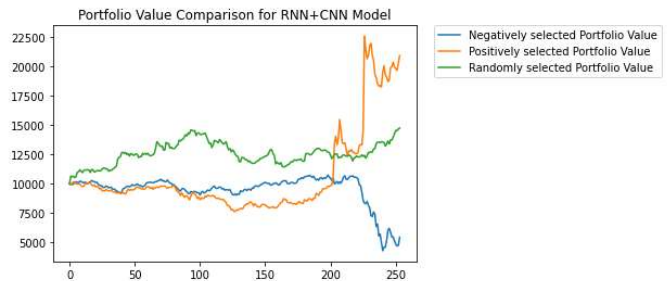


Figure 22: RNN+CNN model results

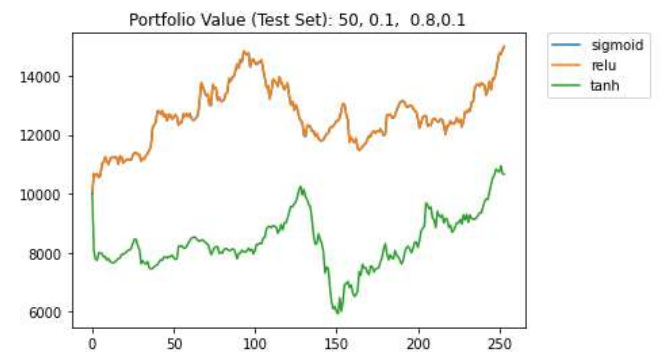


Figure 23: RNN results with activation functions

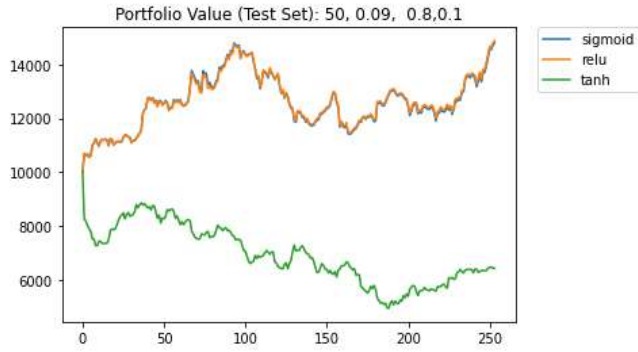


Figure 24: RNN+CNN results with activation functions

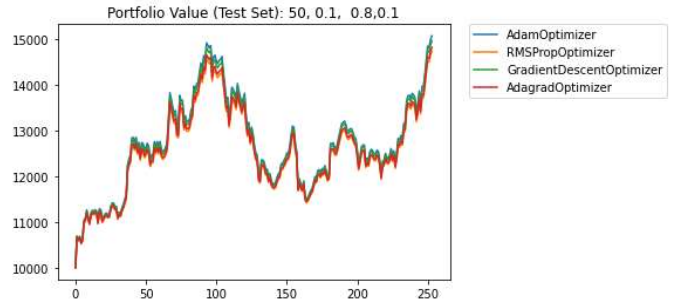


Figure 28: CNN results with different optimizers

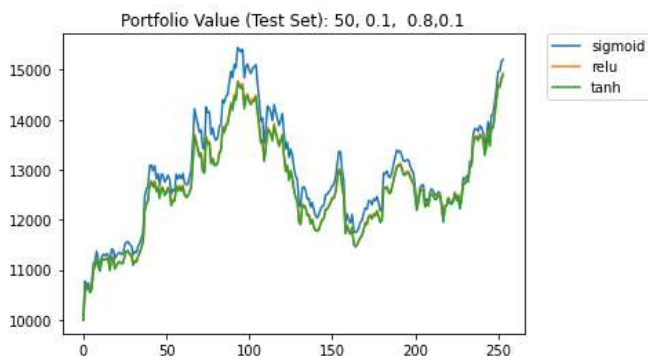


Figure 25: CNN results with activation functions

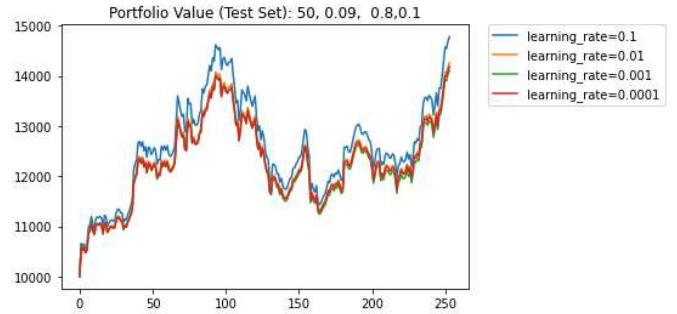


Figure 29: RNN results with different learning rates

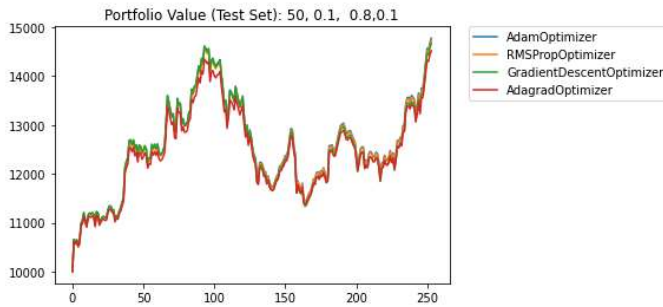


Figure 26: RNN results with different optimizers

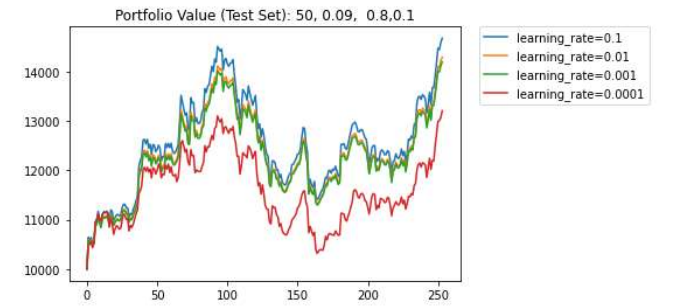


Figure 30: RNN+CNN results with different learning rates

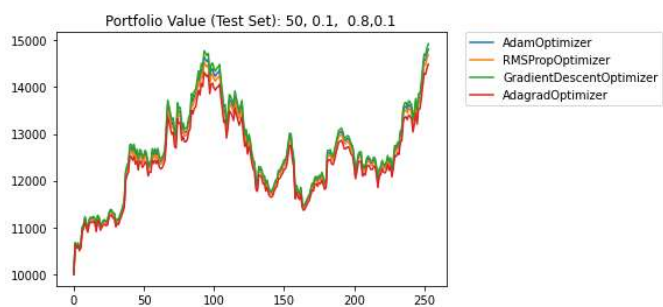


Figure 27: RNN+CNN results with different optimizers

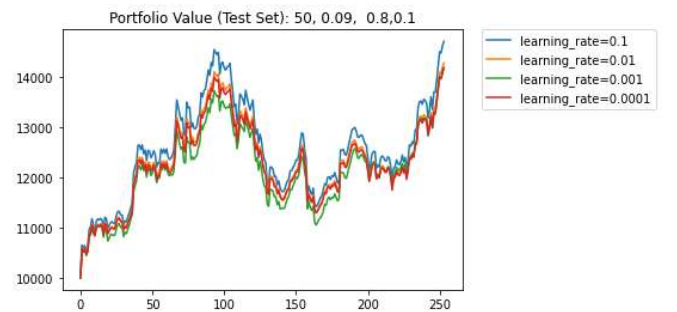


Figure 31: CNN results with different learning rates

6 Conclusion

This paper implemented deep reinforcement learning method to allocate weights for stocks and cryptocurrency in portfolio management and compared different policy network models. The training results from various policy networks constructed with convolutional neural network, recurrent neural network, or the mixed one show similar results. In this research, three different portfolios were considered, and they include stocks that have been negatively influenced by the Covid-19, and the stocks that have been positively influenced by Covid-19, and the stocks and cryptocurrency that were randomly selected. Among the three portfolios, it was found that all the agents have the best performance on the randomly selected portfolio, which presents overall stable moving-up portfolio values as time goes on. The average return rate for all three deep reinforcement learning agents reached 47.38%. It is known that the financial market is filled with uncertainties and unpredictability due to many different factors, so the training results from this research also show the instability when dealing with positively influenced portfolio and negatively influenced portfolio. The linear regression model for portfolio optimization shows poor results, but the weight allocation changed more when compared with deep reinforcement learning agents. Overall, the linear regression model does not meet the goal of financial portfolio management. On the other hand, the deep reinforcement learning agents show a better performance when cryptocurrency was added to the portfolio.

