



Department of Informatics  
King's College London  
United Kingdom

7CCSMPRJ - Individual Project

# Explainable Deep Learning for Portfolio Optimisation

---

Name: **Ingrid Pérez Aguilera**

Student Number: K24087939

Course: Computational Finance M.Sc.

Supervisor: **Riaz Ahmad**



# Abstract

It is a précis of the report (normally in one page), which should include:

- A brief introduction to the project objectives
- A brief description of the main work of the project
- A brief description of the contributions, major findings, results achieved and principal conclusion of the project

# Acknowledgements

It is a short paragraph to thank those who have contributed to the project work.

# Table of Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
1.1	Objectives . . . . .	3
1.2	Report Structure . . . . .	4
<b>2</b>	<b>Background</b>	<b>6</b>
2.1	Portfolio Optimisation . . . . .	6
2.1.1	Modern Portfolio Theory . . . . .	7
2.2	Deep Reinforcement Learning . . . . .	9
2.2.1	Reinforcement Learning . . . . .	10
2.2.2	Deep Reinforcement Learning Algorithms . . . . .	16
<b>3</b>	<b>Methodology</b>	<b>23</b>
3.1	Problem Definition . . . . .	23
<b>4</b>	<b>Results</b>	<b>25</b>
<b>5</b>	<b>Legal, Social, Ethical and Professional Issues</b>	<b>26</b>
<b>6</b>	<b>Conclusion</b>	<b>27</b>
	<b>References</b>	<b>i</b>
<b>A</b>	<b>DRL Algorithms</b>	<b>vii</b>

# List of Figures

2.1	Efficient Frontier in Risk-Return Space. [1]	9
2.2	Agent interaction with environment	11
2.3	Markov Decision Process with policy, transition, and reward functions	12
2.4	Taxonomy of Reinforcement Learning Algorithms [2]	17

# List of Tables

2.1	Comparison of Deep Reinforcement Learning Algorithms . . . . .	22
-----	----------------------------------------------------------------	----

# List of Algorithms

1	Advantage Actor-Critic (A2C) Pseudo-code . . . . .	viii
2	Proximal Policy Optimization (PPO) Pseudo-code . . . . .	ix
3	Deep Deterministic Policy Gradient (DDPG) Pseudo-code . . . . .	x
4	Twin Delayed Deep Deterministic Policy Gradient (TD3) Pseudo-code . .	xi
5	Soft Actor-Critic (SAC) Pseudo-code . . . . .	xii



# Nomenclature

$c$       Speed of light in a vacuum

$h$       Planck constant

# Glossary

**Advantage Actor-Critic** Algorithm that uses both an actor (policy) and a critic (value function) to learn optimal policies by estimating the advantage of actions taken. xi, 3, 16, 17

**Algorithmic Trading** Use of computer algorithms to automate trading decisions and execute trades in financial markets. 1

**Artificial Intelligence** Simulation of human intelligence processes by machines, especially computer systems, enabling them to perform tasks that typically require human intelligence, such as visual perception, speech recognition, decision-making, and language translation. xi, 2, 9

**Bellman Equations** Set of equations that describe the relationship between the value of a state or action and the values of subsequent states or actions in dynamic programming and reinforcement learning, used to compute optimal policies. 15, 20

**Deep Deterministic Policy Gradient** Algorithm that uses deep neural networks to learn policies for continuous action spaces, combining the benefits of deep learning and policy gradient methods. x, xi, 3, 16, 19

**Deep Learning** Subset of machine learning that uses neural networks with many layers

to learn from large amounts of data, enabling the model to automatically learn complex patterns and representations. xi, 2, 4, 6, 10, 17

**Deep Neural Network** Type of artificial neural network with multiple layers that can learn complex representations and patterns in data. xi, 10

**Deep Reinforcement Learning** Combination of deep learning and reinforcement learning, where deep neural networks are used to approximate the value function or policy in reinforcement learning tasks. xi, 3, 6, 10, 23, 24

**Efficient Frontier** Set of optimal portfolios that offer the highest expected return for a given level of risk, or the lowest risk for a given level of expected return, in the context of portfolio optimisation. 8

**Entropy** Measure of uncertainty or randomness in a system, often used in the context of information theory and reinforcement learning to encourage exploration. 21

**Explainable Artificial Intelligence** Methods and techniques in the application of artificial intelligence that make the results of the models understandable by humans, providing insights into how decisions are made. xii, 2

**Feature Importance** Technique for determining the contribution of each feature in a machine learning model to its predictions, helping to identify which features are most influential. 4, 6

**Financial Markets** Marketplaces where people trade financial securities, commodities, and other fungible items of value at low transaction costs and at prices that reflect supply and demand. 1

**Hyper-parameter** Parameter that is set before the learning process begins and control the learning process of a machine learning model, such as learning rate, batch size, and number of layers in a neural network. 2

**Local Interpretable Model-agnostic Explanations** Technique for explaining the predictions of any machine learning model by approximating it with a locally interpretable model, allowing users to understand the model’s behaviour in a specific instance. xi, 4–6

**Machine Learning** Subset of artificial intelligence that enables systems to learn from data and improve their performance over time without being explicitly programmed. xi, 1, 6, 9

**Markov Decision Process** Mathematical framework for modelling decision-making in situations where outcomes are partly random and partly under the control of a decision maker, characterised by states, actions, rewards, and transition probabilities. ix, xi, 11–13, 23

**Modern Portfolio Theory** Investment theory that aims to construct a portfolio of assets that maximises expected return for a given level of risk, or minimises risk for a given level of expected return, through diversification. xi, 7

**Partially Observable Markov Decision Process** Extension of Markov Decision Process where the agent does not have access to the complete state information, requiring the use of belief states or observations to make decisions. xi, 13

**Portfolio Optimisation** Process of selecting the best distribution of assets in a portfolio to achieve specific investment goals, such as maximising returns or minimising risk, while considering constraints and preferences. 1, 6, 23

**Proximal Policy Optimisation** Algorithm that optimises policies by ensuring that updates to the policy are not too large, maintaining a balance between exploration and exploitation. xii, 3, 16, 18

**Reinforcement Learning** Subset of machine learning where an agent learns to make decisions by taking actions in an environment to maximise cumulative reward. xii, 2, 4, 6, 10, 16, 17

**Reward Function** Function that defines the feedback signal received by an agent in reinforcement learning, guiding the agent's learning process by providing rewards or penalties based on its actions. 2

**SHapley Additive exPlanations** Method for interpreting machine learning models by assigning each feature an importance value for a particular prediction, based on cooperative game theory. xii, 4, 6

**Soft Actor-Critic** Algorithm that combines the benefits of off-policy learning and entropy regularisation, allowing for more exploration and better stability in learning policies for continuous action spaces. xii, 3, 16, 21

**Supervised Learning** Machine learning task where a model is trained on labelled data to learn a mapping from inputs to outputs. 10

**Twin Delayed Deep Deterministic Policy Gradient** Extension of Deep Deterministic Policy Gradient that addresses the overestimation bias in value function estimation by using two critic networks and delaying policy updates. xii, 3, 16, 20

**Unsupervised Learning** Machine learning task where a model learns patterns or structures in unlabelled data without explicit supervision. 10

# Acronyms

**A2C** Advantage Actor-Critic. 3, 4, 16–18, 22

**A3C** Asynchronous Advantage Actor-Critic. 16, 18

**AI** Artificial Intelligence. 2, 9

**DDPG** Deep Deterministic Policy Gradient. 3, 4, 16, 19, 20, 22

**DL** Deep Learning. 2, 4, 6, 10, 17

**DNN** Deep neural networks. 10, 17

**DPG** Deterministic Policy Gradient. 19

**DQN** Deep Q-Network. 16, 19, 20

**DRL** Deep Reinforcement Learning. 3–6, 10, 17, 22–24

**LIME** Local Interpretable Model-agnostic Explanations. 4–6

**MDP** Markov Decision Process. 11–14, 23

**ML** Machine Learning. 1–6, 9, 10

**MPT** Modern Portfolio Theory. 7, 9

**POMDP** Partially Observable Markov Decision Process. 13

**PPO** Proximal Policy Optimisation. 3, 4, 16, 18, 19, 22

**RL** Reinforcement Learning. 2, 4, 6, 10, 11, 15–17, 21, 22

**SAC** Soft Actor-Critic. 3, 4, 16, 21, 22

**SHAP** SHapley Additive exPlanations. 4–6

**TD3** Twin Delayed Deep Deterministic Policy Gradient. 3, 4, 16, 20–22

**TRPO** Trust Region Policy Optimisation. 18, 19

**XAI** Explainable Artificial Intelligence. 2

# Chapter 1

## Introduction

Financial markets are highly complex systems influenced by numerous factors, including financial and political events, social trends and technological advancements. Moreover, their evolving and stochastic nature requires using the most advance computational developments to model the financial environment. The tasks of financial time series prediction and Portfolio optimisation are considerably intricate, due to the semi-strong form of market efficiency and the high level of noise. [3]

Algorithmic trading focuses on the application of analytical methods to automatically execute trading actions based on an algorithm without human intervention. In its early days, the field mainly studied the usage of a computer program to follow a predefined strategy [4]. Nonetheless, in recent years, algorithmic trading has evolved to a problem in which environment perception entails learning feature representation from highly non-stationary and noisy financial time series data and decision-making requires the algorithm to explore the environment and simultaneously make correct decisions in an online manner without supervision [5].

