



Department of Informatics
King's College London
United Kingdom

7CCSMPRJ - Individual Project

Explainable Deep Reinforcement Learning for Portfolio Optimisation of Financial Assets

Name: **Ingrid Pérez Aguilera**

Student Number: K24087939

Course: Computational Finance M.Sc.

Supervisor: **Riaz Ahmad**

Word count: 11521

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

1	Introduction	1
1.1	Objectives	3
1.2	Report Structure	4
2	Background	6
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
2.3	Explainable Artificial Intelligence	23
2.3.1	Feature Importance	25
2.3.2	Local Interpretable Model-agnostic Explanations	25
2.3.3	Shapley Additive Explanations	27
2.4	State of the Art of Deep Learning and Explainability for Portfolio Optimisation	32
3	Methodology	36
3.1	Problem Definition	36
3.2	Markov Decision Process Model	37

3.3	Deep Reinforcement Learning Algorithms	40
3.3.1	Hyper-parameter tuning	41
3.4	Post-hoc Explainability	42
3.4.1	Surrogate Model Explainability	42
3.4.2	Direct Model Explainability	44
4	Results	45
4.1	Dataset	45
4.2	Experiment Design	47
4.3	Evaluation	50
4.3.1	Performance Metrics	50
4.3.2	Benchmark Strategies	52
5	Legal, Social, Ethical and Professional Issues	54
5.1	Legal Issues	54
5.2	Social Issues	55
5.3	Ethical Issues	56
5.4	Professional Issues	57
6	Conclusion	58
	References	i
A	Algorithms	xi
A.1	Deep Reinforcement Learning Algorithms	xi
A.2	Explainability Algorithms	xvii
B	State Representation	xviii
B.1	Technical Indicators	xviii
B.2	Macroeconomic Indicators	xxiii

C Hyper-parameter tuning	xxv
D Datasets	xxvii
D.1 Equities	xxvii
D.1.1 Dow Jones 30	xxvii
D.1.2 Euro Stoxx 50	xxix
D.1.3 FTSE 100	xxxiii
D.2 Commodities	xxxviii
D.3 Currencies	xxxix
E Evaluation Metrics	xl

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
2.5	Taxonomy of Explainable Artificial Intelligence Methods [3]	24

List of Tables

2.1	Comparison of Deep Reinforcement Learning Algorithms	23
4.1	Train-Validation-Test Split for each dataset	47
4.2	Summary of experiments conducted.	50
C.1	Hyper-parameter tuning configurations for different RL algorithms. . . .	xxvi
D.1	Constituents of the Dow Jones Industrial Average index as of April 2025. [4]	xxix
D.2	Constituents of the Euro Stoxx 50 index as of April 2025. [5]	xxxiii
D.3	Constituents of the FTSE100 index as of April 2025. [6]	xxxviii
D.4	6 Commodities Futures Contracts	xxxviii
D.5	10 Currency Pairs	xxxix

List of Algorithms

1	Advantage Actor-Critic (A2C) Pseudo-code	xii
2	Proximal Policy Optimisation (PPO) Pseudo-code	xiii
3	Deep Deterministic Policy Gradient (DDPG) Pseudo-code	xiv
4	Twin Delayed Deep Deterministic Policy Gradient (TD3) Pseudo-code . .	xv
5	Soft Actor-Critic (SAC) Pseudo-code	xvi
6	Shapley Value Approximation	xvii

Nomenclature

Numbers and Arrays

a	Scalar
\mathbf{a}	Vector
\mathbf{A}	Matrix
$\mathbf{1}$	Matrix of 1's

Sets

\mathbb{A}	Set
\mathbb{R}	Set of real numbers
$[a, b]$	Real interval including a and b
(a, b)	Real interval excluding a and b
$[a, b)$	Real interval including a but excluding b
$(a, b]$	Real interval excluding a but including b
$\{0, 1\}$	Set containing 0 and 1
$\{0, 1, \dots, n\}$	Set of all integers between 0 and n inclusive

\emptyset	Empty set
$\mathbb{A} \subset \mathbb{B}$	Set \mathbb{A} is a subset of set \mathbb{B}
$\mathbb{A} \subseteq \mathbb{B}$	Set \mathbb{A} is a subset or equal to set \mathbb{B}
$\mathbb{A} \cap \mathbb{B}$	Intersection of sets \mathbb{A} and \mathbb{B}
$\mathbb{A} \cup \mathbb{B}$	Union of sets \mathbb{A} and \mathbb{B}
$\mathbb{A} \setminus \mathbb{B}$	Set difference of sets \mathbb{A} and \mathbb{B}

Indexing

a_i	Element i of vector \mathbf{a} , with indexing start at 1
$A_{i,j}$	Element i, j of matrix \mathbf{A} , with indexing start at 1

Algebra Operations

\mathbf{A}^T	Transpose of matrix \mathbf{A}
\mathbf{A}^{-1}	Inverse of matrix \mathbf{A}

Probability

$P(x)$	Probability distribution over a random variable x
$E[x]$	Expected value of random variable x
$Var[x]$	Variance of random variable x
$Cov[x, y]$	Covariance of random variables x and y
μ	Mean
σ	Standard deviation
σ^2	Variance

