

Deep reinforcement

by Aditya Major

Submission date: 07-Dec-2018 03:39PM (UTC+0530)

Submission ID: 1052583409

File name: Report.pdf (2.36M)

Word count: 7758

Character count: 38486

CHAPTER 1

OUTLINE OF THE REPORT

This report is constructed in a way to explain how the financial domain works and the kind of algorithms we have applied on the data procured from these domains. This flow of contents is made to help the reader understand the work and the pre-requisites required to understand our research and implementation.

Chapter 1 indicates the outline of the report as well as an overview about the work he have done and the research we have put forward.

Chapter 2 lies as an introduction towards our project as well as providing some pre-requisite knowledge. Chapter 2.1 shows our motivation towards taking this project as well as the scope of implementing reinforcement learning on various domains, along with a brief literature review. Chapter 2.2 gives an overview of the financial trade domain, where Chapter 2.2.1 represents Price Chart, Chapter 2.2.2 represents Trade History and Chapter 2.2.3 represents Order Book.

Chapter 3 familiarizes the basics of reinforcement learning to the reader which is a pre-requisite required to understand our model and its implementation which is further introduced in Chapter 3.1. We understand how the model learns to optimize rewards in Chapter 3.1.1. Policy Search is introduced in Chapter 3.2 and the Policy Gradient algorithm is initiated in Chapter 3.2.1. Markov Chains present in Chapter 3.3.1 is present under the subset of Markov Decision Processes present in Chapter 3.3, which is further explained in Chapter 3.3.2. Temporal Difference Learning is introduced in Chapter 3.4 with Q-Learning Algorithm and SARSA Algorithm being a subset of it in

²⁹

Chapter 3.4.1 and Chapter 3.4.2 respectively. Chapter 3.4.3 introduces the concept of Exploration Policies.

Chapter 4 introduces the idea of Deep Reinforcement Learning. We present a brief introduction of this topic in Chapter 4.1 and then move on to introduce the Deep Q-Network algorithm in Chapter 4.2. Chapter 4.2.1 explores the problems and shortcomings of Deep Q-Networks and the solutions to these problems are mentioned in Chapter 4.2.2. We briefly give an overview of the implementation details of Deep Reinforcement Learning algorithms in Chapter 4.2.3. Double DQNs are explored in the next Chapter 4.3. The concept of Deep Recurrent Q-Networks is introduced in Chapter 4.4. We give an overview of Recurrent Neural Networks (RNN) and Long Short Term Memory (LSTM) in Chapter 4.4.1 and Chapter 4.4.2 respectively. Finally, we explore the usefulness of DRQN in Chapter 4.4.3.

CHAPTER 2

INTRODUCTION

2.1 MOTIVATION AND LITERATURE REVIEW

The motivation for considering this project comes from the fact that the Deep Learning researchers and community has stayed away from modelling on Financial Markets. Some of the reasons causing this problem seems to be disinterest towards financial markets from research perspective and the difficult in obtaining data. We believe that this subject has not received enough attention from the researchers despite the potential to create state-of-the art models.

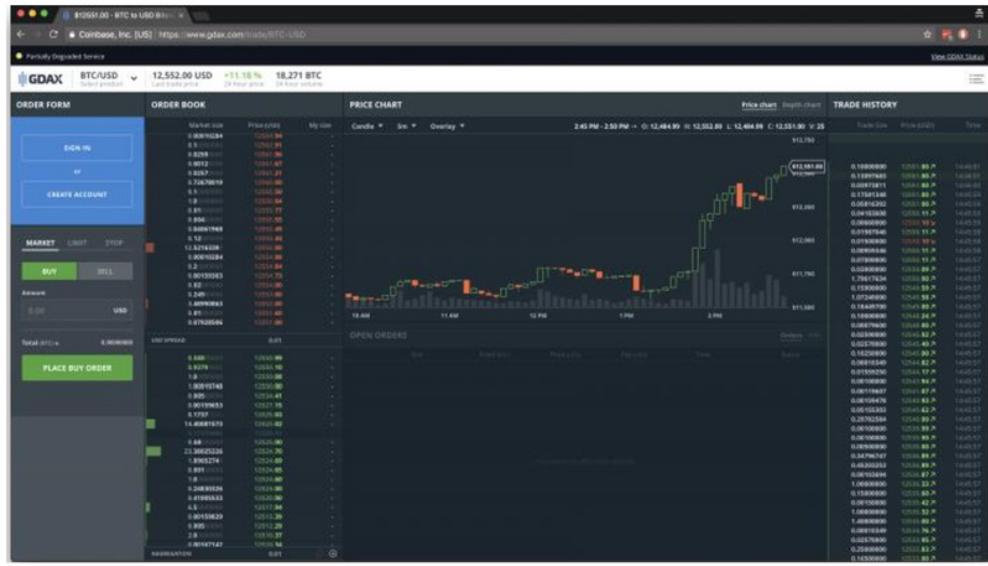
Another fact that leads to us working on this project is the algorithm we are going to implement on the financial market. Among the three dominant machine learning techniques, reinforcement learning is the most sparsely used within various domains, other techniques being supervised learning and unsupervised learning. Taking this project helps us to learn more about the field of reinforcement learning and its applications towards financial markets as well as other domains.

The scope of reinforcement learning in financial domain is vast. Since this is the least applied out of all the major machine learning techniques, a lot of research work can be done to bring better results on domains primarily working on supervised and unsupervised learning.

Work done by Huang, Chien-Yi [1] is one of the only papers which have a baseline model running in the trading domain for reinforcement learning algorithms. Chiao-Ting; An-Pin; Szu-Hao also implemented a agent-based reinforcement learning algorithm [2] to produce cloning strategies from trading records. Yue Deng et al formed a Deep Direct Reinforcement Learning Algorithm [3] for representation of financial signals and its trading. Zhengyao Jiang; Jinjun Liang worked on Cryptocurrency data for Portfolio Management [4].

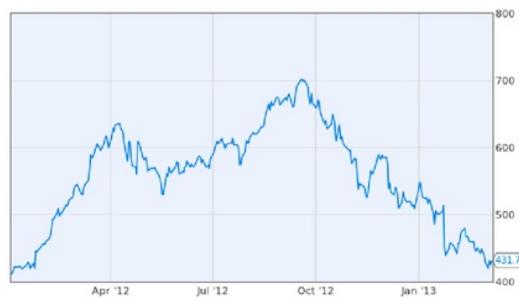
2.2. FINANCIAL MARKETS

Financial trading as well as cryptocurrency trading in most markets happen in an auction called continuous double auction with an open order book on an exchange. This also means that there are buyers and sellers present who get matched and can trade assets amongst themselves. There are various exchanges present in the financial markets which may have different products like U.S Dollar vs Euros or Bitcoin vs U.S Dollars. Apart from that, their interface and their data is very similar.



2.2.1 Price Chart

The current price represents the price of the most recent trade. It depends on the kind of action taken by the trade like buy or sell. The price chart is usually shown as a candlestick chart which for a given time window shows 4 different types of prices; Open (O), High (H), Low (L) and Close(C) prices. Volume (V), which is the total volume of the trades happened within the given time period, is represented by the bars below the given price chart. The volume of a trade represents the market liquidity. It also indicates the quality and veracity of the trend prices, as high volume is usually the consensus of trade participants.



2.2.2 Trade History

The history of all the trades can be viewed on the right hand side. Each trade is represented by a direction (sell or buy), timestamp, price and size. The two parties in the trade are referred as maker and taker.

2.2.3 Order Book

The order book contains information about the selling and buying prices which people are willing to pay. It contains two sides, Bids and Asks. Bids are people who want to buy and Asks are people who want to sell. Best bid is the highest price one is willing to buy and the best sell is the lowest price one is willing to sell. Generally, the best ask is larger than the best bid. Spread is the difference between the best ask and the best bid.

Learn more'. There is an 'Advanced' link with a right-pointing arrow. At the bottom, it shows 'Total (USD) ≈ 10010.00' and a large green 'PLACE BUY ORDER' button."/>

MARKET LIMIT STOP

BUY SELL

Amount
1 BTC

Limit Price
10010 USD

Execution
Post Only Allow Taker

Your order will only execute as a maker order. [Learn more](#)

Advanced

Total (USD) ≈ 10010.00

PLACE BUY ORDER

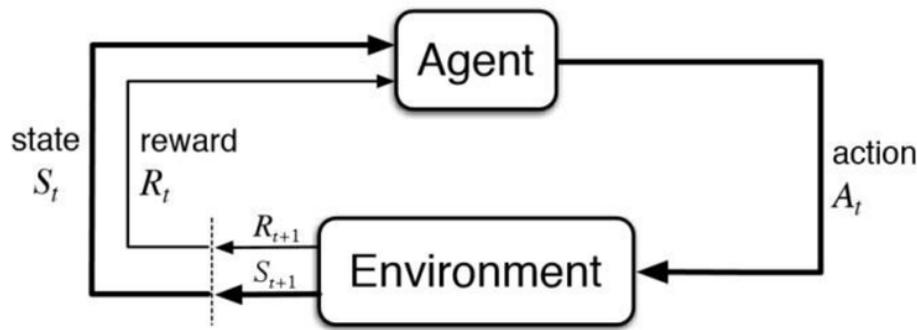
CHAPTER 3

BASICS OF REINFORCEMENT LEARNING

3.1. INTRODUCTION

Reinforcement Learning is one of the oldest as well as one of the most exciting field of machine learning at this moment. It started around 1950s, with various applications over the years, typically related to games like TD-Gammon which was created to solve Backgammon as well as machine control, but it was never highly utilised. In 2013, a start-up called DeepMind made a system which could learn to play any Atari games, revolutionised the field of reinforcement learning. It only required raw pixels as inputs and could beat most humans without any prior knowledge about the rules of the games.

This was followed by a series of remarkable models, including one of the most famous in the domain of reinforcement learning, the victory of AlphaGo over Lee Sedol, who at that time was the world champion of the game Go. Before that, no program had ever come close to beating the champion of a game at his own speciality. In 2014, Google bought DeepMind for 500 million dollars signifying the success of this new wave of Reinforcement Learning.