Machine Learning (ML) is at an advantage for the task given its capability to learn from



historical data and make predictions about the future state of an environment. In the past years, research has explored the application of Deep Learning (DL) in future price prediction of financial assets [3, 6, 7, 8]. However, its main disadvantage is the inability to directly deal with trading, requiring an additional step to convert the predictions into actionable strategies. In contrast, Reinforcement Learning (RL) would allow the algorithm to learn a trading strategy directly from the environment, without the need for a separate step [9, 10]. In this case, there are two main approaches: first, the algorithm can learn the amount of assets to buy, sell or hold at each time step [11], or second, the algorithm can learn the optimal portfolio allocation and automatically rebalance the portfolio weights at each time step [12].

Despite the potential of RL in portfolio optimisation, its widespread adoption in the financial industry remains limited. This is primarily due to following challenges [13]:

1. difficulty in finding the appropriate algorithm with a suitable reward function and hyper-parameters to ensure efficiency and performance,
2. challenge of testing the algorithm in a real-world environment, and
3. lack of transparency of ML models, often referred to as black boxes, making it increasingly complex to interpret the algorithm's decisions.

In recent years, the rise in popularity of Artificial Intelligence (AI) and its widespread use have led to concerns regarding its decisions due to its black-box nature. The concept of explainability in AI, known as Explainable Artificial Intelligence (XAI), refers to a model's ability to provide details and reasons to make itself understandable [14]. The term was first coined in 2016 to describe the need for users to effectively understand, trust and manage artificial intelligence applications [15]. The need for explainability becomes particularly relevant in the context of financial usage, where the regulatory framework requires transparency and accountability in automated decision-making. Various relevant

applications, including volatility models [16], credit risk assessment [17] and portfolio construction [13] have explored the concept of explainability in financial applications. This highlights its importance and the need to explore its advantages in more complex models, without inadvertently increasing the complexity of the overall methodology.

Consequently, this thesis will focus on addressing the aforementioned challenges by exploring the application of Deep Reinforcement Learning (DRL) to portfolio optimisation and implementing post-hoc explainability techniques.

## 1.1 Objectives

The objective of this thesis is to develop an explainable Deep Reinforcement Learning model for portfolio optimisation. A DRL model has the ability to leverage historical financial data to learn an investment strategy that efficiently allocates financial assets while maximising expected returns and minimising risk. Moreover, the incorporation of advanced explainability techniques enhances the interpretability and transparency of the model’s decision-making. This project aims to bridge the gap between cutting-edge machine learning techniques and their practical application in finance by addressing the challenges of algorithm selection, simulation of real-world scenarios, and black box nature of ML models.

First, DRL models, such as Advantage Actor-Critic (A2C), Proximal Policy Optimisation (PPO), Deep Deterministic Policy Gradient (DDPG), Twin Delayed Deep Deterministic Policy Gradient (TD3) and Soft Actor-Critic (SAC), will be implemented to learn the optimal portfolio allocation from high-dimensional environment representations. The algorithms will be trained on historical financial data, including technical and macroeconomic indicators, with the goal of capturing the complex market dynamics. Each of the algorithms is better suited to a particular scenario, for instance, DDPG

encourages maximum returns, while A2C reduces the variance.

Second, post-hoc explainability techniques: Feature Importance, SHapley Additive explanations (SHAP) values and Local Interpretable Model-agnostic Explanations (LIME) analysis, will be applied to interpret the model’s decisions. The goal of these is to understand which market conditions, represented by financial data, lead to the actions/decisions, encoded as portfolio weights.

Finally, the performance of the DRL models will be analysed in different scenarios, including the impact of a larger financial environment representation, portfolio size and asset composition. The results will be compared with traditional portfolio optimisation methods to evaluate the effectiveness of the proposed approach.

## 1.2 Report Structure

This report is organised into six chapters, each of which focuses on a concrete area related to the problem. Additional material, together with source code, is included in the appendices.

The current chapter, 1, presents the motivation and the objectives of this thesis. It gives an overview of the potential of DRL in portfolio optimisation and its main challenges, particularly the lack of transparency of ML models.

Chapter 2 provides an overview of the theoretical background of the project, including financial markets and machine learning concepts. The problem of portfolio optimisation in the financial domain is outlined and the potential of DRL in this context is discussed. The chapter provides a comprehensive background explanation of the fundamentals of Deep Learning and Reinforcement Learning, including the main algorithms (A2C, PPO, DDPG, TD3, SAC) and techniques in their intersection, DRL. In addition, it gives an overview of the post-hoc explainability techniques: feature importance, SHapley Addi-

tive exPlanations (SHAP) values and Local Interpretable Model-agnostic Explanations (LIME) analysis, which will be used to interpret the model's decisions. Finally, it includes a comprehensive in-depth literature review on the topics of ML applied to portfolio optimisation and relevant applications of explainability techniques.

The methodology chapter 3 describes the techniques and methods used to solve the problem and outlines the implementation of the proposed solution. The chapter provides a detailed explanation of the architecture and components of the proposed DRL model, including the state representation, reward function, and training process.

The results of the experiments are presented in chapter 4, which analyses and evaluates the results obtained from the proposed implementation, while critically discussing the findings. It provides a detailed comparison of the proposed DRL strategies with traditional portfolio optimisation methods. Furthermore, it consists of an in-depth analysis of the model's decisions using post-hoc explainability techniques, in particular, SHAP values, feature importance and LIME analysis.

Chapter 5 discusses the legal, social, ethical and professional implications within the context of the project. By addressing these issues, the project aims to ensure that the proposed solution adheres to industry standards, while considering the implications of the technology.

Finally, the report concludes with a summary of the main points of the work, the contributions made, the results achieved as well as potential applications and future work in chapter 6.

## Chapter 2

# Background

This chapter provides an overview of the problem of portfolio optimisation in the financial domain, followed by a comprehensive explanation of the fundamentals of Deep Learning (DL) and Reinforcement Learning (RL), including the relevant algorithms in the field of Deep Reinforcement Learning (DRL). In addition, it discusses the need for explainability in Machine Learning (ML) and the main techniques used to achieve it: SHapley Additive exPlanations (SHAP), Local Interpretable Model-agnostic Explanations (LIME) and Feature Importance. Finally, it presents the state of the art in portfolio optimisation using DRL and the recent advancements in explainability techniques in the field.

### 2.1 Portfolio Optimisation

Portfolio optimisation is the process of selecting optimal weights for a portfolio of assets in order to maximise expected returns for a given level of risk, or conversely, to minimise risk for a given level of expected returns [18]. In mathematical terms, the problem requires finding a solution to the specified objective function, which is typically a function of the expected returns and the risk associated with the portfolio [19]. The task becomes

further complicated if a time dimension is introduced, as the portfolio weights need to be adjusted over time to capture the changes in market conditions and asset prices [20].

### 2.1.1 Modern Portfolio Theory

There exist several traditional frameworks that formalise the problem of portfolio allocation. Markowitz's Modern Portfolio Theory (MPT) was proposed in 1952 [21] and it provides a mathematical framework where investors choose optimal portfolios based on risk and return, by either minimising the risk given a specified return or, maximising the return given a specified risk [22]. The theory extends the concept of diversification by suggesting that owning financial assets of different kinds is less risky than owning assets of the same kind, due to the correlations between assets.

The main assumptions in MPT are:

- investors are risk-averse, rational, and seek to maximise return for a given risk;
- returns are normally distributed;
- markets are frictionless, meaning there are no transaction costs; and
- assets are infinitely divisible.