Σ	Covariance matrix
$\mathcal{N}(\mu, \sigma^2)$	Normal distribution with mean μ and variance σ^2
$\mathcal{N}(0, 1)$	Standard normal distribution

Functions

$f : \mathbb{A} \rightarrow \mathbb{B}$	Function f with domain \mathbb{A} and Range \mathbb{B}
$f(x)$	Function f evaluated at x
$f(x; \theta)$	Function of x with parameters θ
$f'(x)$	Derivative of function f at x
$\int f(x) dx$	Integral of function f
$\min f(x)$	Minimum of function f
$\min_x f(x)$	Minimum of function f with respect to x
$\max f(x)$	Maximum of function f
$\max_x f(x)$	Maximum of function f with respect to x
$\log x$	Natural logarithm of x
$x!$	Factorial of x

Markov Decision Process

t	Time step
\mathcal{S}	State space
s_t	Observation of the state at time step t
\mathcal{A}	Action space

a_t	Action taken at time step t
$R(s, a, s')$	Reward function
r_t	Reward received at time step t
$T(s, a, s')$	Transition function
γ	Discount factor
π	Policy
G	Long-term reward
V^π	State Value function for policy π
Q^π	State-Action Value function for policy π
A^π	Advantage function for policy π

Other Symbols

\forall	For all
\in	Element of
∇	Gradient
\sum	Summation symbol

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.

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

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

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.

Black Box Term used to describe a system or model whose internal workings are not transparent or easily understood.

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.

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.

Deep Neural Network Type of artificial neural network with a succession of layers of non-linear transformations capable of learning a representation of the data with various levels of abstraction.

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.

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.

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

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.

Exploitation Process of leveraging known information to make decisions or take actions that are expected to yield the highest reward.

Exploration Process of investigating and experimenting with different actions or strategies to discover their effects and improve decision-making.

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.

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.

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.

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.

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

Markov Decision Process Mathematical model for sequential-decision making where the outcomes are uncertain, characterised by states, actions, rewards and a transition function.

Mean-Variance Optimisation Investment strategy 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.

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.

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.

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.

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.

Python Interpreted, object-oriented, high-level programming language.

Reinforcement Learning Subset of machine learning where an agent learns to make decisions by taking actions in an environment to maximise cumulative reward.

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.

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.

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.

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

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.

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

Acronyms

A2C Advantage Actor-Critic.

A3C Asynchronous Advantage Actor-Critic.

ADX Average Directional Index.

AI Artificial Intelligence.

AI Act Artificial Intelligence Act.

ATR Average True Range.

AUD Australian Dollar.

BCS British Computer Society.

BXY British Pound Currency Index.

CAD Canadian Dollar.

CCI Commodity Channel Index.

CHF Swiss Franc.

CNY Chinese Yuan.

DDPG Deep Deterministic Policy Gradient.

DJIA Dow Jones Industrial Average.

DL Deep Learning.

DNN Deep Neural Network.

DPG Deterministic Policy Gradient.

DQN Deep Q-Network.

DRL Deep Reinforcement Learning.

DW30 Dow Jones 30 Index.

DX Directional Movement Index.

DXY U.S. Dollar Index.

EMA Exponential Moving Average.

EU European Union.

EUR Euro.

Euro Stoxx 50 Euro Stoxx 50 Index.

EXY Euro Index.

FCA Financial Conduct Authority.

FTSE 100 Financial Times Stock Exchange 100 Index.

FVX 5-Year Treasury Yield.

GBP British Pound.

HKD Hong Kong Dollar.

IET The Institution of Engineering and Technology.

INR Indian Rupee.

IRX 3-Month Treasury Yield.

JPY Japanese Yen.

KRW South Korean Won.

LIME Local Interpretable Model-agnostic Explanations.

LLM Large Language Model.

LSTM Long Short-Term Memory.

MACD Moving Average Convergence Divergence.

MAPE Mean Absolute Percentage Error.

MDI Negative Directional Index.

MDP Markov Decision Process.

MiFID II Markets in Financial Instruments Directive II.

ML Machine Learning.

MPT Modern Portfolio Theory.

MSE Mean Squared Error.

MVO Mean-Variance Optimisation.

OHLCV Open, High, Low, Close, Volume.

PCA Principal Component Analysis.

PDI Positive Directional Index.

POMDP Partially Observable Markov Decision Process.

PPO Proximal Policy Optimisation.

RL Reinforcement Learning.

ROC Rate of Change.

RSI Relative Strength Index.

S&P 500 Standard and Poor's 500 Index.

SAAC Synchronous Advantage Actor Critic.

SAC Soft Actor-Critic.

SHAP SHapley Additive exPlanations.

SMA Simple Moving Average.

TD3 Twin Delayed Deep Deterministic Policy Gradient.

TNX 10-Year Treasury Yield.

TRPO Trust Region Policy Optimisation.

U.S. United States.

UK United Kingdom.

USD United States Dollar.

VaR Value at Risk.

VIX Volatility Index.

VXD Volatility DW30.

XAI Explainable Artificial Intelligence.