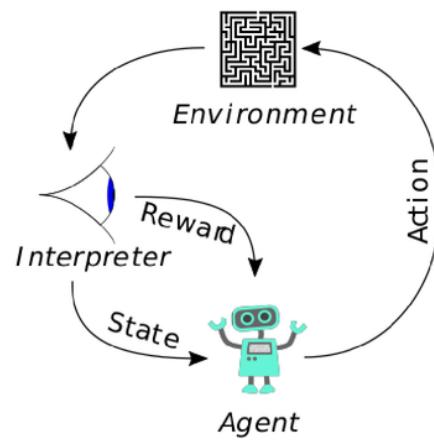
3.1.1 Learning to optimise rewards

In reinforcement learning, actions are taken by a software agent which makes observations within an environment, which returns rewards back to the agent. The main objective is to learn actions that will increase and hence maximize the expected long-term rewards. Therefore, to summarize, the agent; through trial and error, learns to maximize rewards in an environment. Reinforcement learning can be applied to a wide domain of work. Some of the examples are:-

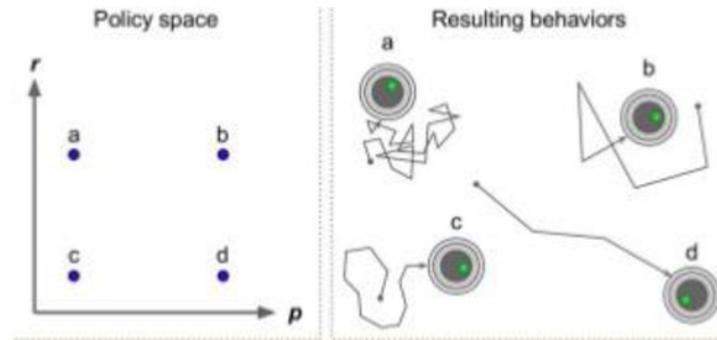
1. The stock market prices are observed by the agent, which decides when to sell or buy, like the implementation by Jae Won Lee.^[5]
2. A board game agent can be programmed like Go.
3. Ms Pac Man controlling agent can be programmed. Here, the simulation of the game is the environment of the game, screenshots are the observations, the joystick positions are the rewards and the game points are rewards.
4. The agent can be the walking robot controller. Here, the agent can observe the environment using sensors like touch sensors and cameras, the signals sent in order to activate monitors will be its actions. A positive reward would occur when the robot approaches the target destination and a negative rewards would occur in the opposite case.

3.2 POLICY SEARCH

The algorithm which is used by the agent to determine actions is known as Policy. For example, observations could act as inputs to a policy which outputs the actions which the agent has to take. It is not necessary that the policy has to be deterministic. If we consider a robotic vacuum cleaner where the reward would be the amount of dirt it cleans in 15 minutes, its policy could be to rotate right or left with probability ‘p’ or move forward with a probability ‘p-1’. The rotation would lie between -r and +r degrees, due to which it is called a stochastic policy.



In the above mentioned case, there are two ways to train the robot; by tweaking either the angle range r or the probability p . One way of exploring the policy space, is to use the brute force policy search method in which the algorithm tries as many different values of the parameters as possible and pick the most rewarding combination. This is not a good method for parameter exploration since its complexity would be very high. Genetic Algorithms can be another way to explore the policy space where we generate a population of policies, then kill the worst policies as well as crossover the best policies to create offsprings. This way we can iterate through various generations in order to get the best policy.



Using optimization techniques is another approach to do policy search. This is done by evaluating the gradients of the rewards with regards to the policy parameters and changing these parameters by following the gradient towards higher rewards (gradients ascent). This technique is known as Policy Gradients (PG).

3.2.1 Policy Gradients

As mentioned earlier, Policy Gradient optimizes the policy parameters by following the gradients towards higher rewards. REINFORCE ²¹ algorithm is one of the most popular policy gradient algorithm. There are several steps we have to follow in order to optimize the parameters.

1. We will let the policy play the game several times and the compute the gradient which would make the chosen action even more likely without actually applying the gradients.
2. After running several episodes, compute the score of each action.
3. If the score of the action is negative, then it represents a bad action and we should apply gradients in the opposite direction to make sure that this action occurs less likely in the future. On the other hand, if the score of the action is positive, then it represents a good action and we should apply gradients in the same direction as the rewards to make the action even more likely to be selected in the future.

- Finally, we will compute the mean of all the resulting gradient vectors and use these vectors to perform a Gradient Descent iteration.

Algorithm

```

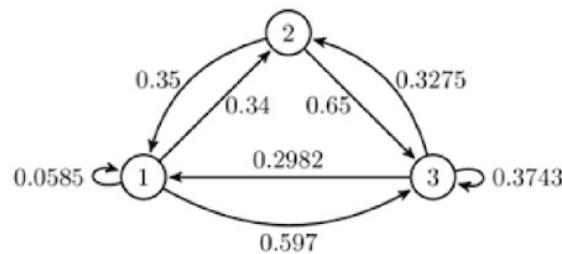
function REINFORCE
    Initialise  $\theta$  arbitrarily
    for each episode  $\{s_1, a_1, r_2, \dots, s_{T-1}, a_{T-1}, r_T\} \sim \pi_\theta$  do
        for  $t = 1$  to  $T - 1$  do
             $\theta \leftarrow \theta + \alpha \nabla_\theta \log \pi_\theta(s_t, a_t) v_t$ 
        end for
    end for
    return  $\theta$ 
end function

```

3.3 MARKOV DECISION PROCESS

3.3.1 Markov Chains

In the early 20th century, Markov chains were discovered by Andrey Markov, who studied stochastic processes having no memory. This process has a fixed number of states where it randomly evolves one state to another state at each step. This transition from a state s to s' has a fixed probability and it only depends on the pair(s, s') not on the past states, i.e. memoryless system.

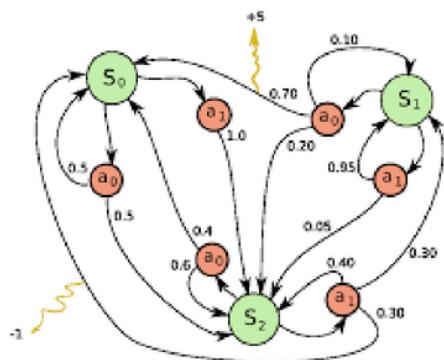


0.0585	0.34	0.597
0.35	0	0.65
0.2982	0.3275	0.3743

The first figure represents the Markov chain as a computational graph having 3 states with probabilities of transforming into different states. The table shows the transition probability matrix.

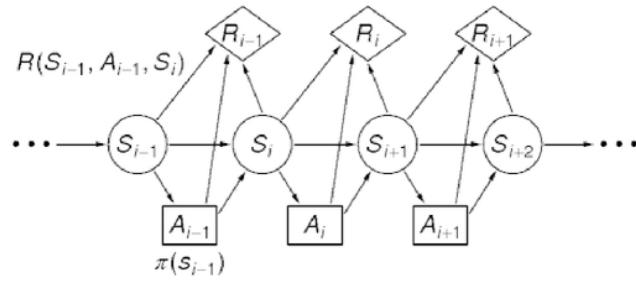
3.3.2 Markov Decision Processes (MDPs)

Markov Decision Process were first described by Richard Bellman in the 1950s. They were very similar to Markov chains. The agent could one of the possible actions at a step where the transition probabilities would also depend on the actions. The only difference would be that the environment returns positive or negative rewards for each action and the agent's goal would be to maximize these rewards over a long period of time.



In the above figure, there are 3 states (S_0, S_1, S_2) and 2 actions (a_0, a_1). If the agent starts at S_0 then it can choose among the two actions. If it chooses a_0 then it has equal probability to stay at S_0 or to move towards S_2 otherwise it will always move to S_2 by choosing the action a_1 . At state S_1 , if it choose action a_0 , then there is a 70% chance to move into state S_0 by getting a reward of +5.

Similarly in state S_2 , if it selects action a_1 , there is a 30% chance to move towards S_0 and getting a reward of -1 from the environment.



$V^*(s)$ is known as the optimal state value of any state s which was estimated by Bellman. He formed the Bellman Optimality Equation where $V^*(s)$ is the sum of all discounted future rewards the agent can expect on average after it reaches a state s if it acts optimally.

$$V^*(s) = \max_a \left\{ \sum_{s'}^{14} T(s, a, s') [R(s, a, s') + \gamma \cdot V^*(s')] \right\} \quad \forall s$$

- $T(s,a,s')$ is referred as the transition probability from state s to s' when the chosen action is a .
25
- $R(s,a,s')$ is the reward which the agent gets due to choosing action a which causes a transition from state s to s' .
22
- γ is the discount rate.

3
This recursive equation states that the optimal value of the present state is equal to the reward it will get on average after taking one optimal action, plus the expected optimal value of all possible next states that this action can lead to, provided that the agent acts optimally.

This equation directly leads to an algorithm known as Value Iteration algorithm which can precisely estimate the optimal state value of each and every state. We first initialize all the state value estimates to zero and then update them using the above mentioned algorithm iteratively. These estimates always converge to the optimal state values depending on the optimal policy, given enough time.

Algorithm

Output: π^* , the optimal policy
12

1. Initialize $v(s)$ arbitrarily for all $s \in S$
2. **repeat**
3. $\Delta \leftarrow 0$
4. **for** each $s \in S$ **do**
5. $v_{old} \leftarrow v(s)$
14
6. $v(s) = \max_a \{ \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma \cdot v(s')] \}$
7. $\Delta \leftarrow \max\{\Delta, |v_{old} - v(s)|\}$
8. **end for**
9. **until** $\Delta < \epsilon$
10. **for** each $s \in S$ **do**
14
11. $\pi(s) \leftarrow \operatorname{argmax}_{a \in A} \{ \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma \cdot v(s')] \}$

```
12. end for  
13. return  $\pi$ 
```

It is very useful to know the optimal state values especially in the case of policy evaluation, but it doesn't explicitly tell the agent what to do. To counter this problem, Bellman also found another similar algorithm to estimate the optimum state-action values, called the Q-values. $Q^*(s,a)$ is known as the optimal Q-Value of the state-action pair which is the sum of discounted future rewards the agent can expect on average once it chooses the action a after reaching the state s . A difference is that the agent chooses the action without actually seeing its outcome given that it acts optimally after the action.