Under these assumptions, portfolio risk and return can be modelled as an optimisation problem. Let  $\mathbf{w} = (w_1, w_2, \dots, w_N)^T$  denote the portfolio weight vector, where each  $w_i$  indicates the proportion of capital allocated to asset  $i$ , subject to the budget constraint:

$$\sum_{i=1}^N w_i = 1 \quad \Leftrightarrow \quad \mathbf{w}^T \mathbf{1} = 1 \quad (2.1)$$

with  $\mathbf{1} \in \mathbb{R}^N$  being a vector of ones, and subject to the non-negativity constraint,

meaning that short-selling is not allowed:

$$w_i \geq 0 \quad \forall i = 1, 2, \dots, N. \quad (2.2)$$

Let  $\boldsymbol{\mu} = (R_1, R_2, \dots, R_N)^T$  represent the vector of expected returns, and  $\Sigma \in \mathbb{R}^{N \times N}$  the covariance matrix of asset returns. The expected return of the portfolio is then given by:

$$R_p = \mathbf{w}^T \boldsymbol{\mu}, \quad (2.3)$$

and the portfolio risk is quantified by the variance of returns:

$$\sigma_p^2 = \mathbf{w}^T \Sigma \mathbf{w}. \quad (2.4)$$

This formulation provides the foundation for solving the mean-variance optimisation problem, by either:

- minimising portfolio variance  $\sigma_p^2$  subject to a target expected return  $R_p$ , or
- maximising expected return  $R_p$  subject to a risk constraint  $\sigma_p$ .

The Markowitz mean-variance optimisation problem can be expressed as:

$$\begin{aligned} \min_{\mathbf{w}} \quad & \mathbf{w}^T \Sigma \mathbf{w} \\ \text{subject to} \quad & \begin{cases} \mathbf{w}^T \boldsymbol{\mu} = R_p \\ \mathbf{w}^T \mathbf{1} = 1 \\ \mathbf{w} \geq 0 \end{cases} \end{aligned} \quad (2.5)$$

Solving the mean-variance optimisation problem for varying levels of target return leads to a set of optimal portfolios that form the efficient frontier. It is typically visualised

in a risk-return space, where the x-axis represents the risk (standard deviation) and the y-axis represents the expected return, as shown in Figure 2.1. Portfolios below the curve are suboptimal, while those on the frontier represent the best achievable combinations of risk and return.

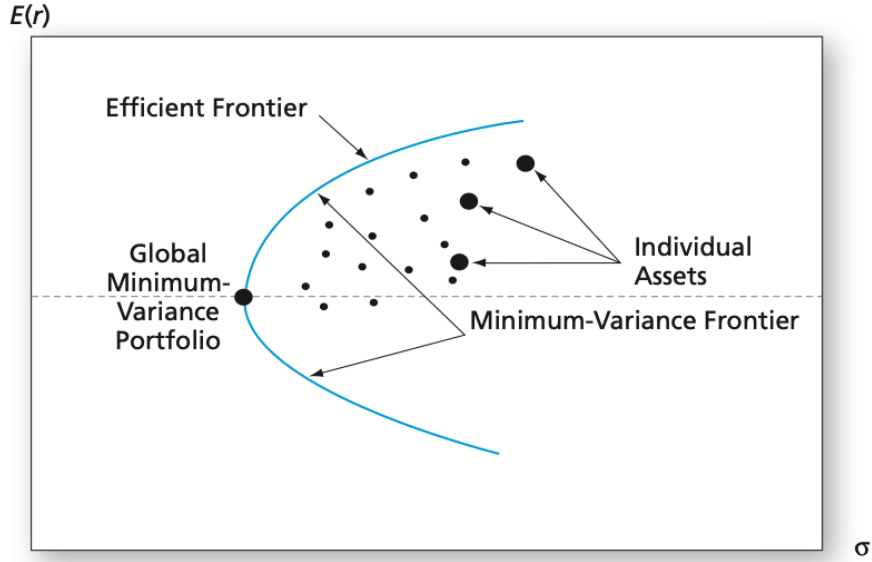


Figure 2.1: Efficient Frontier in Risk-Return Space. [1]

Despite the simplicity in the formulation of MPT, its assumptions do not reflect the behaviour of real markets. Moreover, modern markets are dynamic, non-stationary, and feature non-linear relationships, which have driven research into other approaches better suited to capture the complexities of modern financial markets.

## 2.2 Deep Reinforcement Learning

Machine Learning is a branch of Artificial Intelligence (AI) that focuses on the use of data and algorithms to imitate the way humans learn, gradually improving their accuracy over time [23]. There are three main tasks in ML [24]:



- **Supervised Learning:** Task of training a classification or regression model from labelled training data, where the model learns to map inputs to outputs based on examples.
- **Unsupervised Learning:** Task of drawing inferences from datasets consisting of unlabelled input data, where the model learns to identify patterns or structures in the data.
- **Reinforcement Learning:** Task of training an agent to sequentially make decisions by taking actions in an environment with the goal of maximising cumulative reward, using feedback from the environment to learn an optimal strategy.

Deep Learning is a set of methods and techniques to solve such ML tasks, specially in supervised and unsupervised learning tasks. DL focuses on the use of Deep neural networks (DNN) [25], which are characterised by a succession of layers of non-linear transformation that allow the model to learn a representation of the data with various levels of abstraction.

Therefore, Deep Reinforcement Learning (DRL) combines Deep Learning (DL) and Reinforcement Learning (RL) to solve sequential decision-making problems with high-dimensionality in the environment representation. This approach has gained significant attention in recent years due to its success in various applications, including robotics [26] and game playing [27, 28].

### 2.2.1 Reinforcement Learning

As mentioned, RL is a type of ML that solves the problem of sequential decision-making through continuous interaction with an environment. The agent learns to take actions given a representation of the environment’s state with the goal of optimising a pre-defined notion of reward. The agent learns by successively adjusting its policy based on

its observations and interactions with the environment.

The RL problem can be formalised as a discrete-time stochastic control process where an agent interacts with the environment. At each time step  $t$ , the agent observes the state of the environment  $s_t \in \mathcal{S}$ , takes an action  $a_t \in \mathcal{A}$  to obtain a reward  $r_t \in \mathbb{R}$  and transition to a new state  $s_{t+1} \in \mathcal{S}$ , where  $\mathcal{S}$  is the state space and  $\mathcal{A}$  is the action space [24]. The agent's interaction with the environment is visually represented in Figure 2.2.

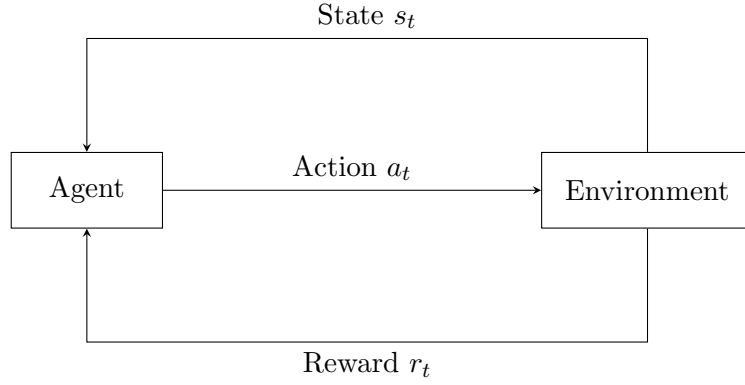


Figure 2.2: Agent interaction with environment

A discrete time stochastic control process can be formalised as a Markov Decision Process (MDP), if it fulfils the Markov Property.

**Definition 2.2.1** (Markov Property). A discrete time stochastic control process satisfies the Markov Property if:

$$P(s_{t+1}|s_t, a_t) = P(s_{t+1}|s_t, a_t, s_{t-1}, a_{t-1}, \dots, s_0, a_0) \quad (2.6)$$

$$P(r_t|s_t, a_t) = P(r_t|s_t, a_t, s_{t-1}, a_{t-1}, \dots, s_0, a_0) \quad (2.7)$$

where  $P(s_{t+1}|s_t, a_t)$  is the transition probability of moving to state  $s_{t+1}$  given the current state  $s_t$  and action  $a_t$ , and  $P(r_t|s_t, a_t)$  is the reward function that defines the expected reward received at time  $t$  given the current state and action.

This implies that the state  $s_{t+1}$  at a future time step  $t + 1$  only depends on the current state  $s_t$  and action  $a_t$ . Similarly, the reward  $r_t$  at time step  $t$  only depends on the current state and action and not on the history of previous states and actions. Consequently, a Markov Decision Process [29] is a discrete time stochastic control process defined as:

**Definition 2.2.2** (Markov Decision Process). An MDP is a tuple  $\mathcal{M} = (\mathcal{S}, \mathcal{A}, T, R, \gamma)$ , where:

- $\mathcal{S}$  is the state space:  $s_t \in \mathcal{S}$ ,
- $\mathcal{A}$  is the action space:  $a_t \in \mathcal{A}$ ,
- $T : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow [0, 1]$  is the transition function,
- $R : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow \mathcal{R}$  is the reward function, where  $\mathcal{R} \in [0, R_{max}]$  is the set of all possible rewards bounded by  $R_{max} \in \mathbb{R}^+$ , and
- $\gamma \in [0, 1)$  is the discount factor.