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. [7]

Algorithmic trading focuses on the application of analytical methods to automatically execute trading actions based on an algorithm without human intervention. Initially, the field mainly studied the usage of a computer program to follow a predefined strategy [8]. More recently, algorithmic trading involves environment perception by learning feature representation from highly non-stationary and noisy financial time series data, and decision-making by exploring the environment and simultaneously taking actions without supervision [9].

Machine Learning (ML) has an advantage for the task due to its ability to learn from historical data and make predictions about the future state of an environment. Research

has explored the application of Deep Learning (DL) in future price prediction of financial assets [10, 11, 7, 12]. 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) allows the algorithm to learn a trading strategy directly from the environment, without the need for a separate step [13, 14]. 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 [15], or second, the algorithm can learn the optimal portfolio allocation and automatically rebalance the portfolio weights at each time step [16].

Despite the potential of RL in portfolio optimisation, its widespread adoption in the financial industry remains limited. This is primarily due to following challenges [17]:

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 [18]. The term was first coined in 2016 to describe the need for users to effectively understand, trust and manage artificial intelligence applications [19]. The need for explainability becomes particularly relevant within the financial domain, where the regulatory framework requires transparency and accountability in automated decision-making. Various relevant applications, including volatility models [20], credit risk assessment [21] and portfolio

construction [17], have explored the concept of explainability in financial applications.

Consequently, this thesis will focus on addressing the aforementioned challenges by exploring the application of Deep Reinforcement Learning (DRL) to portfolio optimisation and implementing 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 and PPO is better at balancing exploration versus exploitation.

Second, post-hoc explainability techniques: Feature Importance, Local Interpretable

Model-agnostic Explanations (LIME) and SHapley Additive exPlanations (SHAP) 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 portfolio allocation and 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). In addition, it gives an overview of the post-hoc explainability techniques: Feature Importance, Local Interpretable Model-agnostic Explanations (LIME) analysis and SHapley Additive exPlanations (SHAP) values, 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 tackle 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 and reward functions. Furthermore, given the trained algorithms, it describes the implementation of the post-hoc explainability techniques.

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, feature importance, LIME and SHAP.

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 relevant post-hoc explainability techniques to achieve interpretable and transparent models. 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 [22]. 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 [23]. 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 [24].

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 [25] 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 [26]. 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 σ_p^2 subject to a target expected return R_p , or
- maximising expected return R_p subject to a risk constraint σ_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 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 sub-optimal, while those on the frontier represent the best achievable combinations of risk and return.

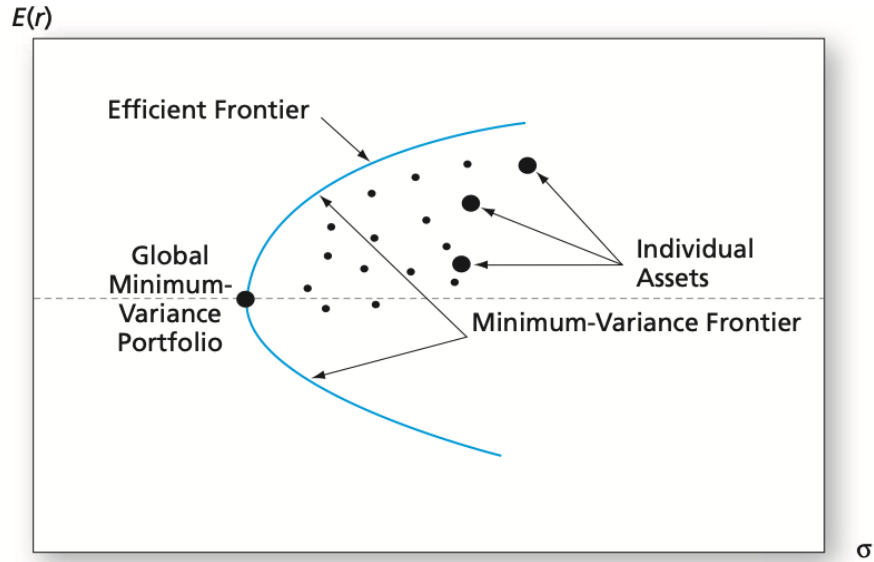


Figure 2.1: Efficient Frontier in Risk-Return Space. [1]

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

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 by gradually improving their accuracy over time [27]. There are three main tasks in ML [28]:

- **Supervised Learning:** Task of training a classification or regression model from

labelled data, where the model learns to map inputs to outputs based on examples.

- **Unsupervised Learning:** Task of drawing inferences from datasets consisting of unlabelled 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 (DNNs) [29], which are characterised by a succession of layers of non-linear transformations 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 [30] and game playing [31, 32].

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 [28]. The agent's interaction with the environment is visually represented in Figure 2.2.