Algorithm

```
Initialize  $V(s)$  to arbitrary values  
Repeat  
  For all  $s \in S$   
    For all  $a \in \mathcal{A}$   
       $Q(s,a) \leftarrow E[r|s,a] + \gamma \sum_{s' \in S} P(s'|s,a)V(s')$   
       $V(s) \leftarrow \max_a Q(s,a)$   
Until  $V(s)$  converge
```

3.4 TEMPORAL DIFFERENCE LEARNING

As seen before, Reinforcement Learning problems with discretized actions can be modelled as Markov Decision processes. There is one problem with this process, the transition probabilities ($T(s,a,s')$) are unknown in the beginning of the process. Also, the rewards ($R(s,a,s')$) are unknown as well. Therefore, the model should experience each transition and state to know the rewards at least one time before the iterations. In order

to estimate the transition probabilities well, it must experience these states and transitions multiple times.

Just like the Value Iteration algorithm, the Temporal Difference algorithm is also very similar except it takes care of the problems we mentioned above. Usually, we assume that the software agent only knows about the possible actions and states initially. Exploration policy is then used by the agent to explore the Markov Decision Process. Then, based on the rewards and transitions that are actually observed, the Temporal Difference Learning algorithm updates the estimates of the state values.

Algorithm

```
Input: the policy  $\pi$  to be evaluated
Initialize  $V(s)$  arbitrarily (e.g.,  $V(s) = 0, \forall s \in S^+$ )
Repeat (for each episode):
    Initialize  $S$ 
    Repeat (for each step of episode):
         $A \leftarrow$  action given by  $\pi$  for  $S$ 
        Take action  $A$ ; observe reward,  $R$ , and next state,  $S'$ 
         $V(S) \leftarrow V(S) + \alpha[R + \gamma V(S') - V(S)]$ 
         $S \leftarrow S'$ 
    until  $S$  is terminal
```

From the algorithm we can see that it keeps the running average of the immediate rewards which the agent receives when it leaves the state, plus the rewards it expects to get later.

3.4.1 Q-LEARNING

Similar to TD Learning algorithm, the Q-learning algorithm is a likewise adaptation of the Q-value iteration algorithm. Here, the rewards as well as the transition probabilities are initially not given.

Q-learning: Learn function $Q : \mathcal{X} \times \mathcal{A} \rightarrow \mathbb{R}$

Require:

Sates $\mathcal{X} = \{1, \dots, n_x\}$
Actions $\mathcal{A} = \{1, \dots, n_a\}, \quad A : \mathcal{X} \Rightarrow \mathcal{A}$
Reward function $R : \mathcal{X} \times \mathcal{A} \rightarrow \mathbb{R}$
Black-box (probabilistic) transition function $T : \mathcal{X} \times \mathcal{A} \rightarrow \mathcal{X}$
Learning rate $\alpha \in [0, 1]$, typically $\alpha = 0.1$
Discounting factor $\gamma \in [0, 1]$
procedure QLEARNING($\mathcal{X}, A, R, T, \alpha, \gamma$)
 Initialize $Q : \mathcal{X} \times \mathcal{A} \rightarrow \mathbb{R}$ arbitrarily
 while Q is not converged **do**
 Start in state $s \in \mathcal{X}$
 while s is not terminal **do**
 Calculate π according to Q and exploration strategy (e.g. $\pi(x) \leftarrow \arg \max_a Q(x, a)$)
 $a \leftarrow \pi(s)$
 $r \leftarrow R(s, a)$ ▷ Receive the reward
 $s' \leftarrow T(s, a)$ ▷ Receive the new state
 $Q(s', a) \leftarrow (1 - \alpha) \cdot Q(s, a) + \alpha \cdot (r + \gamma \cdot \max_{a'} Q(s', a'))$
 return $\overset{s \leftarrow s'}{Q}$

This algorithm keeps track of a running average of the rewards r for each state-action pair (s, a) which the agent gets by deploying action a and leaving the state s , plus the rewards it expects to get later. Maximum of the Q-Value estimates are taken for the next state, given that the target policy would act optimally. This algorithm will always converge to the optimal Q-values if ran for a long time. In this case, the policy being executed and the policy being trained are not the same. This is known as off-policy algorithm.

8 3.4.2 STATE-ACTION-REWARD-STATE-ACTION (SARSA)

SARSA is very similar to Q-learning. The main difference between Q-Learning and SARSA is that Q-learning is an off-policy algorithm whereas SARSA is an on-policy algorithm. It means that SARSA learns the Q-value based on the action determined by the current policy rather than the greedy policy.

SARSA(λ): Learn function $Q : \mathcal{X} \times \mathcal{A} \rightarrow \mathbb{R}$

Require:

Sates $\mathcal{X} = \{1, \dots, n_x\}$

Actions $\mathcal{A} = \{1, \dots, n_a\}$, $\mathcal{A} : \mathcal{X} \Rightarrow \mathcal{A}$

Reward function $R : \mathcal{X} \times \mathcal{A} \rightarrow \mathbb{R}$

Black-box (probabilistic) transition function $T : \mathcal{X} \times \mathcal{A} \rightarrow \mathcal{X}$

Learning rate $\alpha \in [0, 1]$, typically $\alpha = 0.1$

Discounting factor $\gamma \in [0, 1]$

$\lambda \in [0, 1]$: Trade-off between TD and MC

procedure QLEARNING($\mathcal{X}, \mathcal{A}, R, T, \alpha, \gamma, \lambda$)

 Initialize $Q : \mathcal{X} \times \mathcal{A} \rightarrow \mathbb{R}$ arbitrarily

 Initialize $e : \mathcal{X} \times \mathcal{A} \rightarrow \mathbb{R}$ with 0. ▷ eligibility trace

while Q is not converged **do**

 Select $(s, a) \in \mathcal{X} \times \mathcal{A}$ arbitrarily

while s is not terminal **do**

$r \leftarrow R(s, a)$

$s' \leftarrow T(s, a)$ ▷ Receive the new state

 Calculate π based on Q (e.g. epsilon-greedy)

$a' \leftarrow \pi(s')$

$e(s, a) \leftarrow e(s, a) + 1$

$\delta \leftarrow r + \gamma \cdot Q(s', a') - Q(s, a)$

for $(\tilde{s}, \tilde{a}) \in \mathcal{X} \times \mathcal{A}$ **do**

$Q(\tilde{s}, \tilde{a}) \leftarrow Q(\tilde{s}, \tilde{a}) + \alpha \cdot \delta \cdot e(\tilde{s}, \tilde{a})$

$e(\tilde{s}, \tilde{a}) \leftarrow \gamma \cdot \lambda \cdot e(\tilde{s}, \tilde{a})$

$s \leftarrow s'$

return $\overset{a}{\overleftarrow{Q}}$

From the above algorithm, we can see that the two action selections always follow the current policy unlike the Q-learning algorithm in which the algorithm has no constraint over the next action, as long as the next state Q-value is maximized.

3.4.3 EXPLORATION POLICIES

Exploration policies are very important since the Q-learning models depends a lot on the exploration of the Markov Decision Process. We can use a random policy as our exploration policy since it is guaranteed to visit all the transitions and states many times. But, there is a serious problem with this policy and that is its complexity. Therefore, we use a different option which is known as the ϵ -greedy policy in which the policy acts randomly with a probability ϵ or chooses the action with the highest Q-value (greedily) with the probability $1 - \epsilon$. The main advantage of ϵ -greedy policy is the fact that it spends more time exploring the useful parts of the MDP while still exploring some unknown parts of the environment as the Q-estimates improve.

Another approach which can be used as exploration policy is to increase the chances of trying actions it hasn't tried much instead of relying on chance for exploration. This can be implemented as a small tweak to the Q-Value estimates.

$$Q(s', a') \xleftarrow{23} (1 - \alpha) \cdot Q(s, a) + \alpha \cdot (r + \gamma \cdot \max_a \{f(Q(s', a'), N(s', a))\})$$

- $N(s', a')$ counts the number of times state s' chose action a' .
- $f(q, n)$ is an exploration function.

CHAPTER 4

DEEP REINFORCEMENT LEARNING

4.1 INTRODUCTION

Q-Learning is an extremely powerful algorithm that provides a robust solution to Reinforcement Learning problems modeled as MDPs. However, there is a major drawback to this approach. In Q Learning, we maintain a Q-value table that maps <state, action> pairs to real valued probable rewards.

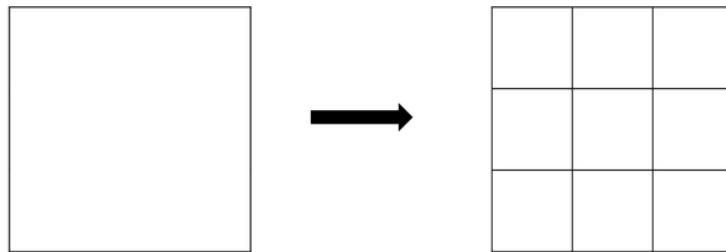
$$Q: S \times A \rightarrow R$$

These values are stores in a Q-table and iteratively tuned during the training phase. This very much resembles dynamic programming. The space complexity of the algorithm increases in direct proportion to the state space of the environment. Hence, for situations where the state space or the action space is very large, Q-Learning becomes quite difficult to implement in practice.