At each time step, the probability of advancing to the next state  $s_{t+1}$  is given by the transition function  $T(s_t, a_t, s_{t+1})$  and the reward  $r_t$  is given by the reward function  $R(s_t, a_t, s_{t+1})$ . This can be visualised in Figure 2.3.

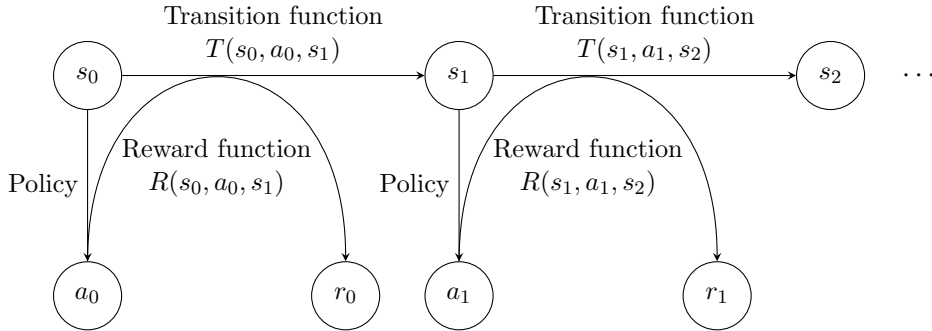


Figure 2.3: Markov Decision Process with policy, transition, and reward functions

The agent's objective is to learn a policy  $\pi : \mathcal{S} \rightarrow \mathcal{A}$  that maps states to actions, in order to maximise the expected cumulative reward over time. Policies can be categorised as:

- deterministic:  $\pi(s) : \mathcal{S} \rightarrow \mathcal{A}$ , at a given state  $s$ , the policy specifies the only available action to take, or
- stochastic:  $\pi(s, a) : \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$ , at a given state  $s$ , the policy specifies the probability of taking action  $a$ .

Markov Decision Process are based on the idea that the current state is fully representative of the environment. However, in most real world scenarios, the agent does not have access to the complete state. In such cases, Partially Observable Markov Decision Process (POMDP) can be used to model the uncertainty in the agent's observations and actions.

The goal of the agent is to maximise the cumulative long-term reward  $G_t$ , which is defined as the sum of discounted rewards over time:

$$G_t = r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \dots = r_{t+1} + \gamma G_{t+1} \quad (2.8)$$

where  $\gamma \in [0, 1)$  is the discount factor and is used to balance the importance between immediate and future rewards. If the discount factor is set to 0, the agent is myopic and only maximises the immediate reward; whereas, as  $\gamma$  approaches 1, the agent becomes more far-sighted and places greater importance on future rewards.

The expected cumulative reward is defined as the state value function  $V^\pi(s) : \mathcal{S} \rightarrow \mathbb{R}$ , which is the expected return when starting from state  $s$  and following policy  $\pi$ :

$$V^\pi(s) = \mathbb{E}_\pi [G_t | s_t = s] \quad (2.9)$$

where  $\mathbb{E}_\pi$  denotes the expectation over the policy  $\pi$ .

Similarly, the state-action value function  $Q^\pi(s, a) : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$  is defined as the expected

return when starting from state  $s$ , taking action  $a$ , and then following policy  $\pi$ :

$$Q^\pi(s, a) = \mathbb{E}_\pi [G_t | s_t = s, a_t = a] \quad (2.10)$$

The state value function  $V^\pi(s)$  and the state-action value function  $Q^\pi(s, a)$  are related as follows:

$$V^\pi(s) = \sum_{a \in \mathcal{A}} \pi(a|s) Q^\pi(s, a) = \mathbb{E}_\pi [Q^\pi(s, a) | s_t = s] \quad (2.11)$$

Moreover, the advantage function  $A^\pi(s, a)$  combines both the state value function  $V^\pi(s)$  and the state-action value function  $Q^\pi(s, a)$ , and is defined as:

$$A^\pi(s, a) = Q^\pi(s, a) - V^\pi(s) \quad (2.12)$$

A policy  $\pi$  is said to be optimal if the policy's value function is the optimal value function of the MDP, defined as:

$$V^*(s) = \max_{\pi'} V^{\pi'}(s), \forall s \in \mathcal{S} \quad (2.13)$$

$$Q^*(s, a) = \max_{\pi'} Q^{\pi'}(s, a), \forall s \in \mathcal{S}, a \in \mathcal{A} \quad (2.14)$$

The optimal policy  $\pi^*$  is:

$$\pi^*(s) = \arg \max_{a \in \mathcal{A}} Q^*(s, a) \quad (2.15)$$

As a result, the optimal policy is the greedy policy that performs the optimal actions at each time step as determined by the optimal value functions. This framework enables the agent to determine optimal actions that maximise long-term returns by evaluating immediate information, without requiring knowledge of the values of future states and actions.

The Bellman equations [30] provide a recursive relation between the value functions in terms of the future state/action values. There are four main Bellman equations, classified in two groups: the Bellman expectation equations and the Bellman optimality equations. The Bellman expectation equations are defined as follows:

$$V^\pi(s) = \sum_{a \in \mathcal{A}} \pi(a | s) \sum_{s' \in \mathcal{S}} T(s, a, s') [R(s, a) + \gamma V^\pi(s')] \quad (2.16)$$

$$Q^\pi(s, a) = \sum_{s' \in \mathcal{S}} T(s, a, s') \left[ R(s, a) + \gamma \sum_{a' \in \mathcal{A}} \pi(a' | s') Q^\pi(s', a') \right] \quad (2.17)$$

and the Bellman optimality equations are defined as:

$$V^*(s) = \max_{a \in \mathcal{A}} \sum_{s' \in \mathcal{S}} T(s, a, s') [R(s, a) + \gamma V^*(s')] \quad (2.18)$$

$$Q^*(s, a) = \sum_{s' \in \mathcal{S}} T(s, a, s') \left[ R(s, a) + \gamma \max_{a' \in \mathcal{A}} Q^*(s', a') \right] \quad (2.19)$$

Although explicitly solving the Bellman equations would lead to the optimal policy, it is often intractable due to the size of the state and action spaces. Therefore, in RL algorithms, the goal is to learn an approximation of the optimal value functions, which can be used to derive the optimal policy. Another problem that arises is that of balancing exploration and exploitation [31]. Theoretically, following the greedy action yields the optimal policy, but this is only true if all the action values are known. In practice, the agent at each time step and given state chooses either an action whose value is higher, thus exploiting its current knowledge, or picks an action at random, thus exploring the environment and gaining more information about the state-action space, leading to potentially discovering a better action than the greedy one at the current time.

### 2.2.2 Deep Reinforcement Learning Algorithms

A Reinforcement Learning agent includes one or more of the following components [24]:

- a representation of the value function that provides a prediction of the value of each state or state-action pair,
- a direct representation of the policy  $\pi$ , and
- a model of the environment, consisting of estimates of transition and reward functions.

Depending on the components, the main Reinforcement Learning paradigms are:

- **Model-free** algorithms do not learn a representation of the environment, but focus on either the value function or the policy.
  - **Value-based** algorithms learn an approximation of the value function, which is used to compute the state or state-action values. The policy is not learnt explicitly but can be derived from the value function. Examples include Deep Q-Network (DQN) [32] and C51 [33].
  - **Policy-based** algorithms learn a direct representation of the policy, which is used to select actions. Advantage Actor-Critic (A2C), Asynchronous Advantage Actor-Critic (A3C) [34] and Proximal Policy Optimisation (PPO) [35] are examples of policy-based algorithms.
  - **Actor-Critic** algorithms combine both value-based and policy-based approaches, where the actor learns the policy and the critic learns the value function. The critic provides feedback to the actor to improve the policy. Some examples are Deep Deterministic Policy Gradient (DDPG) [36], Twin Delayed Deep Deterministic Policy Gradient (TD3) [37], and Soft Actor-Critic (SAC) [38].

- **Model-based** algorithms include a model of the environment, which can be used to simulate future states and rewards. For example, World Models [39] learns a model of the environment and AlphaZero [40] has a representation of the model.

The taxonomy of the RL algorithms is shown in Figure 2.4.