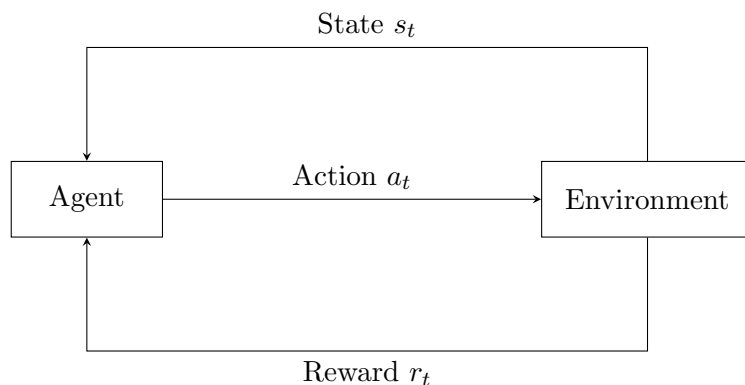


Figure 2.2: Agent interaction with environment

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 only depends on the current state and

action and not on the history of previous states and actions. Consequently, a Markov Decision Process [33] 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.

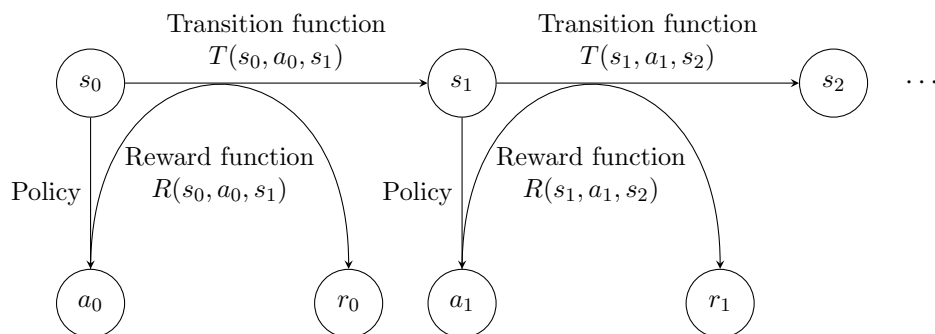


Figure 2.3: Markov Decision Process with policy, transition, and reward functions

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, a 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 γ 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 π :

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

where \mathbb{E}_π denotes the expectation over the policy π .

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 π :

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

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 π 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 π^* 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 [34] 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:

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

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

and the Bellman optimality equations are defined as:

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

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

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 [35]. 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.

2.2.2 Deep Reinforcement Learning Algorithms

A Reinforcement Learning (RL) agent includes one or more of the following components [28]:

- 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, and
- a model of the environment, consisting of estimates of transition and reward functions.

Depending on the components, the main RL paradigms are:

- **Model-free** algorithms do not learn a representation of the environment, but focus on the value function, the policy or both. These algorithms can be further divided into:
 - **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) [36] and C51 [37].
 - **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) [38] and Proximal Policy Optimisation (PPO) [39] 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) [40], Twin Delayed Deep Deterministic Policy Gradient (TD3) [41], and Soft Actor-Critic (SAC) [42].

- **Model-based** algorithms include a model of the environment, which can be used to simulate future states and rewards. For example, World Models [43] learns a model of the environment and AlphaZero [44] has a representation of the model.

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

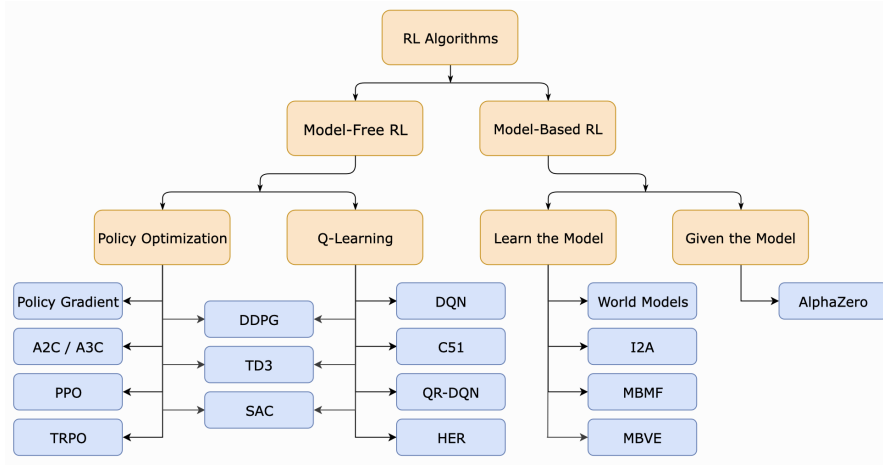


Figure 2.4: Taxonomy of Reinforcement Learning Algorithms [2]

The full potential of RL 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 Deep Neural Networks to approximate the value functions, or use gradient ascent to find the optimal policy parameters. This thesis will focus 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* [38]. 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 [45].

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 γ 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.1.

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* [39] with the objective of constraining policy updates. The algorithm balances sufficiently large 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) [46] by using a clipped surrogate objective function that ensures that the probability ratio r_t is bounded within a range of $[1 - \epsilon, 1 + \epsilon]$, where ϵ is a hyperparameter that controls the clipping range. The objective function is defined as:

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

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 π_θ 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.1.

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* [40]. The algorithm arises to solve the challenges of applying Deep Q-Network [36] 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:

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

where \mathcal{D} is the replay buffer that stores the agent’s experiences, θ_μ are the parameters of the actor network, and θ_Q are the parameters of the critic network, and the critic network is updated using the Bellman equations:

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

where θ'_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 [47], 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.1.