Also, since the agent makes decision after referring to the Q-table, for states which are not present in the Q-table, the agent will have no clue what to do. The agent will make no assumptions and no approximations. Scenarios where the state space is not discrete are not uncommon in practical world setup. For example, most robotic applications of

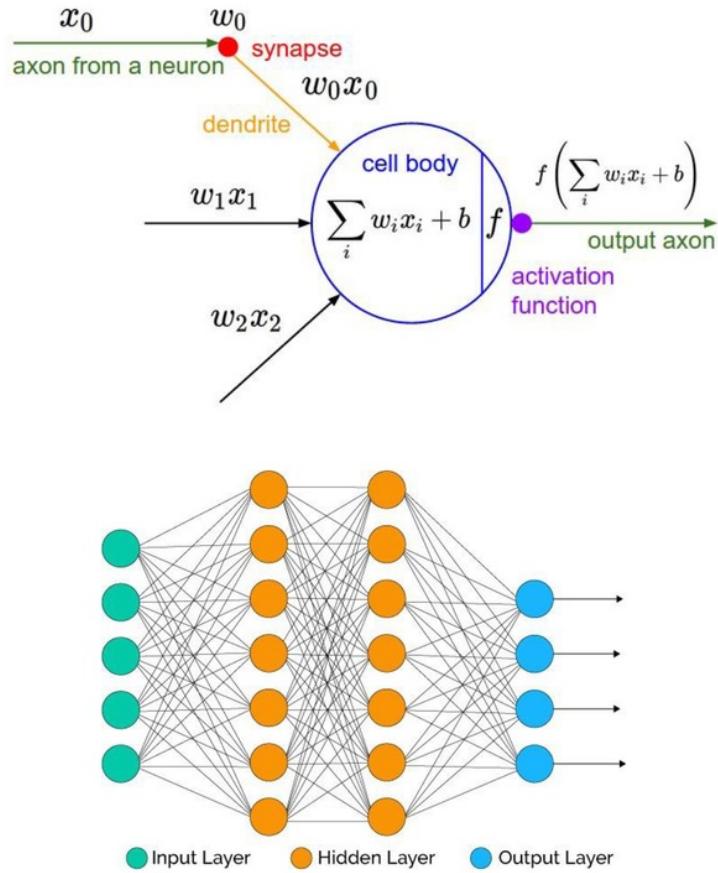
reinforcement learning possess state spaces that are continuous, defined by the means of continuous variables such as position, angular displacement, velocity, torque, etc.

The usual approach has been to discretize the continuous variables, which more often than not, leads to combinatorial explosion and the well-known “Curse of Dimensionality”.



The above figure shows the discretization of a continuous state space. It has been observed that Q-learning algorithm where the states and actions have been discretised, scale poorly. As we increase the number of state and action variables, the size of the table that is used to store Q-values explodes exponentially. For accurate control, the variables have to be quantized very fine, but as these systems fail to generalise between similar states and actions as they require very large quantities of training data. We can avoid such problems that arise due to the discretization of state space by using methods that can directly deal with continuous states.

4.2 DEEP REINFORCEMENT LEARNING

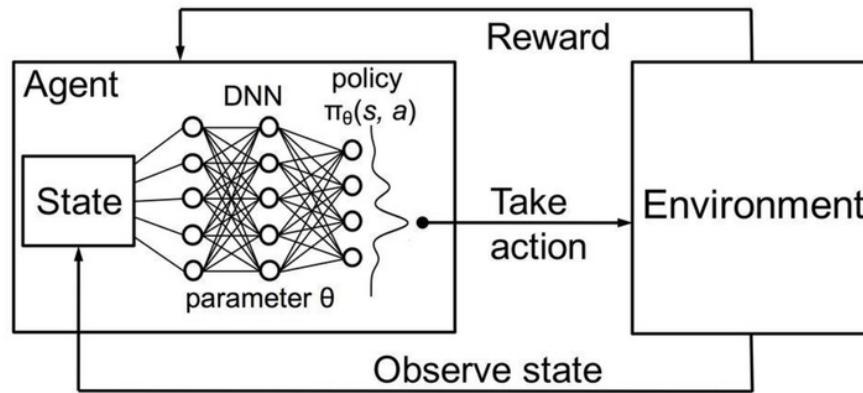


One of the most popular and effective approach to overcome this shortcoming has been to use neural networks to approximate the Q-values of <state, action> pairs for continuous state space. This is strongly supported by the [Universal Approximation Theorem](#). In the mathematical theory of artificial neural networks, the theorem states that under certain soft assumptions on the activation function, a neural network

16

containing a finite number of neurons can approximate continuous functions on compact subsets of R^n .

A neural network architecture is designed to estimate the Q-values in case of continuous state space. The input to the network is the current state of the agent and the output is an n-dimensional vector where each dimension represents the Q-value of that particular action; given there are 'n' possible actions. The algorithm is called Deep Q-Network (DQN).^[6]



Algorithm

Start with $Q_0(s, a)$ for all s, a .

Get initial state s

For $k = 1, 2, \dots$ till convergence

 Sample action a , get next state s'

 If s' is terminal:

$$\text{target} = R(s, a, s')$$

 Sample new initial state s'

 else:

$$\text{target} = R(s, a, s') + \gamma \max_{a'} Q_k(s', a')$$

$$\theta_{k+1} \leftarrow \theta_k - \alpha \nabla_{\theta} \mathbb{E}_{s' \sim P(s'|s, a)} [(Q_{\theta}(s, a) - \text{target}(s'))^2] |_{\theta=\theta_k}$$

$$s \leftarrow s'$$

Chasing a nonstationary target!

Updates are correlated within a trajectory!

4.2.1 Problems with DQN

Using target values, the above equations seem to be similar to those encountered in supervised learning. However, we want the input data to be I.I.D (Independently and Identically Distributed), i.e.

- Samples should be randomized among batches and the distribution of data in every batch must be same.
- Samples in a batch should be independent of each other.

¹If not, the model may overfit for some class of samples at different time and the solution will not be generalized. In addition, for the same input, its label should not change over time. These stability conditions are required for the supervised learning algorithms to work well.

In reinforcement learning, both the input and the target change constantly during the process and hence make training unstable.

- **Target unstable:** We are using the same parameters (weights) for estimating the target ³ and the Q value. As a consequence, there is a big correlation between the target and the parameters that are being changed. Therefore, at every step of training, our Q values change but the target values also change. So, we're getting closer to our target but the target is also moving. It's like chasing a moving target! This leads to amplified oscillations in training.
- **Non I.I.D.:** One of the problems related to the inter-trajectory correlations is that within a training iteration, we update model parameters to move Q function closer to the ground truth and these updates to the network will influence further estimations. When we use the Q-values in such a network, the Q-values in the surrounding states will be pulled up also like a net. For example, if we just receive a new reward and tune the Q -network to incorporate it, and then make another move. The new state will look very similar to the previous state in particular since that environment does not change drastically and hence the states are correlated. The newly estimated $Q(s', a')$ will be comparatively higher and hence the new target for Q will be higher, regardless of the merit of the new action taken. Upon modifications of the network with a sequence of actions in the same trajectory, the updates are greatly magnified. This destabilizes the learning process.

In Reinforcement Learning, we are often dependent on the value functions or the policy functions to sample actions. However, this changes frequently as we get to know better what to explore. As we proceed further in the exploration, we get to know better about the ground truth values of actions and states. So, the target outputs too change. Hence, we are trying to learn a mapping for a constantly changing input and output!

It has been experimentally discovered that both input and output can converge if the changes in the input and the output are slowed down sufficiently, we can model f while allowing it to evolve.

4.2.2 Solutions

- **Experience replay:** The past few transitions of the agent in the environment are put into a memory-buffer and mini-batch of samples is randomly sampled from this buffer to train the deep network. This forms an input dataset which is stable enough for training. As the data is sampled randomly from the replay memory buffer, the data is more independent of each other and closer to i.i.d.
- **Target network:** Two deep neural networks having parameters θ^- and θ are created. The first one is used to retrieve the Q values while all the updates during the training process are made in the second network. After a fixed number of updates, we synchronize θ^- and θ . The objective is to update the Q -value targets temporarily so that we don't have to evaluate with a target that is constantly changing. Also, parameter changes do not affect θ^- immediately and therefore even though the input may not be completely i.i.d., it is not expected to incorrectly magnify the effects as mentioned above.

$$L_i(\theta_i) = \mathbb{E}_{s,a,s',r \sim D} \left(\underbrace{r + \gamma \max_{a'} Q(s', a'; \theta_i^-)}_{\text{target}} - Q(s, a; \theta_i) \right)^2$$

On using both experience replay and the target network, we have a more stable target values and hence the training of the network resembles supervised learning more.