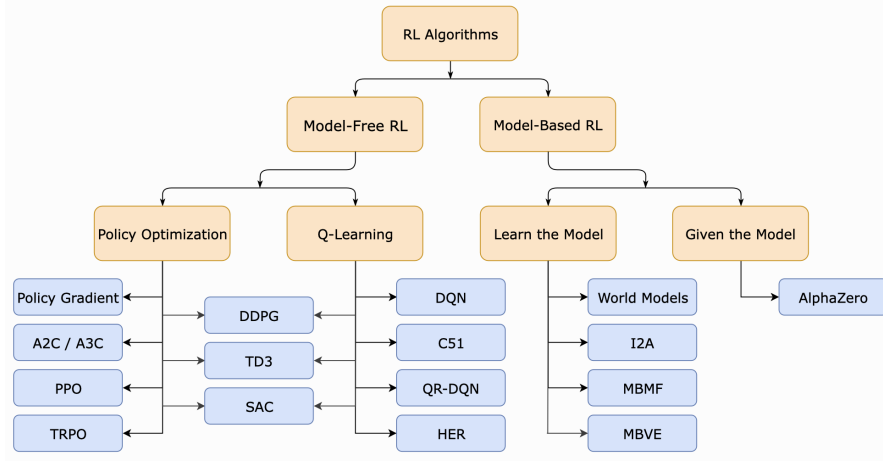


Figure 2.4: Taxonomy of Reinforcement Learning Algorithms [2]

The full potential of Reinforcement Learning algorithms is achieved when leveraging the power of Deep Learning to solve dynamic stochastic control problems that have high-dimensionality in their representation of the state and action spaces. DRL algorithms use DNNs to approximate the value functions, or use gradient ascent to find the optimal policy parameters. This thesis will focus on the following DRL algorithms on policy-based and actor-critic algorithms.

#### 2.2.2.1 Advantage Actor-Critic (A2C)

The Advantage Actor-Critic (A2C) algorithm was developed by Mnih et al. (2016) in their paper on *Asynchronous methods for deep reinforcement learning* [34]. The main contribution is the usage of the advantage function to address the variance issues present



in policy gradient methods. The A2C is the synchronous version of A3C, and is preferred due to its better performance in terms of training time and cost-effectiveness [41].

The algorithm consists of a dual-network architecture, where the actor network learns a stochastic policy  $\pi(a_t | s_t; \theta)$  and the critic network learns the value function  $V(s_t; \theta_v)$ . The policy and the value function are updated after every  $t_{max}$  actions or when a terminal state is reached. The advantage function is estimated as follows:

$$A(s_t, a_t; \theta, \theta_v) = \sum_{i=0}^{k-1} \gamma^i r_{t+i} + \gamma^k V(s_{t+k}; \theta_v) - V(s_t; \theta_v) \quad (2.20)$$

where  $k$  represents the  $n$ -step return and is upper-bounded by the maximum number of steps  $t_{max}$ , and  $\gamma$  is the discount factor. The algorithm's objective function is defined as:

$$J(\theta, \theta_v) = \mathbb{E}_{\pi} [\log \pi(a_t | s_t; \theta) A(s_t, a_t; \theta, \theta_v)] \quad (2.21)$$

The pseudo-code for the algorithm for A2C is outlined in appendix A.

#### 2.2.2.2 Proximal Policy Optimisation (PPO)

Proximal Policy Optimisation (PPO) is a policy gradient algorithm that was introduced by Schulman et al. (2017) in their paper on *Proximal Policy Optimization Algorithms* [35] with the objective of constraining policy updates. The algorithm balances sufficiently large policy updates in order to improve the policy, while avoiding excessively large changes that could hinder performance.

PPO improves the performance of Trust Region Policy Optimisation (TRPO) [42] by using a clipped surrogate objective function, by ensuring that the probability ratio  $r_t$  is bounded within a range of  $[1 - \epsilon, 1 + \epsilon]$ , where  $\epsilon$  is a hyperparameter that controls the

clipping range. The objective function is defined as:

$$L_t^{CLIP}(\theta) = \mathbb{E}_t \left[ \min \left( r_t(\theta) \hat{A}_t, \text{clip} \left( r_t(\theta), 1 - \epsilon, 1 + \epsilon \right) \hat{A}_t \right) \right] \quad (2.22)$$

where

$$r_t(\theta) = \frac{\pi_\theta(a_t|s_t)}{\pi_{\theta_{old}}(a_t|s_t)} \quad (2.23)$$

is the probability ratio between the current policy  $\pi_\theta$  and the old policy  $\pi_{\theta_{old}}$

Another improvement with respect to TRPO is the use of simple first-order optimisation methods as opposed to the second-order methods used in TRPO. Consequently, PPO maintains the stability and reliability of other trust-region methods, while being easier to implement and more computationally efficient. The pseudo-code for the algorithm for PPO is outlined in appendix A.

### 2.2.2.3 Deep Deterministic Policy Gradient (DDPG)

Deep Deterministic Policy Gradient (DDPG) is an off-policy actor-critic algorithm that was introduced by Lillicrap et al. (2015) in their paper on *Continuous control with deep reinforcement learning* [36]. The algorithm arises to solve the challenges of applying Deep Q-Network [32] to continuous action spaces by applying Deterministic Policy Gradient (DPG), which enables the efficient computation of the policy gradient without the need to integrate over the action space.

The algorithm uses the following networks: the actor network  $\mu(s_t; \theta_\mu)$ , which learns a deterministic policy, the critic network  $Q(s_t, a_t; \theta_Q)$ , which learns the state-action value function, and the actor's and critic's target networks. The actor network is updated using DPG:

$$\nabla_{\theta_\mu} J \approx \mathbb{E}_{s_t \sim \mathcal{D}} \left[ \nabla_a Q(s_t, a_t; \theta_Q) \nabla_{\theta_\mu} \mu(s_t; \theta_\mu) \right] \quad (2.24)$$

where  $\mathcal{D}$  is the replay buffer that stores the agent’s experiences, and the critic network is updated using the Bellman equations:

$$\nabla_{\theta_Q} J \approx \mathbb{E}_{s_t \sim \mathcal{D}} \left[ (r_t + \gamma Q(s_{t+1}, \mu(s_{t+1}; \theta_\mu); \theta'_Q) - Q(s_t, a_t; \theta_Q)) \nabla_{\theta_Q} Q(s_t, a_t; \theta_Q) \right] \quad (2.25)$$

where  $\theta'_Q$  are the parameters of the target critic network.

From Deep Q-Network, the algorithm incorporates a replay buffer to store the agent’s experiences, which allows the agent to learn from past experiences and improve its performance over time. Moreover, DDPG incorporates noise, typically Ornstein-Uhlenbeck noise [43], to the actions taken by the actor network to encourage exploration of the action space. The pseudo-code for the algorithm for DDPG is outlined in appendix A.

#### 2.2.2.4 Twin Delayed Deep Deterministic Policy Gradient (TD3)

Twin Delayed Deep Deterministic Policy Gradient (TD3) is an extension of DDPG introduced by Fujimoto et al. (2018) in their paper on *Addressing Function Approximation Error in Actor-Critic Methods* [37]. The proposal addresses the hyper-parameter sensitivity and overestimation bias present in the critic network of DDPG. There are three main improvements to DDPG:

- **Clipped Double Q-Learning:** The algorithm employs two critic networks and uses their minimum value to compute the target value for the actor network, which reduces the overestimation bias.
- **Delayed Policy Updates:** The actor and critic networks updates are performed at different frequencies. The critic networks are updated every time step, whereas the actor and target networks are updated less frequently. The main benefit is that it allows the critic to improve the accuracy of the value estimates before the actor updates its policy.

- **Target Policy Smoothing:** The algorithm adds noise to the target action during the critic updates to smooth the value function over similar actions, resulting in a lower impact of approximation errors and more robust policies.

The pseudo-code for the algorithm for TD3 is outlined in appendix A.

#### 2.2.2.5 Soft Actor-Critic (SAC)

Soft Actor-Critic (SAC), despite being an off-policy actor-critic algorithm, represents a paradigm shift in RL. The algorithm was presented by Haarnoja et al. (2018) in their paper on *Soft Actor-Critic: Off-Policy Maximum Entropy Deep Reinforcement Learning with a Stochastic Actor* [38] and is based on the maximum entropy framework, which aims to maximise both the expected return and the entropy of the policy, encouraging the agent to succeed at the task while acting as randomly as possible.