This paper considers that cryptocurrency does not have a long history, so only 5 years data were used. For future studies, data from 10 or more years can be collected to research the performance of deep reinforcement learning agents. The paper only applied to policy gradient in reinforcement learning; in the future, deep Q network and Deep Deterministic Policy Gradient can also be added to do a comparison study with the agents in the paper. For financial fields, there are many financial portfolio optimization methods, but in this research only equal weighted allocation strategy was used to compare with deep reinforcement learning agents. To do a further study, other portfolio optimization methods can also be applied and tested.

References

- [1] A. Jonsson, "Deep Reinforcement Learning in Medicine," *Kidney Diseases*, 5:18-22, 2019.
- [2] A. Tianello, *Using a Multi-Factor Regression Model to Determine Asset Allocation*, <https://ufdcimages.uflib.ufl.edu/AA/00/05/75/50/00001/altianello-Thesis.pdf>.
- [3] C. Adjasi and N. Biekpe, "Stock Market Development and Economic Growth: The Case of Selected African Countries," *African Development Review*, 18:144-161, 2006.
- [4] G. Ban, N. El Karoui, and A. Lim, "Machine Learning and Portfolio Optimization," *Management Science*, 64:1136-1154, 2018.
- [5] G. Lample, and D. Chaplot, "Playing FPS Games with Deep Reinforcement Learning," *Thirty-First AAAI Conference On Artificial Intelligence, San Francisco, California, USA*, 2017.
- [6] J. Pineiro-Chousa, M. Vizcaino-Gonzalez, and A. Perez-Pico, "Influence of Social Media over the Stock Market," *Psychology and Marketing*, 34: 101-108, 2017.
- [7] L. Kaiser, M. Babaeizadeh, P. Milos, B. Osinski, R. Campbell, K. Czechowski, D. Erhan, C. Finn, P. Kozakowski, and S. Levine, "Model-Based Reinforcement Learning for Atari," *ArXiv Preprint ArXiv:1903.00374*, 2019.
- [8] L. Schultz, and V. Sokolov, "Deep Reinforcement Learning for Dynamic Urban Transportation Problems," *ArXiv Preprint ArXiv:1806.05310*, 2018.
- [9] L. Wang, and A. Kutan, "The Impact of Natural Disasters on Stock Markets: Evidence from Japan and the US," *Comparative Economic Studies*, 55:672-686, 2013.
- [10] N. Kanwar, *Deep Reinforcement Learning-based Portfolio Management*, <https://rc.library.uta.edu/uta-ir/bitstream/handle/10106/28108/KANWAR-THESIS-2019.pdf?sequence=1&isAllowed=y>, 2019.
- [11] N. Kumar, S. Rahman, and N. Dhakad, "Fuzzy Inference Enabled Deep Reinforcement Learning-Based Traffic Light Control for Intelligent Transportation System," *IEEE Transactions On Intelligent Transportation Systems*, 22:4919-4928, 2020.
- [12] P. Henderson, R. Islam, P. Bachman, J. Pineau, D. Precup, and D. Meger, "Deep Reinforcement Learning that Matters," *Proceedings Of The AAAI Conference On Artificial Intelligence, Louisiana, USA*, 392:3207-3214, 2018.
- [13] P. Santa-Clara, and R. Valkanov, "The Presidential Puzzle: Political Cycles and the Stock Market," *The Journal Of Finance*, 58:1841-1872, 2003.
- [14] P. Yu, J. Lee, I. Kulyatin, Z. Shi, and S. Dasgupta, "Model-Based Deep Reinforcement Learning for Dynamic Portfolio Optimization," *ArXiv Preprint ArXiv:1901.08740*, 2019.
- [15] S. Amrouni, A. Moulin, and P. Mizrahi, *A Deep Reinforcement Learning Framework for the Financial Portfolio Management Problem*, <https://selimamrouni.github.io/portfolio/air-elle.html>, 2018.
- [16] S. Bartram, G. Brown, and R. Stulz, "Why are US Stocks more Volatile?" *The Journal Of Finance*, 67:1329-1370, 2012.
- [17] V. Mnih, K. Kavukcuoglu, D. Silver, A. Graves, I. Antonoglou, D. Wierstra, and M. Riedmiller, "Playing Atari with Deep Reinforcement Learning," *ArXiv Preprint ArXiv:1312.5602*, 2013.
- [18] X. Liu, Z. Ding, S. Borst, and A. Walid, "Deep Reinforcement Learning for Intelligent Transportation Systems," *ArXiv Preprint ArXiv:1812.00979*, 2018.