Algorithm

```

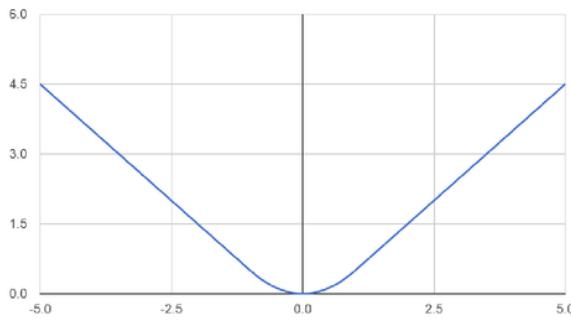
Initialize replay memory  $D$  to capacity  $N$ 
Initialize action-value function  $Q$  with random weights  $\theta$ 
Initialize target action-value function  $\hat{Q}$  with weights  $\theta^- = \theta$ 
For episode = 1,  $M$  do
    Initialize sequence  $s_1 = \{x_1\}$  and preprocessed sequence  $\phi_1 = \phi(s_1)$ 
    For  $t = 1, T$  do
        With probability  $\varepsilon$  select a random action  $a_t$ 
        otherwise select  $a_t = \operatorname{argmax}_a Q(\phi(s_t), a; \theta)$ 
        Execute action  $a_t$  in emulator and observe reward  $r_t$  and image  $x_{t+1}$ 
        Set  $s_{t+1} = s_t, a_t, x_{t+1}$  and preprocess  $\phi_{t+1} = \phi(s_{t+1})$ 
        Store transition  $(\phi_t, a_t, r_t, \phi_{t+1})$  in  $D$ 
        Sample random minibatch of transitions  $(\phi_j, a_j, r_j, \phi_{j+1})$  from  $D$ 
        Set  $y_j = \begin{cases} r_j & \text{if episode terminates at step } j+1 \\ r_j + \gamma \max_{a'} \hat{Q}(\phi_{j+1}, a'; \theta^-) & \text{otherwise} \end{cases}$ 
        Perform a gradient descent step on  $(y_j - Q(\phi_j, a_j; \theta))^2$  with respect to the
        network parameters  $\theta$ 
        Every  $C$  steps reset  $\hat{Q} = Q$ 
    End For
End For

```

4.2.3 Implementation Details

- **Loss Function**

¹ DQN uses Huber Loss function, where the loss is quadratic for small values of x , and linear for large values.



$$L_\delta(a) = \begin{cases} \frac{1}{2}a^2 & \text{for } |a| \leq \delta, \\ \delta(|a| - \frac{1}{2}\delta), & \text{otherwise.} \end{cases}$$

- **Optimization**

Training of reinforcement learning agents is sensitive to optimization methods used. To handle the changes in the input during the training, a simple learning rate schedule does not suffice because it is not dynamic enough. Therefore, RMSProp or Adam optimizers are often used while training a Deep Q-Network.

4.3 DOUBLE DQN

In many environments that possess stochasticity, the well-known Q-learning reinforcement learning algorithm performs somewhat poorly. This poorness of the performance can be attributed to the large overestimations of the action values. Such overestimations are often resulted from a positive bias that is often introduced into the system because of the use of maximum action value by the Q-learning as an approximation for the maximum expected action value (Q-value).

³ Double Deep Q-Networks, or double Learning, was introduced by Hado van Hasselt, which handles the problem of the overestimation of Q-values. This is how we used to calculate the Temporal Difference (TD) Target –

$$Q(s, a) = \underbrace{r(s, a)}_{\text{Q target}} + \underbrace{\gamma \max_a Q(s', a)}_{\substack{\text{Reward} \\ \text{of taking} \\ \text{that action} \\ \text{at that state}}} \quad \begin{array}{l} \text{Discounted max q} \\ \text{value among all.} \\ \text{possible actions} \\ \text{from next state.} \end{array}$$

³ However, we cannot be sure whether the best action for the next state is the action with the highest Q-value. The accuracy of Q-values depends on what actions were tried and what neighbouring states were explored.

As a consequence, at the beginning of the training we don't have enough information with us to know about the best action to take. Therefore, taking the maximum Q-value (which is noisy due to random initializations) as the best action to take may lead to false positives. If we assign a higher Q-values to non-optimal actions than the optimal best action, the learning will tend to be complicated.^[7]

³ The solution is to decouple the action selection from the target Q-value generation by using two networks simultaneously when we compute the Q target.

- Use DQN network to select the best possible action to be taken for the next state (which is the action that has the highest Q value).³
- Use target network to calculate the target Q value of taking that action at the next state.

$$Q(s, a) = r(s, a) + \gamma \underbrace{Q(s', \operatorname{argmax}_a Q(s', a))}_{\substack{\text{TD target} \\ \text{DQN Network choose} \\ \text{action for next state}}}$$

Target network calculates the Q value of taking that action at that state

Therefore, Double DQN helps us reduce the overestimation of Q-values and, as a consequence, helps us train faster and have more stable learning process.

Algorithm

```

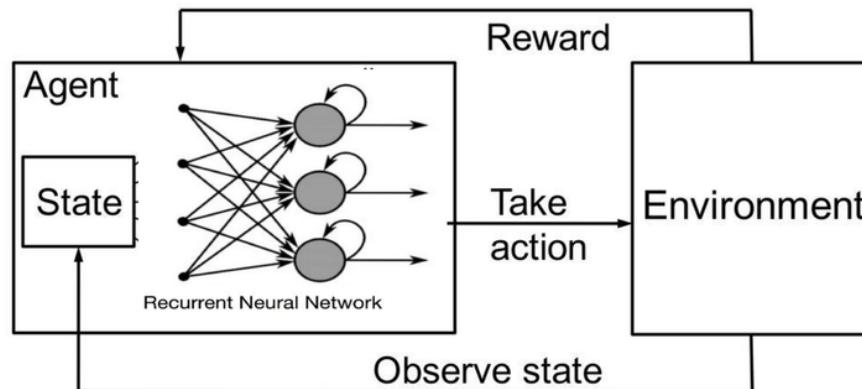
1: Initialize  $Q^A, Q^B, s$ 
2: repeat
3:   Choose  $a$ , based on  $Q^A(s, \cdot)$  and  $Q^B(s, \cdot)$ , observe  $r, s'$ 
4:   Choose (e.g. random) either UPDATE(A) or UPDATE(B)
5:   if UPDATE(A) then
6:     Define  $a^* = \arg \max_a Q^A(s', a)$ 
7:      $Q^A(s, a) \leftarrow Q^A(s, a) + \alpha(s, a)(r + \gamma Q^B(s', a^*) - Q^A(s, a))$ 
8:   else if UPDATE(B) then
9:     Define  $b^* = \arg \max_a Q^B(s', a)$ 
10:     $Q^B(s, a) \leftarrow Q^B(s, a) + \alpha(s, a)(r + \gamma Q^A(s', b^*) - Q^B(s, a))$ 
11:   end if
12:    $s \leftarrow s'$ 
13: until end

```

4.4 DEEP RECURRENT Q-NETWORKS

⁹ Deep Q-Networks are limited in their scope as they learn a functional mapping from a discrete, limited number of past states. Thus they will be unable to master the processes that require them to remember even further events in the past. Such processes do not only depend on the current state of the agent. Hence, the process becomes Partially Observable Markov Decision Process (POMDP) instead of a Markov Decision Process (MDP). Many real world tasks often feature such incomplete and noisy state information resulting from partial observability of the states.

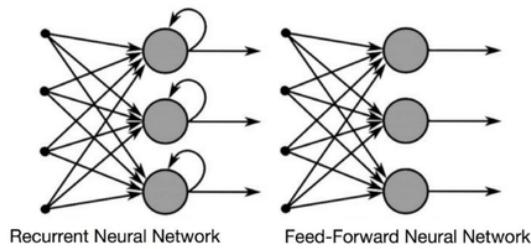
Deep Q-Networks' performance declines when it is made to estimate Q-functions from partially observable states, due to the lack of memory. However, the existing DQN can be updated and modified to better deal with POMDP by grasping the power of Recurrent Neural Networks (RNN).



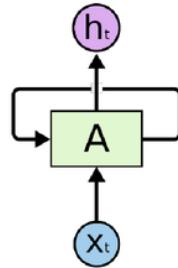
4.4.1 Recurrent Neural Networks

Recurrent Neural Networks (RNNs) are an improvement over the existing neural network architecture that possess the ability to capture the past trends in the input. This property is of utmost importance when the data that we have is sequential in format. Some problems that require such an architecture are image captioning, financial data modelling, named entity recognition, etc. Famous applications include Apple's digital personal voice assistant Siri, Google Translate, etc.

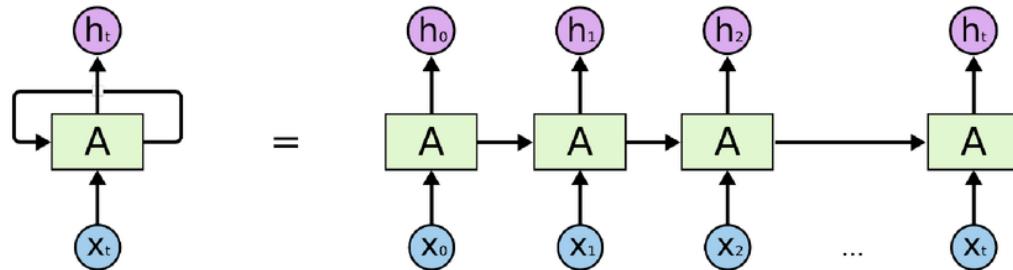
RNNs are a very powerful class of algorithms that are built on top of the existing neural networks by introducing the concept of network 'memory', that remembers the past states as well.



The information cycles through a loop in an RNN. It takes into consideration the inputs that it has received what it has received previously, along with the current input when it makes a decision.



7 In the above diagram, a collection of neural networks, A , takes some input x_t and outputs the value h_t . The loops shown in the network diagram above allows the information to be passed from one step of the network to the next. The loops may make recurrent neural networks seem sort of strange and mysterious. However, 7 on careful observation, it can be seen that they aren't much different from the vanilla normal neural network. A recurrent neural network (RNN) can be assumed of as multiple copies of the same network, each passing a message to a successor over time. Consider what happens if we unroll the loop of the RNN:

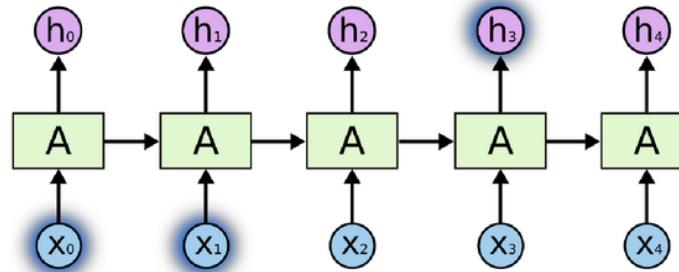


This appears to be a stacked layers of various neural networks, with the exception that the layers are stacked across time, not across space. The values h_t are the hidden units

that represent the memory of the recurrent neural network. This allows the network to capture and extract information from historical input data and use it further down the time.