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* [41]. The proposal addresses the hyper-parameter sensitivity and overestimation bias present in the critic network of DDPG. The main improvements are:

- **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.1.

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* [42] 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 α 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 θ_{Q_i} are the parameters of the two critic networks. The temperature parameter α 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.1.

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 Function Estimation
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

2.3 Explainable Artificial Intelligence

Explainable Artificial Intelligence (XAI) refers to a set of processes, methods and techniques that enable human users to comprehend and trust the outcomes and decisions made by Artificial Intelligence (AI) systems. The goal is to make AI algorithms more transparent, interpretable, and accountable, allowing end-users to understand the outputs of the models. This is particularly important in Deep Learning (DL) where the models are often considered black boxes due to their complexity, non-linearity and high-dimensionality, making it difficult to understand how they arrive at their decisions. With explainable systems, the benefits are numerous ranging from informed decision making to increased user adoption and better governance [48].

Although there are many possible classifications of XAI methods, they can be broadly categorised into two main categories. First, transparent algorithms are inherently understandable and interpretable, such as linear regression, decision trees and rule-based systems. Second, post-hoc explanations are methods that require the usage of an additional algorithm to clarify the decisions made by a model. Examples include saliency

maps [49] and interaction data [50].

Another classification relates to whether the explanation method depends on the type of model it is being applied to. On the one hand, model-agnostic methods can be applied to any model regardless of its internal architecture, effectively treating all types of models as black boxes. The analysis is conducted by understanding the input-output behaviour and how small perturbations to the input impact the output. Examples include Local Interpretable Model-agnostic Explanations (LIME) [51] and SHapley Additive exPlanations (SHAP) [52]. On the other hand, model-specific methods are tailored to a particular architecture and leverage its components as part of the explanation process. In the case of neural networks, this can be achieved by analysing the learned features or by incorporating attention mechanisms [53]. The model-agnostic methods can be further sub-divided between global and local explanations. The goal of global explainability is to provide an understanding of the model’s overall behaviour across an entire dataset, while local explainability focuses on explaining the individual predictions.

The taxonomy of the XAI methods is shown in Figure 2.5.

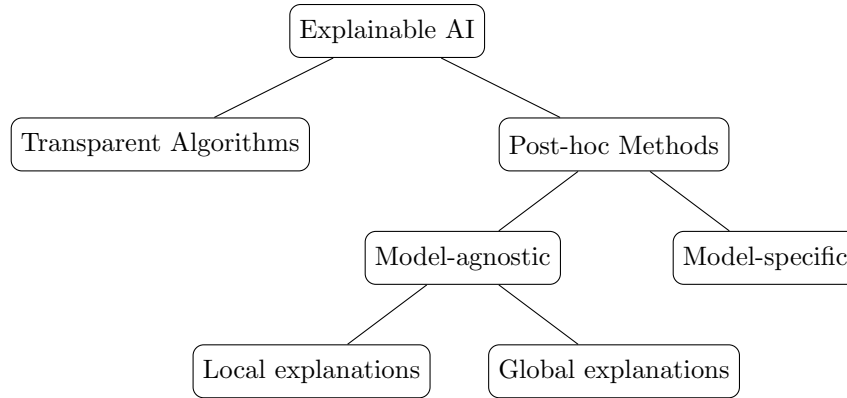


Figure 2.5: Taxonomy of Explainable Artificial Intelligence Methods [3]

2.3.1 Feature Importance

Feature Importance is a global model-specific method that quantifies the contribution of each feature to the model’s predictions. In particular, permutation feature importance [54] measures the contribution of each feature by evaluating the model’s performance on the original dataset and comparing it to the performance on a permuted version of the dataset, where a specific feature is randomly shuffled. The process allows to understand how much the model’s relies on a particular feature for its prediction, by measuring the decrease in predictive power if the input feature’s values vary.

An alternative method, well-suited for tree-based models, is impurity-based feature importance. Random forests split their features with the goal of reducing an impurity measure at each node, normally Gini impurity for classification tasks or Mean Squared Error (MSE) for regression. As such, a variable responsible for a split with a large decrease in impurity is considered important [55]. As a result, the impurity importance of a feature is the sum of all impurity reductions across all nodes and trees in the forest where the particular feature was used to split the data. While powerful, this method suffers from bias towards features with high cardinality and they do not necessarily reflect the ability of a feature to make useful predictions on the test set, given that importance is computed on the training set.

2.3.2 Local Interpretable Model-agnostic Explanations

Local Interpretable Model-agnostic Explanations (LIME) is a model-agnostic explanation method that interprets individual predictions of a Machine Learning model by analysing the model locally around the prediction of interest. The method, introduced by Ribeiro et al. (2016) in their paper on *Why should I trust you? Explaining the predictions of any classifier* [51], utilises the black box model to understand what happens

to the outputs when the input data is slightly modified, then fits a simpler, interpretable model to the perturbed data to approximate the black box model’s behaviour in that local region.