The algorithm uses two critic networks to estimate the state-action value function, and a stochastic actor network that learns a policy that maximises the expected return while also maximising the entropy of the policy. The critic networks are updated using the Bellman equations, and the actor network is updated using the policy gradient method. The objective function for the actor network is defined as:

$$J(\theta) = \mathbb{E}_{s_t \sim \mathcal{D}} [\mathbb{E}_{a_t \sim \pi_\theta} [\alpha \log \pi_\theta(a_t|s_t) + Q(s_t, a_t; \theta_Q)]] \quad (2.26)$$

where  $\alpha$  is a temperature parameter that controls the trade-off between exploration and exploitation, and  $\mathcal{D}$  is the replay buffer that stores the agent's experiences. The critic networks are updated using the Bellman equations:

$$J(\theta_Q) = \mathbb{E}_{s_t \sim \mathcal{D}} \left[ \left( r_t + \gamma \min_{i=1,2} Q(s_{t+1}, \mu_\theta(s_{t+1}) + \mathcal{N}(0, \sigma); \theta_{Q_i}) - Q(s_t, a_t; \theta_{Q_i}) \right)^2 \right] \quad (2.27)$$

where  $\mathcal{N}(0, \sigma)$  is the noise added to the action to encourage exploration, and  $\theta_{Q_i}$  are the parameters of the two critic networks. The temperature parameter  $\alpha$  is learned automatically by maximising the entropy of the policy, which is defined as:

$$J(\alpha) = \mathbb{E}_{s_t \sim \mathcal{D}} [\alpha \log \pi_{\theta}(a_t | s_t)] \quad (2.28)$$

This parameter is used to control the trade-off between the entropy and the reward terms, effectively controlling the stochasticity of the policy.

The pseudo-code for the algorithm for SAC is outlined in appendix A.

#### 2.2.2.6 Algorithm comparison

The five algorithms presented above are among the most widely used in the field of DRL, each addressing specific challenges in policy estimation and value function approximation. A2C is an on-policy actor-critic method that reduces variance through advantage estimation, while PPO changed on-policy learning by introducing a clipped surrogate objective that ensures stable policy updates. Although DDPG pioneered continuous control through deterministic policies, TD3 addressed the overestimation bias of its predecessor via twin critics and delayed updates. SAC revolutionised the field of DRL by incorporating maximum entropy principles.

An overview of the algorithms and their key characteristics is summarised in Table 2.1.

Algorithm	Policy Type	On/Off Policy	Key Idea
A2C	Stochastic	On-Policy	Advantage Actor-Critic
PPO	Stochastic	On-Policy	Clipped Policy Updates
DDPG	Deterministic	Off-Policy	Actor-Critic + Replay buffer
TD3	Deterministic	Off-Policy	Double critics + Delayed updates
SAC	Stochastic	Off-Policy	Max-entropy RL

Table 2.1: Comparison of Deep Reinforcement Learning Algorithms

## Chapter 3

# Methodology

This chapter covers the methodology and framework established to provide an explainable Deep Reinforcement Learning (DRL) model capable of optimising a portfolio. The chapter is structured as follows: first, it describes the architecture and components of the proposed DRL model, including the state representation, reward function and training process. Second, it discusses the evaluation metrics and experimental setup used to assess the performance of the proposed solution. Finally, it outlines the implementation of the post-hoc explainability techniques used to interpret the model’s decisions.

### 3.1 Problem Definition

The problem of portfolio optimisation is the task of finding an optimal allocation of financial assets in a portfolio to maximise expected returns while minimising risk. Thus, it is necessary to decide how to rebalance the portfolio at each time step in a highly stochastic and complex financial market. This can be formulated using a Markov Decision Process (MDP) framework, where the agent interacts with the environment by deciding the optimal allocation based on the state of the environment at each time step

to maximise the expected cumulative reward over time. Deep Reinforcement Learning (DRL) gives the agent the ability to learn the optimal policy directly from the environment by taking actions and receiving rewards.

## Chapter 4

# Results



## Chapter 5

# Legal, Social, Ethical and Professional Issues

## Chapter 6

## Conclusion

# References

- [1] Z. Bodie, A. Kane, and A. J. Marcus, *Investments*. McGraw-Hill Education, 10 ed., 2014.
- [2] J. Achiam, “Part 2: Kinds of rl algorithms,” 2018.
- [3] J. Shen and M. O. Shafiq, “Short-term stock market price trend prediction using a comprehensive deep learning system,” *Journal of Big Data*, vol. 7, pp. 1–33, 12 2020.
- [4] K. Lei, B. Zhang, Y. Li, M. Yang, and Y. Shen, “Time-driven feature-aware jointly deep reinforcement learning for financial signal representation and algorithmic trading,” *Expert Systems with Applications*, vol. 140, p. 112872, 2 2020.
- [5] C. Ma, J. Zhang, J. Liu, L. Ji, and F. Gao, “A parallel multi-module deep reinforcement learning algorithm for stock trading,” *Neurocomputing*, vol. 449, pp. 290–302, 8 2021.
- [6] I. K. Nti, A. F. Adekoya, and B. A. Weyori, “A comprehensive evaluation of ensemble learning for stock-market prediction,” *Journal of Big Data*, vol. 7, pp. 1–40, 12 2020.

- [7] J. M. T. Wu, Z. Li, N. Herencsar, B. Vo, and J. C. W. Lin, “A graph-based cnn-lstm stock price prediction algorithm with leading indicators,” *Multimedia Systems*, vol. 29, pp. 1751–1770, 6 2023.
- [8] M. Hasan, M. Z. Abedin, P. Hajek, K. Coussement, M. N. Sultan, and B. Lucey, “A blending ensemble learning model for crude oil price forecasting,” *Annals of Operations Research*, pp. 1–31, 1 2024.
- [9] J. Moody and M. Saffell, “Learning to trade via direct reinforcement,” *IEEE Transactions on Neural Networks*, vol. 12, pp. 875–889, 7 2001.
- [10] H. Yang, X. Y. Liu, S. Zhong, and A. Walid, “Deep reinforcement learning for automated stock trading: An ensemble strategy,” *ICAIF 2020 - 1st ACM International Conference on AI in Finance*, 10 2020.
- [11] X.-Y. Liu, Z. Xiong, S. Zhong, H. Yang, and A. Walid, “Practical deep reinforcement learning approach for stock trading,” *NeurIPS 2018 AI in Finance Workshop.*, 11 2018.
- [12] M. Guan and X. Y. Liu, “Explainable deep reinforcement learning for portfolio management: An empirical approach,” *ICAIF 2021 - 2nd ACM International Conference on AI in Finance*, 11 2021.
- [13] D. G. Cortés, E. Onieva, I. Pastor, L. Trinchera, and J. Wu, “Portfolio construction using explainable reinforcement learning,” *Expert Systems*, vol. 41, p. e13667, 11 2024.
- [14] A. B. Arrieta, N. Díaz-Rodríguez, J. D. Ser, A. Bennetot, S. Tabik, A. Barbado, S. Garcia, S. Gil-Lopez, D. Molina, R. Benjamins, R. Chatila, and F. Herrera, “Explainable artificial intelligence (xai): Concepts, taxonomies, opportunities and challenges toward responsible ai,” *Information Fusion*, vol. 58, pp. 82–115, 10 2019.

- [15] D. Gunning, M. Stefik, J. Choi, T. Miller, S. Stumpf, and G. Z. Yang, “Xai—explainable artificial intelligence,” *Science Robotics*, vol. 4, 12 2019.
- [16] D. Brigo, X. Huang, A. Pallavicini, and H. S. de Ocariz Borde, “Interpretability in deep learning for finance: a case study for the heston model,” *SSRN Electronic Journal*, 4 2021.
- [17] R. García-Céspedes, F. J. Alias-Carrascosa, and M. Moreno, “On machine learning models explainability in the banking sector: the case of shap,” *Journal of the Operational Research Society*, 2025.
- [18] Y. Sato, “Model-free reinforcement learning for financial portfolios: A brief survey,” *arXiv preprint arXiv:1904.04973*, 4 2019.
- [19] B. Bruce and J. Greene, *Chapter 4 - Portfolio Construction*, pp. 133–178. Academic Press, 2014.
- [20] X. Li, Y. Li, Y. Zhan, and X.-Y. Liu, “Optimistic bull or pessimistic bear: Adaptive deep reinforcement learning for stock portfolio allocation,” *arXiv preprint arXiv:1907.01503*, 6 2019.
- [21] H. Markowitz, “Portfolio selection,” *The Journal of Finance*, vol. 7, pp. 77–91, 3 1952.
- [22] F. M. Clinic, “Mean variance portfolio theory.”
- [23] IBM, “What is machine learning (ml)?,” 9 2021.
- [24] V. François-Lavet, P. Henderson, R. Islam, M. G. Bellemare, and J. Pineau, “An introduction to deep reinforcement learning,” *Foundations and Trends in Machine Learning*, vol. 11, pp. 219–354, 12 2018.
- [25] I. Goodfellow, Y. Bengio, and A. Courville, *Deep Learning*. MIT Press, 2016. <http://www.deeplearningbook.org>.