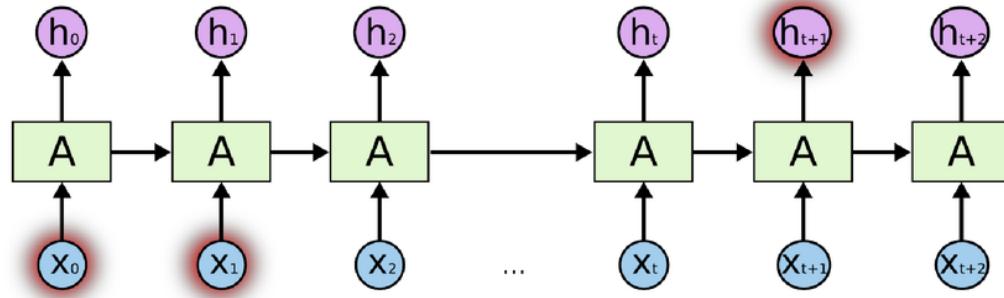
4.4.2 Long Short Term Memory (LSTM)

One of the most appealing ideas of RNNs is that they might be able to attach and employ previous information to be used in the present task, such as utilizing the previous stock movement ticks can provide information about the understanding of the scenario. In cases where the relative time gap between the relevant information and the place that it's required is small, RNNs can learn to grasp and leverage the past information present in the hidden units.



However, as the above mentioned gap grows, RNNs increasingly become unable to learn to use the past information in the hidden units and connect it with the present scenario. In theory, RNNs have been proven to be capable of managing such “long-term dependencies.” Parameters could be carefully chosen for the RNN to solve the problems of this form. Sadly, it has been observed that in practice that RNNs don't learn such a

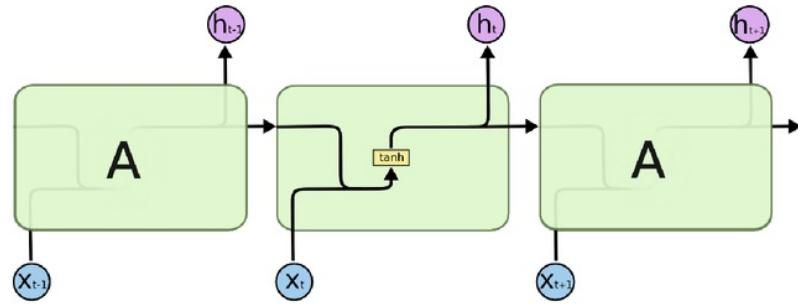
⁴ technique easily. This was explored in depth by Hochreiter (1991) and Bengio, et al. (1994), who found some pretty fundamentally strong reasons for the difficulty.



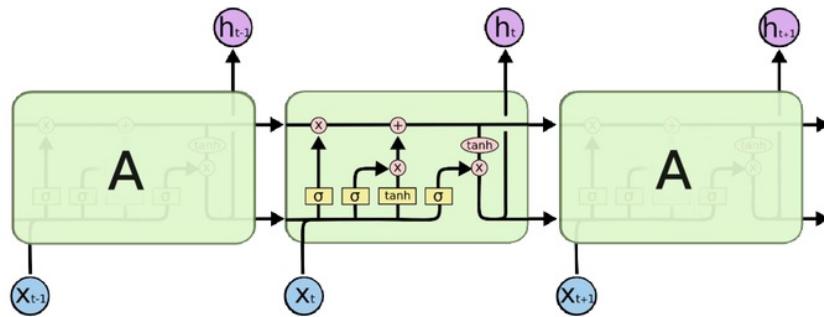
²¹ Long Short Term Memory (LSTM) networks, are used ¹¹ to solve this problem. “LSTMs” are a special kind of RNN which are designed to be capable of learning such long-term dependencies. They were first introduced by Hochreiter & Schmidhuber in 1997, and were popularized and refined by many people over time. They work extremely well on a large variety of problems, and are now being widely used in a number of AI applications.

LSTMs were explicitly designed to overcome this long-term dependency problem. It's their default behaviour to remember information over time periods that are spread over long durations, ⁴ not something they have to struggle to learn!

All recurrent neural networks (RNNs) have the structure of a chain of repeating neural network modules over time. In standard RNNs, this repeating module has a very simple and lucid structure.



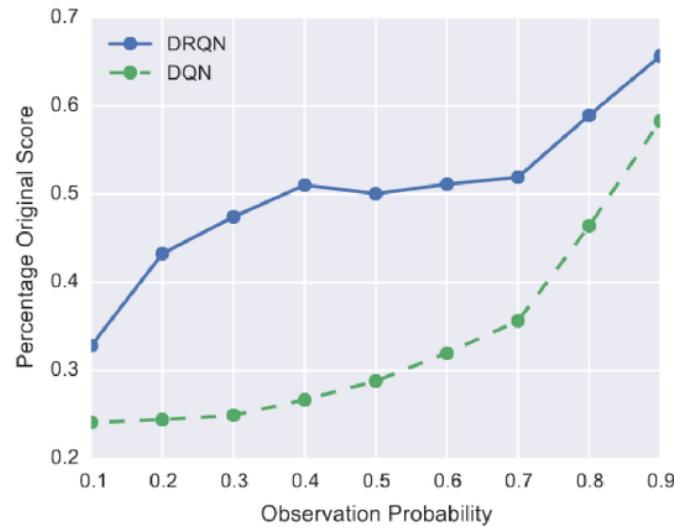
This chain like structure is also present in LSTMs, but the structure of the repeating module is all that makes a difference. Instead of having a single neural network layer, there are four separate networks, each interacting with each other in a very different way.



4.4.3 Deep Recurrent Q-Networks (DRQN)

We modify the existing DQN network and append an LSTM layer in the end of the dense layers to capture the long term dependencies.

A recurrent Q-network could also be trained on a standard Markov Decision Process environment and then generalized to a Partially Observable Markov Decision Process at runtime. Matthew Hausknecht and Peter Stone showed in 2017 that it can indeed be done and the model developed by incorporating RNN in Deep Q-Learning algorithm.^[8]



CHAPTER 5

EXPERIMENTAL IMPLEMENTATION

5.1 DATA

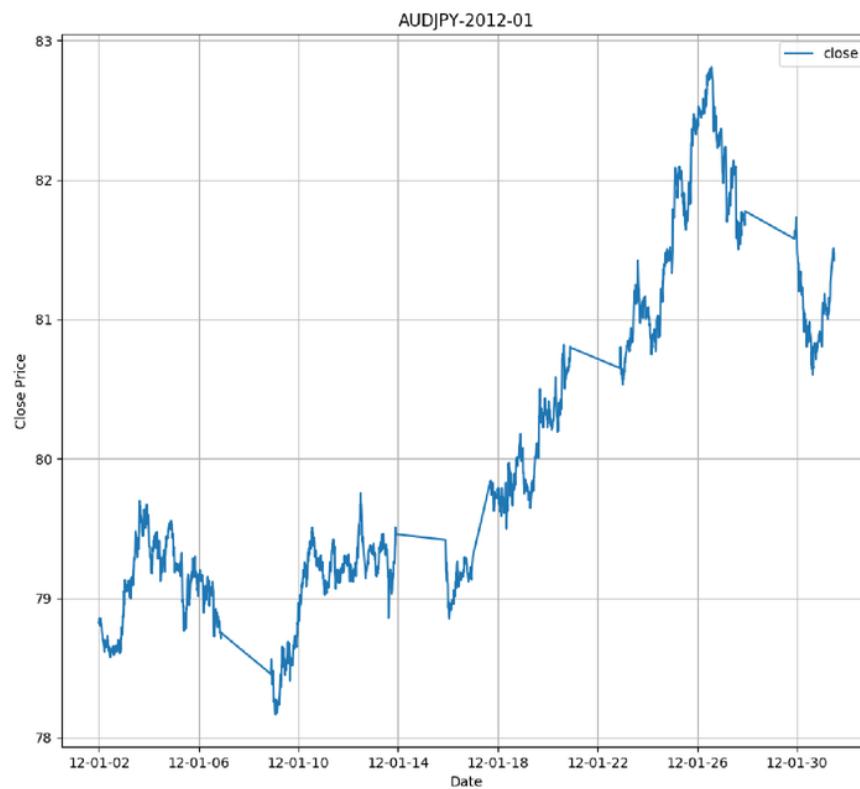
5.1.1 Source

We decide to operate the model on forex data obtained from TrueFX.com. The data collected is highly precise tick-by-tick information of bid and ask prices of 15 currency pairs, given below:

AUDJPY, AUDNZD, AUDUSD,
CADJPY, CHFJPY, EURCHF,
EURGBP, EURJPY, EURUSD,
GBPJPY, GBPUSD, NZDUSD,
USDCAD, USDCHF, USDJPY

The average time interval between ticks ranges between a few milliseconds, with a mean volume of around 2000 ticks per 15-minutes.

The resolution of the captured data is seen in the following closing-price chart.



The data was obtained in an automated manner, by iteratively fetching and extracting archives from the online repository provided by TrueFX. In its entirety, the input data set comprises of 1080 files taking 26 GB of compressed disk space. The availability of this breadth and depth of data allows the model to learn temporal correlations between pairs and the intrinsic periodicity of each asset over time.

5.1.2 Pre-processing

To obtain inferential features from this data, we resample the given data based on time intervals of 15-minutes each.

The first order features obtained from this resampling contain the Open, High, Low and Close prices of the base currency in relation to the counter currency along with the tick volume, i.e. the number of price movements within the interval.

This resampling allows us to visualise the price movement over larger intervals of time, as is shown below.



This pre-processed data is then used to calculate the features required to compute the input state of the financial trading MDP, as discussed below.

5.2 ENVIRONMENT

The training and evaluation of the agent are done in an environment modelling the financial trading task as an MDP, with the following attributes:

5.2.1 Action Space

¹² The action space describes the set of actions available to the agent at any time step during the episode. In a trading setting, these correspond to the decisions of buying, selling or holding a given asset at a given time. We encode this choice as an integer, described as 0 (hold), 1 (buy), 2 (sell).

The position held on the asset currently is denoted by a one-hot encoding of this index.