For a more rigorous definition, let $\mathcal{L}(f, g, \pi_x)$ be a measure of how unfaithful an explanation model g is in approximating the model $f : \mathbb{R}^d \rightarrow \mathbb{R}$ in the neighbourhood of the instance x defined as π_x . The goal is to find an interpretable model $g \in G$ that minimises the following objective function:

$$\xi(x) = \mathcal{L}(f, g, \pi_x) + \Omega(g) \quad (2.29)$$

where $\Omega(g)$ is a regularisation term that penalises the complexity of the explanation model g . This encourages interpretability, meaning a qualitative understanding of the relationship between inputs and outputs, and promotes local fidelity, ensuring the explanation accurately approximates the model’s behaviour in that region.

The method works as follows:

- **Instance Selection:** Start with a specific prediction instance that requires explanation.
- **Perturbation:** Generate synthetic data points by perturbing the original instance.
- **Model Querying:** Obtain predictions from the black box model for these perturbed instances.
- **Weighting:** Assign weights to the synthetic data points based on their proximity to the original instance.
- **Surrogate Training:** Train a simple, interpretable model (typically linear regression) on the weighted synthetic dataset.

- **Explanation:** Use the surrogate model to explain the original prediction by interpreting its coefficients or structure.

Its main advantages are its ability to work across different types of data and models, due to its model-agnostic nature, and the fact that its explanations are human-friendly and include a fidelity measure that indicates how well the explanation approximates the black box model’s local behaviour. Nonetheless, the method offers only local explanations, which may not generalise well to the entire dataset; it relies on the choice of neighbourhood; and the complexity of the surrogate model needs to be defined in advance and can compromise the interpretability of the explanation [3].

2.3.3 Shapley Additive Explanations

Another model-agnostic technique is SHapley Additive exPlanations (SHAP), which was introduced by Lundberg and Lee (2017) in their paper on *A unified approach to interpreting model predictions* [52]. The method is based on cooperative game theory and the concept of Shapley values [56].

The Shapley value arises in cooperative game theory to answer the question of how to fairly distribute the contribution of each player in a coalition. In this framework, a coalitional game is represented as a tuple (N, v) , where $N = \{1, 2, \dots, n\}$ is the set of players and $v : 2^N \rightarrow \mathbb{R}$ is the characteristic function that assigns a value $v(S)$ to each possible coalition $S \subseteq N$, with the convention that $v(\emptyset) = 0$. The characteristic function $v(S)$ represents the total value that coalition S can achieve through cooperation.

Mathematically, the Shapley value for a feature i is defined as:

$$\phi_i = \sum_{S \subseteq N \setminus \{i\}} \frac{|S|!(|N| - |S| - 1)!}{|N|!} [v_{S \cup \{i\}}(x_{S \cup \{i\}}) - v_S(x_S)] \quad (2.30)$$

where $S \subseteq N$ is a subset of all features N , $v_{S \cup \{i\}}$ is marginal contribution of player i

to coalition S , and v_S is the value of coalition S without player i . This formula can be understood as computing the weighted average of marginal contributions of feature i across all possible coalitions.

The advantage of Shapley values is that it is the only attribution method that results in a fair payout. In particular, it satisfies the four fundamental properties that define a fair allocation:

- **Efficiency:** The sum of all Shapley values equals the value of the grand coalition:

$$\sum_{i \in N} \phi_i = v(N)$$

- **Symmetry:** The contributions of two players i and j should be equal if they make identical marginal contributions to all possible coalitions:

$$v(S \cup \{i\}) = v(S \cup \{j\}), \forall S \subseteq N \setminus \{i, j\} \implies \phi_i = \phi_j$$

- **Dummy:** If a player contributes nothing to any coalition they join, their Shapley value should be zero:

$$v(S \cup \{i\}) = v(S), \forall S \subseteq N \setminus \{i\} \implies \phi_i = 0$$

- **Additivity:** For two games (N, v_1) and (N, v_2) , the Shapley value of the combined game $(N, v_1 + v_2)$ equals the sum of individual Shapley values:

$$\phi_i(v_1 + v_2) = \phi_i(v_1) + \phi_i(v_2)$$

The algorithm for approximating Shapley values is outlined in Appendix 6. However, the method is computationally expensive, as it requires evaluating the model for all possible

coalitions of features.

The use of Shapley values in ML was not introduced until 2011, when Štrumbelj and Kononenko [57] proposed it as a method for explaining black-box regression models. However, it was popularised in 2017 by Lundberg and Lee [52], who introduced the SHAP framework. Although their proposal relies on the principles of Shapley values, their main contribution is to represent the Shapley value as an additive feature attribution method. The explanation is given as:

$$g(\mathbf{z}') = \phi_0 + \sum_{j=1}^M \phi_j z'_j \quad (2.31)$$

where g is the explanation model, $\mathbf{z}' \in \{0, 1\}^M$ is the coalition vector, M is the maximum coalition size, ϕ_0 is the bias term, and ϕ_j is the feature attribution for feature j .