- [26] C. Tang, B. Abbatematteo, J. Hu, R. Chandra, R. Martín-Martín, and P. Stone, “Deep reinforcement learning for robotics: A survey of real-world successes,” *Annual Review of Control, Robotics, and Autonomous Systems*, vol. 8, pp. 153–188, 8 2024.
- [27] D. Silver, A. Huang, C. J. Maddison, A. Guez, L. Sifre, G. V. D. Driessche, J. Schrittwieser, I. Antonoglou, V. Panneershelvam, M. Lanctot, S. Dieleman, D. Grewe, J. Nham, N. Kalchbrenner, I. Sutskever, T. Lillicrap, M. Leach, K. Kavukcuoglu, T. Graepel, and D. Hassabis, “Mastering the game of go with deep neural networks and tree search,” *Nature* 2016 529:7587, vol. 529, pp. 484–489, 1 2016.
- [28] K. Shao, Z. Tang, Y. Zhu, N. Li, and D. Zhao, “A survey of deep reinforcement learning in video games,” *arXiv preprint arXiv:1912.10944*, 12 2019.
- [29] R. Bellman, “A markovian decision process,” *Journal of Mathematics and Mechanics*, vol. 6, pp. 679–684, 1957.
- [30] R. Bellman, *Dynamic Programming*. Princeton University Press, 1 1957.
- [31] S. Thrun, “Efficient exploration in reinforcement learning,” tech. rep., Carnegie Mellon University, 1 1992.
- [32] 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*, 12 2013.
- [33] M. G. Bellemare, W. Dabney, and R. Munos, “A distributional perspective on reinforcement learning,” *34th International Conference on Machine Learning, ICML 2017*, vol. 1, pp. 693–711, 7 2017.
- [34] V. Mnih, A. P. Badia, L. Mirza, A. Graves, T. Harley, T. P. Lillicrap, D. Silver, and K. Kavukcuoglu, “Asynchronous methods for deep reinforcement learning,” *33rd*

*International Conference on Machine Learning, ICML 2016*, vol. 4, pp. 2850–2869, 2 2016.

- [35] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. K. Openai, “Proximal policy optimization algorithms,” *arXiv preprint arXiv:1707.06347*, 7 2017.
- [36] T. P. Lillicrap, J. J. Hunt, A. Pritzel, N. Heess, T. Erez, Y. Tassa, D. Silver, and D. Wierstra, “Continuous control with deep reinforcement learning,” *4th International Conference on Learning Representations, ICLR 2016 - Conference Track Proceedings*, 9 2015.
- [37] S. Fujimoto, H. V. Hoof, and D. Meger, “Addressing function approximation error in actor-critic methods,” *35th International Conference on Machine Learning, ICML 2018*, vol. 4, pp. 2587–2601, 2 2018.
- [38] T. Haarnoja, A. Zhou, P. Abbeel, and S. Levine, “Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor,” *35th International Conference on Machine Learning, ICML 2018*, vol. 5, pp. 2976–2989, 1 2018.
- [39] D. H. G. B. Tokyo and J. Schmidhuber, “Recurrent world models facilitate policy evolution,” *Advances in Neural Information Processing Systems*, vol. 31, 2018.
- [40] D. Silver, T. Hubert, J. Schrittwieser, I. Antonoglou, M. Lai, A. Guez, M. Lanctot, L. Sifre, D. Kumaran, T. Graepel, T. Lillicrap, K. Simonyan, and D. Hassabis, “Mastering chess and shogi by self-play with a general reinforcement learning algorithm,” *arXiv preprint arXiv:1712.01815*, 12 2017.
- [41] Y. Wu, E. Mansimov, S. Liao, R. Grosse, and J. Ba, “Scalable trust-region method for deep reinforcement learning using kronecker-factored approximation,” *Advances in Neural Information Processing Systems*, vol. 2017-December, pp. 5280–5289, 8 2017.

- [42] J. Schulman, S. Levine, P. Moritz, M. Jordan, and P. Abbeel, “Trust region policy optimization,” *32nd International Conference on Machine Learning, ICML 2015*, vol. 3, pp. 1889–1897, 2 2015.
- [43] G. E. Uhlenbeck and L. S. Ornstein, “On the theory of the brownian motion,” *Physical Review*, vol. 36, p. 823, 9 1930.



## Appendix A

# DRL Algorithms

---

**Algorithm 1** Advantage Actor-Critic (A2C) Pseudo-code

---

**Initialise:**Global shared policy parameters  $\theta$  and value parameters  $\theta_v$ Number of parallel workers  $N$ Global step counter  $T \leftarrow 0$ Hyper-parameters: discount  $\gamma$ , max steps per update  $t_{\max}$ , max total steps  $T_{\max}$ , learning rates  $\alpha_\pi, \alpha_v$ **repeat**Reset gradients:  $d\theta \leftarrow 0, d\theta_v \leftarrow 0$ 

Initialise empty batch storage for all workers

**for** worker  $i = 1$  to  $N$  **do** $t_{\text{start}} \leftarrow t$ Get initial state  $s_t^{(i)}$  from worker  $i$ **repeat**Select action  $a_t^{(i)} \sim \pi_\theta(\cdot | s_t^{(i)})$ Execute  $a_t^{(i)}$ , observe reward  $r_t^{(i)}$  and next state  $s_{t+1}^{(i)}$ Store  $(s_t^{(i)}, a_t^{(i)}, r_t^{(i)})$  in worker  $i$ 's trajectory $t \leftarrow t + 1$ **until** terminal  $s_t^{(i)}$  or  $t - t_{\text{start}} == t_{\max}$ 

$$R^{(i)} = \begin{cases} 0 & \text{if terminal } s_t^{(i)} \\ V_{\theta_v}(s_t^{(i)}) & \text{otherwise} \end{cases}$$

**for**  $j \in \{t - 1, \dots, t_{\text{start}}\}$  **do**

$$R^{(i)} \leftarrow r_j^{(i)} + \gamma R^{(i)}$$

Accumulate gradients w.r.t.  $\theta$ :

$$d\theta \leftarrow d\theta + \nabla_\theta \log \pi_\theta(a_j^{(i)} | s_j^{(i)}) (R^{(i)} - V_{\theta_v}(s_j^{(i)}))$$

Accumulate gradients w.r.t.  $\theta_v$ :

$$d\theta_v \leftarrow d\theta_v + \nabla_{\theta_v} (R^{(i)} - V_{\theta_v}(s_j^{(i)}))^2$$

**end for****end for**

// Synchronous update: wait for all workers to complete

Average gradients:  $d\theta \leftarrow \frac{1}{N}d\theta, d\theta_v \leftarrow \frac{1}{N}d\theta_v$ Update  $\theta \leftarrow \theta + \alpha_\pi d\theta, \theta_v \leftarrow \theta_v + \alpha_v d\theta_v$  $T \leftarrow T + N \times t_{\max}$ **until**  $T > T_{\max}$ 

---

---

**Algorithm 2** Proximal Policy Optimization (PPO) Pseudo-code

---

**Initialise:**Policy parameters  $\theta_0$  and value function parameters  $\phi_0$ Global step counter  $T \leftarrow 0$ Hyper-parameters: discount  $\gamma$ , GAE parameter  $\lambda$ , clipping parameter  $\epsilon$ , learning rates  $\alpha_\pi, \alpha_v$ , epochs per update  $K_{\text{epochs}}$ , minibatch size  $N_{\text{minibatch}}$ , loss coefficients  $c_1, c_2$ **for**  $k = 0, 1, 2, \dots$  **do**Collect set of trajectories  $\mathcal{D}_k = \{\tau_i\}$  by running policy  $\pi_{\theta_k}$  in environmentStore transitions:  $\{(s_t, a_t, r_t, s_{t+1}, \text{done}_t)\}$ **for** each trajectory  $\tau$  in  $\mathcal{D}_k$  **do**Compute value estimates:  $V_t = V_{\phi_k}(s_t)$ Compute TD residuals:  $\delta_t = r_t + \gamma V_{t+1}(1 - \text{done}_t) - V_t$ Compute GAE advantages:  $\hat{A}_t = \sum_{l=0}^{\infty} (\gamma \lambda)^l \delta_{t+l}$ Compute returns:  $\hat{R}_t = \hat{A}_t + V_t$ **end for****for** epoch  $e = 1$  to  $K_{\text{epochs}}$  **do**Shuffle dataset  $\mathcal{D}_k$ **for** each minibatch  $\mathcal{B}$  in  $\mathcal{D}_k$  **do**

$$r_t(\theta) = \frac{\pi_\theta(a_t|s_t)}{\pi_{\theta_k}(a_t|s_t)}$$

$$L^{\text{CLIP}}(\theta) = \frac{1}{|\mathcal{B}|} \sum_{t \in \mathcal{B}} \min \left( r_t(\theta) \hat{A}_t, \text{clip}(r_t(\theta), 1 - \epsilon, 1 + \epsilon) \hat{A}_t \right)$$

$$L^{VF}(\phi) = \frac{1}{|\mathcal{B}|} \sum_{t \in \mathcal{B}} \left( V_\phi(s_t) - \hat{R}_t \right)^2$$

$$L(\theta, \phi) = L^{\text{CLIP}}(\theta) - c_1 L^{VF}(\phi) + c_2 S[\pi_\theta]$$

$$\theta \leftarrow \theta + \alpha_\pi \nabla_\theta L^{\text{CLIP}}(\theta)$$

$$\phi \leftarrow \phi - \alpha_v \nabla_\phi L^{VF}(\phi)$$

**end for****end for**Update policy:  $\theta_{k+1} = \theta$ Update value function:  $\phi_{k+1} = \phi$  $T \leftarrow T + |\mathcal{D}_k|$ **end for**