¹² 5.2.2 State space

The state space is used to capture the current state of the environment as seen by the agent. Here, the information a trader sees is the history of price movements and positions held by the trader at those instances. This information is encapsulated into the following features which make up the state space:

- Periodic Time features: These contain the information about the time session at which the environment is operating. The periodicity is encoded using a sine function for representing the periodic repetition of the minute, hour, the day of the week and month of the year of the tick. The four features are thus computed as

$$f_t = \sin\left(\frac{2\pi t}{T}\right)$$

t = current index of the time measure, T = maximum size of the measure

For example, for the first minute of the hour, $t = 0$ and $T = 60$ similarly, for Tuesday, $t = 2$ and $T = 7$ and so on.

- Market history features: These features convey information regarding the current and recent history of the market. We denote this history using the eight most recent log returns on the closing price and tick volume. Thus, 16 features per currency pair included.
If we choose to include n currency pairs in our observation, the number of market history features then becomes = $16xn^2$
- Trading position feature: This represents the recent history of trades undertaken by the agent. This feature is a one-hot encoding of the trading position as denoted by the action space, adding three more dimensions to the state.

Thus, for an environment operating on n currency pairs, the total number of features is = $7 + n * 16$.

For $n = 1$ (Single currency pair), the dimensionality of the state space is 23.

Thus, at each time step, the agent navigates the trading environment by taking a step from the action space according to its policy, computing the relevant Q-values to choose the best action given the state as described above.

5.2.3 Reward

The reward is the measure of the incentive received by the agent for taking a particular action in the environment. We use the most natural incentive for trading, the return obtained at each time step.

The reward function is the logarithmic portfolio returns for the time step at which the action was taken

$$r_t = \log\left(\frac{v_t}{v_{t-1}}\right)$$

Where v_t = portfolio value at time t , r_t = reward obtained at time t

The value of the portfolio is calculated as follows

$$v_t = v_{t-1} + a_t \cdot c \cdot (c_t - o_t)$$

Where, v_t = portfolio value time t ,

a_t = action taken at time t ,

c_t = closing price of the currency pair at time t

o_t = opening price of currency pair at time t

c = trade volume (1 unit in our case)

5.3 AGENT

We seek to completely exploit the environment with the help of action augmentation. Unsatisfactory results are obtained from random exploration especially in the domain of financial algorithmic trading domain since there are transaction costs with every change of position. Hence, we shall use a simple bypass technique to mitigate the need for random exploration by providing the agent with reward signals if it had taken all the other actions as well. This is easily possible since the reward is easily computable after the price of the current timestamp is revealed.

The agent's position is the only thing that is altered if we take other actions that deviate from the current policy. This is known as the zero market impact hypothesis which states that the action taken from the reinforcement learning agent has no influence on the current market condition.

Now we can update Q-values for all actions. The loss function is defined by –

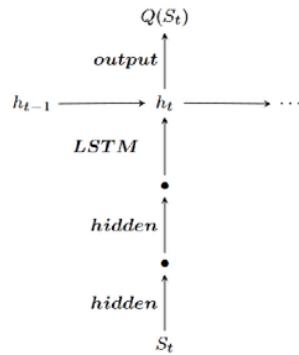
$$L(\theta) = E_{(s,a,r,s') \sim D} [\| r + \gamma \cdot \text{argmax}_{a'} \{ Q_\theta(s', a') \} - Q_\theta(s, a) \|^2]$$

Where Q_{θ^-} denotes the target network.

5.4 MODEL ARCHITECTURE AND TRAINING

5.4.1 Weight/Parameter Initialisation

Proper initialization of weights is extremely critical for training of deep neural networks effectively.² We follow the initialization scheme presented in He et al. (2015) for weight matrices in both hidden layers and input-to-hidden layer in LSTM in our approach. We follow Le et al. (2015) to initialize all hidden-to- hidden weight matrices to be identity. All bias weights in the network are set to be zero except for the forget gates in the LSTM which are set to be 1. We sparsely initialize the output layer weight matrix with Gaussian distribution $N(0, 0.001)$.



5.4.2 Training Scheme

1. A smaller replay memory is used as it is proven to be more effective in DRQN for financial domain. The reason is that in financial trading, recent

data points have more predictive importance than the points which are further in the past.

2. A longer sequence is sampled from the replay memory than the conventionally used values. The reason behind this modification is that, a successful trading strategy involves opening a position at the right time and holding the position for a significantly long period of time before closing it. Long term dependency cannot be captured by a smaller sequence and hence it can't train the network as efficiently.
2. We can mitigate the need for training the network for each step since we are sampling a longer sequence. Hence we now train the network for every T time steps. This significantly reduces the computation since the number of backward passes are reduced by a factor of T.

5.4.3 Complete Learning Algorithm

Adopting the above learning scheme, the following algorithm is finally used for the training of the Deep Recurrent Q-Network. The algorithm has been named as “Financial Deep Recurrent Q-Network (FDRQN)”. The code for the training is written in Python. The algorithm is optimised to capture the nuances of Financial Trading.

Algorithm 1 Financial DRQN Algorithm

- 1: Initialize $T \in \mathbb{N}$, recurrent Q-network Q_θ , target network Q_{θ^-} with $\theta^- = \theta$, dataset \mathcal{D} and environment E , $steps = 1$
- 2: Simulate env E from dataset \mathcal{D}
- 3: Observe initial state s from env E
- 4: **for** each step **do**
- 5: $steps \leftarrow steps + 1$
- 6: Select greedy action w.r.t. $Q_\theta(s, a)$ and apply to env E
- 7: Receive reward r and next state s' from env E
- 8: Augment actions to form $\mathcal{T} = (s, a, r, s')$ and store \mathcal{T} to memory \mathcal{D}
- 9: **if** \mathcal{D} is filled and $steps \bmod T = 0$ **then**
- 10: Sample a sequence of length T from \mathcal{D}
- 11: Train network Q_θ with equation (4) and (5)
- 12: **end if**
- 13: Soft update target network $\theta^- \leftarrow (1 - \tau)\theta^- + \tau\theta$
- 14: **end for**

5.4.4 Hyper-parameters

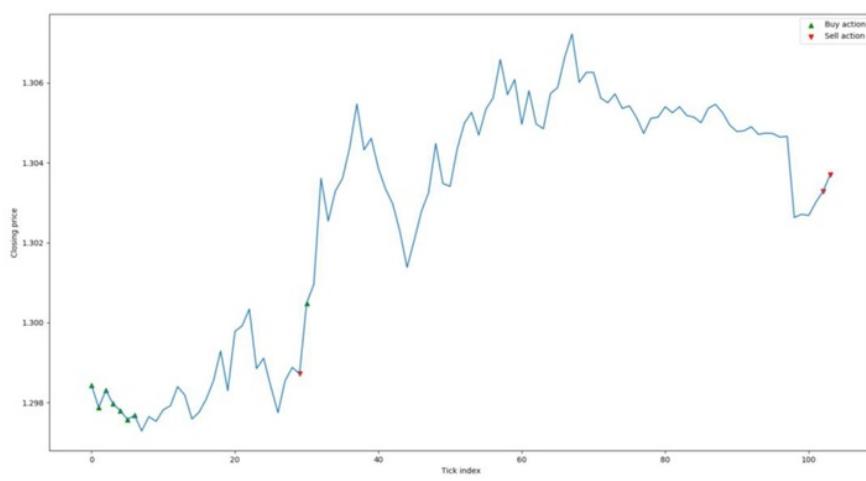
Hyper-parameters	Value
Learning time-step T	96
Replay memory size N	480
Learning rate	0.00025
Optimizer	Adam
Discount factor γ	0.99
Target network τ	0.001

5.5 Model Evaluation and Testing

We use the following model architecture – two dense fully connected layers having 32 neurons each, followed by an LSTM layer having 64 neurons. The agent was able to learn the Q-function from the data sufficiently well. The following are the results obtained during one such episode of model running –

```
\ tradegame/src : python3 evaluate.py EURUSD_2012_norm /working_2012/model_ep10
Using TensorFlow backend.
2018-12-07 07:20:04.982210: I tensorflow/stream_executor/cuda/cuda_gpu_executor.cc:964] successful NUMA node read from SysFS had negative value (-1), but there must be at least one NUMA node, so returning NUMA node zero
2018-12-07 07:20:04.982613: I tensorflow/core/common_runtime/gpu/gpu_device.cc:1432] Found device 0 with properties:
name: Tesla K80 major: 3 minor: 7 memoryClockRate(GHz): 0.8235
pciBusID: 0000:00:04.0
totalMemory: 11.17GiB freeMemory: 11.10GiB
2018-12-07 07:20:04.982645: I tensorflow/core/common_runtime/gpu/gpu_device.cc:1511] Adding visible gpu devices: 0
2018-12-07 07:20:05.292523: I tensorflow/core/common_runtime/gpu/gpu_device.cc:982] Device interconnect StreamExecutor with strength 1 edge matrix:
2018-12-07 07:20:05.292588: I tensorflow/core/common_runtime/gpu/gpu_device.cc:988]      0
2018-12-07 07:20:05.292597: I tensorflow/core/common_runtime/gpu/gpu_device.cc:1001] 0:   N
2018-12-07 07:20:05.292881: I tensorflow/core/common_runtime/gpu/gpu_device.cc:1115] Created TensorFlow device (/job:localhost/replica:0/task:0/device:GPU:0 with 10759 MB memory) -> physical GPU (device: 0, name: Tesla K80, pci bus id: 00:00:04.0, compute capability: 3.7)
Evaluating the model on sampled data points.

-----
Evaluation parameters
-----
Start cash : 10000
Trade size : 1000
-----
Buy: Rate = $1.2984300 | Cost = 1298.43
Buy: Rate = $1.2978800 | Cost = 1297.8799999999999
Buy: Rate = $1.2983100 | Cost = 1298.3100000000002
Buy: Rate = $1.2979700 | Cost = 1297.97
Buy: Rate = $1.2978900 | Cost = 1297.8
Buy: Rate = $1.2975800 | Cost = 1297.58
Buy: Rate = $1.2976800 | Cost = 1297.679999999998
Sell: $1.2987200 | Profit: $0.2900000
Buy: Rate = $1.3004900 | Cost = 1300.49
Sell: $1.3032800 | Profit: $5.4000000
Sell: $1.3037000 | Profit: $5.3900000
-----
Closing all open positions to compute portfolio value
Total Profit: $38.0600000
Final portfolio worth: $10038.0600000
-----
```