2.3.3.1 KernelSHAP

The KernelSHAP algorithm builds on the SHAP framework by using a kernel-based approach to estimate Shapley values. It approximates the Shapley value by sampling different coalitions and using a weighted linear regression model to fit the contributions of each feature. The method relies on not all coalitions contributing equally to the final Shapley value and as such uses kernel weights to identify the most important coalitions. The main steps of the KernelSHAP algorithm are:

- Sample a set of coalitions from the feature space: $\mathbf{z}'_k \in \{0, 1\}^M, k \in \{1, \dots, K\}$, where K is the total number of samples.
- For each coalition \mathbf{z}'_k , compute the model's prediction by first converting \mathbf{z}'_k to the original feature space and then applying the model $\hat{f} : \hat{f}(h_{\mathbf{x}}(\mathbf{z}'_k))$.
- Compute the weight for each coalition with the SHAP kernel: $\pi_{\mathbf{x}}(\mathbf{z}')$
- Fit a weighted linear regression model to the sampled coalitions and their corre-

sponding predictions.

- Return Shapley values ϕ_k as the coefficients of the fitted model.

The sample coalitions are generated by randomly selecting subsets of the features, or equivalently, a vector of 0s and 1s indicating whether a feature is included in the coalition. The sampled coalitions represent the dataset for the regression model whose target is the prediction for a coalition. However, the sampled coalitions are not on the target feature space and, as such, it is necessary to implement another function $h_{\mathbf{x}}(\mathbf{z}') = \mathbf{z}$ that maps the sampled coalitions to the original feature space. In the case of tabular data, this is done by replacing the features that are not included in the coalition with their mean value or a random sample from its possible values. As a result, sampling from the marginal distribution ignores the dependence structure between features.

Although there are similarities between SHAP and LIME, their main difference lies in feature weighting in the regression model. The SHAP weighting assumes that small and large coalitions provide the most information about its isolated effects or its total effects, respectively. Whereas coalitions with half the features add little information about a specific feature's contributions. The proposed weighting function for KernelSHAP is:

$$\pi_{\mathbf{x}}(\mathbf{z}') = \frac{(M-1)}{\binom{M}{|\mathbf{z}'|} |\mathbf{z}'| (M-|\mathbf{z}'|)}, \quad (2.32)$$

where M is the maximum coalition size and $|\mathbf{z}'|$ is the number of features in the coalition \mathbf{z}' .

In addition, since the coalitions with the smallest and the largest number of features are the most informative, the sampling process is biased towards coalitions of size 1 and $M-1$, resulting in $2M$ possible coalitions. Then, the remaining budget $K-2M$ is used to sample coalitions of size 2 to $M-2$. This process continues until the sampled coalitions reach the desired number of samples.

Consequently, given the samples, the KernelSHAP algorithm fits a weighted linear regression model to estimate the Shapley values. The model is defined as:

$$g(\mathbf{z}') = \phi_0 + \sum_{j=1}^M \phi_j z'_j. \quad (2.33)$$

The goal is to minimise the following loss function:

$$\mathcal{L}(\hat{f}, g, \pi_{\mathbf{x}}) = \sum_{\mathbf{z}' \in \mathcal{Z}} \left[\hat{f}(h_{\mathbf{x}}(\mathbf{z}')) - g(\mathbf{z}') \right]^2 \pi_{\mathbf{x}}(\mathbf{z}'), \quad (2.34)$$

where \mathbf{Z} is the dataset of sampled coalitions.

2.3.3.2 TreeSHAP

A variant of their original proposal is TreeSHAP [58], designed specifically for tree-based models, such as decision trees, random forests or gradient-boosted trees. In this case, the method is model-specific and by leveraging the structure of the tree, it improves the computational efficiency by reducing computation time from exponential for KernelSHAP to polynomial.

The method works by exploiting the tree structure and the main steps of the algorithm are as follows:

- For each feature, traverse the tree and compute the contribution of the feature to the prediction at each node.
- For each node, compute the contribution of the feature to the prediction by considering the difference between the prediction at the node and the prediction at the parent node.
- For each feature, compute the Shapley value by averaging the contributions across

all nodes in the tree.

2.4 State of the Art of Deep Learning and Explainability for Portfolio Optimisation

The emergence of Machine Learning (ML) and its rise in popularity have led to a new research direction in various financial applications, due to its ability to learn complex patterns from high-dimensional data and adapt to changing market conditions. Particularly, the main approaches in the financial field include the use of Deep Learning and Reinforcement Learning (RL) for price trend prediction [59] and portfolio optimisation [60].

The topic of portfolio optimisation with multiple risky financial assets has been extensively studied and it still poses a challenging task due to the complex and stochastic nature of financial markets. Traditional approaches to portfolio optimisation, such as the Mean-Variance Optimisation (MVO) [25], have been widely used, but they often rely on assumptions that do not hold in practice, such as the normality of returns and the stability of the covariance matrix. As a result, these methods can lead to suboptimal portfolios and poor performance in real-world scenarios.

The use of Deep Learning architectures is particularly prominent in the context of price prediction, where deep neural networks are employed to learn complex patterns in historical price data. For instance, Long Short-Term Memory (LSTM) architectures are particularly well-suited for financial time series forecasting due to their ability to capture long-term dependencies and temporal patterns in the data. The performance of these architectures can lead to a Mean Absolute Percentage Error (MAPE) as low as 2.72%, as demonstrated by Chaudhary (2025) [61]. Shen et al. (2020) [7] also propose a LSTM architecture for short-term trend prediction. Despite its simple architecture, its main

contributions rely on the feature engineering process. The proposal incorporates feature extension, followed by recursive feature elimination and finally, performs Principal Component Analysis (PCA) to reduce the dimensionality of the input data for the LSTM model. The use of PCA improves the training efficiency of the architecture by 36.8% [7].