---

---

**Algorithm 3** Deep Deterministic Policy Gradient (DDPG) Pseudo-code

---

**Initialise:**

Critic network  $Q_{\theta^Q}(s, a)$  and actor  $\mu_{\theta^\mu}(s)$  with random weights  $\theta^Q, \theta^\mu$

Target networks:  $\theta^{Q'} \leftarrow \theta^Q, \theta^{\mu'} \leftarrow \theta^\mu$

Replay buffer  $\mathcal{B}$

Hyper-parameters: discount  $\gamma$ , soft update rate  $\tau$ , batch size  $N$ , exploration noise process  $\mathcal{N}$ , learning rates  $\alpha_Q, \alpha_\mu$ , total episodes  $M$ , steps per episode  $T$

**for** episode = 1 to  $M$  **do**

    Initialise random process  $\mathcal{N}$  for action exploration

    Receive initial state  $s_1$

**for**  $t = 1$  to  $T$  **do**

        Select action  $a_t = \mu_{\theta^\mu}(s_t) + \mathcal{N}_t$

        Execute  $a_t$ , observe reward  $r_t$  and next state  $s_{t+1}$

        Store  $(s_t, a_t, r_t, s_{t+1})$  in  $\mathcal{B}$

        Sample mini-batch  $\{(s_i, a_i, r_i, s_{i+1})\}_{i=1}^N$  from  $\mathcal{B}$

        Compute target:  $y_i = r_i + \gamma Q_{\theta^{Q'}}(s_{i+1}, \mu_{\theta^{\mu'}}(s_{i+1}))$

        Update critic by minimising:  $L(\theta^Q) = \frac{1}{N} \sum_i (y_i - Q_{\theta^Q}(s_i, a_i))^2$

        Update actor by policy gradient:

$$\nabla_{\theta^\mu} J \approx \frac{1}{N} \sum_i \nabla_a Q_{\theta^Q}(s_i, a)|_{a=\mu_{\theta^\mu}(s_i)} \nabla_{\theta^\mu} \mu_{\theta^\mu}(s_i)$$

        Update target networks:

$$\theta^{Q'} \leftarrow \tau \theta^Q + (1 - \tau) \theta^{Q'}$$

$$\theta^{\mu'} \leftarrow \tau \theta^\mu + (1 - \tau) \theta^{\mu'}$$

**end for**

**end for**

---

---

**Algorithm 4** Twin Delayed Deep Deterministic Policy Gradient (TD3) Pseudo-code
 

---

**Initialise:**

 Critic networks  $Q_{\theta_1}(s, a)$  and  $Q_{\theta_2}(s, a)$  with random weights  $\theta_1, \theta_2$ 

 Actor policy  $\mu_\phi(s)$  with random weights  $\phi$ 

 Target networks:  $\theta'_1 \leftarrow \theta_1, \theta'_2 \leftarrow \theta_2, \phi' \leftarrow \phi$ 

 Replay buffer  $\mathcal{B}$ 

 Hyper-parameters: discount  $\gamma$ , target policy noise  $\tilde{\sigma}$ , noise clip  $c$ , policy delay  $d$ , exploration noise  $\sigma$ , update rate  $\tau$ , batch size  $N$ , total steps  $T$ 
**for**  $t = 1$  **to**  $T$  **do**

 Observe state  $s_t$ 

 Sample action with exploration:  $a_t = \mu_\phi(s_t) + \epsilon, \epsilon \sim \mathcal{N}(0, \sigma)$ 

 Execute  $a_t$ , observe reward  $r_t$  and next state  $s_{t+1}$ 

 Store  $(s_t, a_t, r_t, s_{t+1})$  in  $\mathcal{B}$ 

 Sample mini-batch  $\{(s_i, a_i, r_i, s_{i+1})\}_{i=1}^N$  from  $\mathcal{B}$ 

 Sample clipped noise:  $\tilde{\epsilon} \sim \text{clip}(\mathcal{N}(0, \tilde{\sigma}), -c, c)$ 

 Compute target action:  $\tilde{a}_{i+1} = \mu_{\phi'}(s_{i+1}) + \tilde{\epsilon}$ 

 Compute target Q-value:  $y_i = r_i + \gamma \min_{j=1,2} Q_{\theta'_j}(s_{i+1}, \tilde{a}_{i+1})$ 

 Update each critic by minimising  $L(\theta_j) = \frac{1}{N} \sum_i (y_i - Q_{\theta_j}(s_i, a_i))^2$ 
**if**  $t \bmod d = 0$  **then**

 Update actor by policy gradient:  $\nabla_\phi J \approx \frac{1}{N} \sum_i \nabla_a Q_{\theta_1}(s_i, a) |_{a=\mu_\phi(s_i)} \nabla_\phi \mu_\phi(s_i)$ 

Update target networks:

 $\theta'_j \leftarrow \tau \theta_j + (1 - \tau) \theta'_j$  for  $j = 1, 2$ 
 $\phi' \leftarrow \tau \phi + (1 - \tau) \phi'$ 
**end if**
**end for**


---

---

**Algorithm 5** Soft Actor-Critic (SAC) Pseudo-code

---

**Initialise:**Actor network  $\pi_\theta(a|s)$  with parameters  $\theta$ Two critic networks  $Q_{\phi_1}(s, a)$  and  $Q_{\phi_2}(s, a)$  with parameters  $\phi_1, \phi_2$ Target critic networks:  $\phi'_1 \leftarrow \phi_1, \phi'_2 \leftarrow \phi_2$ Replay buffer  $\mathcal{D}$ Hyper-parameters: discount  $\gamma$ , temperature  $\alpha$  (fixed or learnable), target entropy  $\mathcal{H}$  (if  $\alpha$  is learnable), batch size  $N$ , learning rates  $\lambda_Q, \lambda_\pi, \lambda_\alpha$ , soft update rate  $\tau$ **for** each training step **do**Observe state  $s_t$ Sample action:  $a_t \sim \pi_\theta(\cdot|s_t)$ Execute  $a_t$ , observe reward  $r_t$  and next state  $s_{t+1}$ Store  $(s_t, a_t, r_t, s_{t+1})$  in  $\mathcal{D}$ **if** time to update **then**Sample mini-batch  $\{(s_i, a_i, r_i, s_{i+1})\}_{i=1}^N$  from  $\mathcal{D}$ **Update Critics:****for**  $j = 1, 2$  **do**Sample next actions:  $\tilde{a}_{i+1} \sim \pi_\theta(\cdot|s_{i+1})$ 

Compute target Q-values:

$$y_i = r_i + \gamma \left( \min_{k=1,2} Q_{\phi'_k}(s_{i+1}, \tilde{a}_{i+1}) - \alpha \log \pi_\theta(\tilde{a}_{i+1}|s_{i+1}) \right)$$

Update critic:  $\phi_j \leftarrow \phi_j - \lambda_Q \nabla_{\phi_j} \frac{1}{N} \sum_i (Q_{\phi_j}(s_i, a_i) - y_i)^2$ **end for****Update Actor:**Sample actions with reparametrisation:  $\tilde{a}_i = f_\theta(\epsilon_i; s_i)$  where  $\epsilon_i \sim \mathcal{N}(0, I)$ 

Compute policy loss:

$$J(\theta) = \frac{1}{N} \sum_i [\alpha \log \pi_\theta(\tilde{a}_i|s_i) - \min_{j=1,2} Q_{\phi_j}(s_i, \tilde{a}_i)]$$

Update actor:  $\theta \leftarrow \theta - \lambda_\pi \nabla_\theta J(\theta)$ **if**  $\alpha$  is learnable **then****Update Temperature:**

$$J(\alpha) = \frac{1}{N} \sum_i \alpha (\log \pi_\theta(\tilde{a}_i|s_i) + \mathcal{H})$$

Update temperature:  $\alpha \leftarrow \alpha - \lambda_\alpha \nabla_\alpha J(\alpha)$ **end if****Update Target Networks:** $\phi'_j \leftarrow \tau \phi_j + (1 - \tau) \phi'_j$  for  $j = 1, 2$ **end if****end for**

---