APPENDIX-A

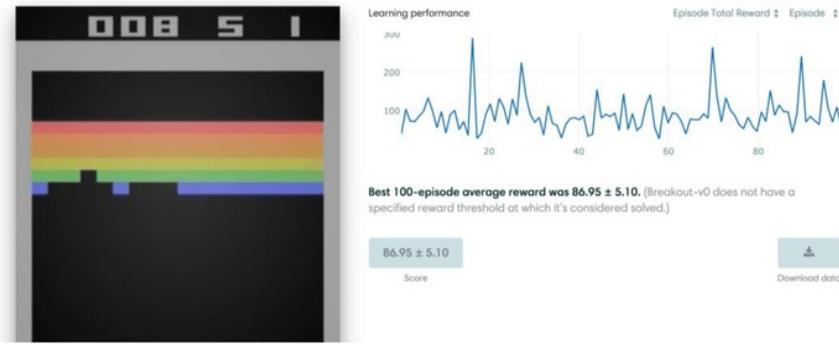
1. OPEN AI GYM

⁵ OpenAI is a not-for-profit artificial intelligence research company that seeks to develop and develop human friendly AI systems in such a way as to benefit the people by enhancing their technological experiences with the help of AI. The company was founded in late 2015 in San Francisco by Elon Musk and Sam Altman. The OpenAI organization aims to “freely collaborate” with other researchers and organisations by making its research and open to the people.

The most important product released by OpenAI is the OpenAI Gym (or simply Gym). ⁵ Gym strives to provide a general-intelligence benchmark with a wide variety of different environments that are somewhat similar to, but broader than, the ImageNet Large Scale Visual Recognition Challenge (ILSVRC) used in supervised learning research – hoping to standardize the way in which environments are defined in AI research publications, so that published research becomes more easily reproducible. The environment is fairly easy to set-up and the user is provided with a simple interface. Currently, ⁵ Gym can only be used with Python.

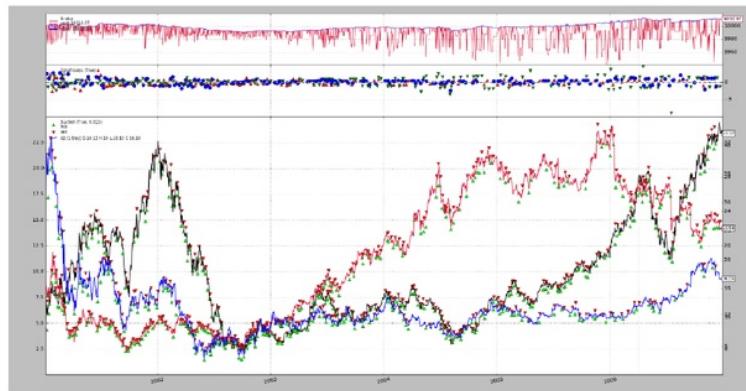
OpenAI ¹⁰ Gym comes with a diverse range of environments that range from easy to difficult and involve many different kinds of data -

- ¹⁰ Classic control and toy text: complete small-scale tasks, mostly from the RL literature.
- Algorithmic: perform computations such as adding multi-digit numbers and reversing sequences.
- Atari: play classic Atari games.
- 2D and 3D robots: control a robot in simulation.



2. BACKTRADER

Backtrader is a feature-rich Python framework for back-testing and trading. Backtrader allows the user to focus on writing reusable trading strategies, indicators and analysers instead of having to spend time building infrastructure. The platform has two main features – ease of use and reproducibility. The user can create a stock trading strategy by deciding the potential adjustable parameters, instantiating the indicators required in the strategy and then writing down the logic for entering/exiting the market.



3. BT GYM

BTGym is an OpenAI Gym-compatible environment for Backtrader back-testing/trading library, designed to provide gym-integrated framework for running reinforcement learning experiments in [close to] real world algorithmic trading environments.



REFERENCES

- [1]. Chien Yi Huang , "Financial Trading as a Game: A Deep Reinforcement Learning Approach", arXiv:1807.02787[q-finTR], 8 July 2018
- [2]. C. T. Chen, A. Chen and S. Huang, "Cloning Strategies from Trading Records using Agent-based Reinforcement Learning Algorithm," *2018 IEEE International Conference on Agents (ICA)* , Singapore, 2018, pp. 34-37.
- [3]. Y. Deng, F. Bao, Y. Kong, Z. Ren and Q. Dai, "Deep Direct Reinforcement Learning for Financial Signal Representation and Trading," in *IEEE Transactions on Neural Networks and Learning Systems* , vol. 28, no. 3, pp. 653-664, March 2017.
- [4]. Z. Jiang and J. Liang, "Cryptocurrency portfolio management with deep reinforcement learning," *2017 Intelligent Systems Conference (IntelliSys)* , London, 2017, pp. 905-913.
- [5]. Jae Won Lee, "Stock price prediction using reinforcement learning," *ISIE 2001. 2001 IEEE International Symposium on Industrial Electronics Proceedings (Cat. No.01TH8570)* , Pusan, South Korea, 2001, pp. 690-695 vol.1.
- [6]. K. Kashihara, "Deep Q learning for traffic simulation in autonomous driving at a highway junction," *2017 IEEE International Conference on Systems, Man, and Cybernetics (SMC)* , Banff, AB, 2017, pp. 984-988.
- [7]. Y. Zhang, P. Sun, Y. Yin, L. Lin and X. Wang, "Human-like Autonomous Vehicle Speed Control by Deep Reinforcement Learning with Double Q-Learning," *2018 IEEE Intelligent Vehicles Symposium (IV)* , Changshu, 2018, pp. 1251-1256.

[8]. J. Zeng, J. Hu and Y. Zhang, "Adaptive Traffic Signal Control with Deep Recurrent Q-learning," *2018 IEEE Intelligent Vehicles Symposium (IV)* , Changshu, 2018, pp. 1215-1220.

Deep reinforcement

ORIGINALITY REPORT



PRIMARY SOURCES

1	medium.com Internet Source	5%
2	arxiv.org Internet Source	4%
3	Submitted to Georgia Institute of Technology Main Campus Student Paper	2%
4	mobiles-han.blogspot.com Internet Source	2%
5	en.wikipedia.org Internet Source	1%
6	wwwsyseng.rsise.anu.edu.au Internet Source	1%
7	Feng Liu, Zhigang Chen, Jie Wang. "Video image target monitoring based on RNN-LSTM", Multimedia Tools and Applications, 2018 Publication	1%
8	Submitted to Amity University Student Paper	1%

9	www.cs.utexas.edu Internet Source	1 %
10	gym.openai.com Internet Source	<1 %
11	Submitted to University of Florida Student Paper	<1 %
12	scholar.sun.ac.za Internet Source	<1 %
13	Submitted to University of Sheffield Student Paper	<1 %
14	oa.upm.es Internet Source	<1 %
15	Submitted to University of Leicester Student Paper	<1 %
16	Submitted to University College London Student Paper	<1 %
17	Meyer, Stephan, Martin Wagener, and Christof Weinhardt. "Politically Motivated Taxes in Financial Markets: The Case of the French Financial Transaction Tax", Journal of Financial Services Research, 2014. Publication	<1 %
18	Kai Arulkumaran, Marc Peter Deisenroth, Miles Brundage, Anil Anthony Bharath. "Deep	<1 %

Reinforcement Learning: A Brief Survey", IEEE Signal Processing Magazine, 2017

Publication

19

Submitted to Technische Universiteit Delft

Student Paper

<1 %

20

Kathy Thi Aung. "Reinforcement learning using Voronoi space division", Artificial Life and Robotics, 12/2010

Publication

<1 %

21

"Intelligent Systems and Applications", Springer Nature America, Inc, 2019

Publication

<1 %

22

Springer Series in Bio-/Neuroinformatics, 2015.

Publication

<1 %

23

Sunil Kumar Rajendran, Qi Wei, Feitian Zhang. "Foveation Control of a Robotic Eye Using Deep Reinforcement Learning", Volume 1: Advances in Control Design Methods; Advances in Nonlinear Control; Advances in Robotics; Assistive and Rehabilitation Robotics; Automotive Dynamics and Emerging Powertrain Technologies; Automotive Systems; Bio Engineering Applications; Bio-Mechatronics and Physical Human Robot Interaction; Biomedical and Neural Systems; Biomedical and Neural Systems Modeling, Diagnostics, and Healthcare, 2018

<1 %

- 24 Mehmet F. Dicle, John Levendis. "Importing Financial Data", The Stata Journal: Promoting communications on statistics and Stata, 2018 <1 %
- Publication
-
- 25 Yongheng Wang, Shaofeng Geng, Hui Gao. "A proactive decision support method based on deep reinforcement learning and state partition", Knowledge-Based Systems, 2018 <1 %
- Publication
-
- 26 "Workshop on ADASIVA at IIIT-A.(Education)", The Times of India <1 %
- Publication
-
- 27 Armando Vieira, Bernardete Ribeiro. "Introduction to Deep Learning Business Applications for Developers", Springer Nature, 2018 <1 %
- Publication
-
- 28 Yi Zhang, Ping Sun, Yuhan Yin, Lin Lin, Xuesong Wang. "Human-like Autonomous Vehicle Speed Control by Deep Reinforcement Learning with Double Q-Learning", 2018 IEEE Intelligent Vehicles Symposium (IV), 2018 <1 %
- Publication
-
- 29 www.culturalpolicies.net <1 %
- Internet Source
-

Exclude quotes On Exclude matches < 10 words

Exclude bibliography On