- [19] Yahoo Finance. <https://finance.yahoo.com/>.
- [20] Y. Liu, B. Logan, N. Liu, Z. Xu, J. Tang, and Y. Wang, "Deep Reinforcement Learning for Dynamic Treatment Regimes on Medical Registry Data," *2017 IEEE International Conference On Healthcare Informatics (ICHI), Park City, Utah, USA*, pp. 380-385, 2017.
- [21] Y. Zhang, P. Zhao, Q. Wu, B. Li, J. Huang, and M. Tan, "Cost-Sensitive Portfolio Selection via Deep Reinforcement Learning," *IEEE Transactions On Knowledge And Data Engineering*, 34: 236-248, 2020.
- [22] Z. Liang, H. Chen, J. Zhu, K. Jiang, and Y. Li, "Adversarial Deep Reinforcement Learning in Portfolio Management," *ArXiv Preprint ArXiv:1808.09940*, 2018.
- [23] Z. Liang, K. Jiang, H. Chen, J. Zhu, and Y. Li, "Deep Reinforcement Learning in Portfolio Management," *ArXiv Preprint ArXiv:1808.09940*, 2018.
- [24] Z. Liu, C. Yao, H. Yu, and T. Wu, "Deep Reinforcement Learning with its Application for Lung Cancer Detection in Medical Internet of Things," *Future Generation Computer Systems*. 97: 1-9, 2019.
- [25] Z. Jiang, D. Xu, and J. Liang, "A Deep Reinforcement Learning Framework for the Financial Portfolio Management Problem," *ArXiv Preprint ArXiv:1706.10059*, 2017.
- [26] Z. Jiang and J. Liang, "Cryptocurrency Portfolio Management with Deep Reinforcement Learning," *2017 Intelligent Systems Conference (IntelliSys), London, United Kingdom*, pp. 905-913, 2017.



Yue Ma is a PhD student at Northern Illinois University, where she is pursuing her passion for computer science. She received her Master's degree from Southeast Missouri State University, bringing a wealth of knowledge and experience to her studies. Her research interests include data visualization, augmented reality (AR), and machine learning. She is dedicated to furthering her knowledge,

and actively working towards making more publications and contributing to her research fields in meaningful ways.



Ziping Liu received her PhD in Engineering Science from Southern Illinois University at Carbondale in 1999 and began her computing career at Motorola, where she developed software for mobile phones. Currently, she is a professor of Computer Science at Southeast Missouri State University, where she has been teaching since 2001. Her research interests encompass a wide range of

topics, including machine learning, cloud computing, secured software design, wireless ad hoc networks and sensor networks, distributed computing and game development.



Charles D. McAllister is currently a professor of business analytics at Southeast Missouri State University. He has previously served the university in such roles as vice provost and dean of graduate studies,

interim dean of business, interim chairperson of computer science, and interim director of institutional research. He received a Ph.D. in industrial engineering and operations research from The Pennsylvania State University in 2002, M.S. in industrial engineering from Purdue University in 1999, and B.S. in industrial and management systems engineering from the University of Nebraska in 1998. His research interests include optimization, analytics, simulation, pedagogy, and higher education administration.