The use of ensemble methods is also prominent for price prediction, where techniques such as stacking, blending, boosting and bagging are employed to combine the predictions of multiple models. For example, Nti et al. (2020) [11] explore twenty-five different ensemble classifiers and regressors based on Decision Trees, Support Vector Machines and Neural Networks. Their research shows that although stacking and boosting ensemble algorithms provide better results in terms of accuracy, they are the most computationally expensive.

When it comes to portfolio optimisation, DRL algorithms are particularly better suited as they address the sequential decision making nature of portfolio management. Unlike traditional Supervised Learning approaches that require labelled data, DRL enables agents to learn optimal investment strategies through direct interaction with market environments. This paradigm allows for the dynamic adjustment of portfolio weights without the need for predefined rules, making it suitable to capture market non-linearities and adapt to changing conditions. There are numerous DRL algorithms, as outlined in Section 2.2.2, with each being particularly befitting to different aspects of the portfolio management process, namely, DDPG encourages maximum returns, while A2C reduces the variance.

Liu et al. (2018) [15] propose the use of DDPG for profitable stock trading, where the agent learns to buy, hold and sell stocks based on the historical price data with the goal of maximising the investment return. Their proposal outperforms the Dow Jones Industrial Average (DJIA) and the min-variance portfolio allocation strategies, achieving an annualised return of 25.86% and a Sharpe ratio of 1.79. Their work is further extended

by Liu (2020) [62] by proposing a Python library, FinRL, which provides a comprehensive framework for developing and evaluating DRL algorithms in the context of financial trading. Although the library provides implementation for numerous DRL algorithms including DQN, A2C and SAC, in this paper, they only evaluate the performance of TD3 and DDPG on the Dow Jones Industrial Average (DJIA) constituents for the tasks of multiple stock trading and portfolio allocation. Notwithstanding, an interesting aspect of their work is the incorporation of the turbulence index, that measures extreme asset price fluctuations, to control portfolio risk. Beyond traditional assets, optimal portfolio allocation has also been explored in the context of cryptocurrencies [63] and future contracts [64].

Despite the significant applications of DRL in portfolio optimisation, the field still faces significant challenges that prevent its widespread adoption. When it comes to market conditions, identifying the correct objective function that guarantees both efficiency and accuracy remains an open research question. Finally, one of the most significant concerns is the lack of transparency of ML models. If there is to be widespread use of these models in the financial field, great strides must be made in the area of explainability.

Barriedo et al. (2020) [18] provide a comprehensive overview of Explainable Artificial Intelligence (XAI) methods and their need in different fields. In 2024, the European Union has passed the Regulation 2024/1689, commonly known as the AI Act [65]. The goal is to ensure that the AI systems used in the European Union are safe, transparent and traceable. Given the regulatory requirements, it is crucial for financial institutions to adopt techniques that enhance the interpretability and traceability of ML models.

The application of XAI methods has gained traction in the financial domain, with notable use cases including stock market trend prediction [66] and auditing [67]. However, despite examples of XAI methods in finance, there is still a lack of research in DRL. An interesting study was conducted by González Cortés et al. (2024) [17], whose au-

thors address the black box behaviour of DRL models by proposing an intrinsically transparent algorithm. They use Advantage Actor-Critic (A2C) ¹ enhanced with an attention-layered LSTM to determine the importance of the input features. Although their results show an improvement in performance compared against two variants of the Markowitz model, their implementation does not scale well with the number of assets in the portfolio, as the attention mechanism is independent for each stock.

Another recent study by de-la-Rica-Escudero et al. (2025) [68] proposes post-hoc explainability methods to interpret the decisions of PPO agent optimised for portfolio management. They implement three model-agnostic explainability techniques: SHAP, LIME and feature importance analysis. Despite their proposal’s ability to provide insights into the model’s decision-making process, they do not use the model prediction function directly but instead implement a surrogate model to approximate the decision boundary of the PPO agent, effectively adding an additional layer to the model.

Given the lack of research in the area of explainability for DRL models that perform automated portfolio allocation and the current limitations of existing methods, this thesis proposes to provide a comprehensive study of DRL techniques tailored for portfolio optimisation in financial markets, by exploring the potential of five prominent algorithms (A2C, PPO, DDPG, TD3 and SAC) across various market conditions and asset classes. Moreover, due to the complex and hidden nature of those algorithms, their decision-making process will be explained using XAI methods. The objective is to bridge the gap between the performance of DRL and algorithmic transparency.

¹In the paper [17], they refer to A2C as Synchronous Advantage Actor Critic (SAAC)

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 of financial assets. 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 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.

3.2 Markov Decision Process Model

Due to the dynamic, stochastic and interactive nature of financial markets, a Markov Decision Process is a suitable framework to model the problem. The main elements of the MDP model are defined as follows:

- State space \mathcal{S}
- Action space \mathcal{A}
- Reward function $R : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow \mathbb{R}$
- Transition function $T : \mathcal{S} \times \mathcal{A} \rightarrow \mathcal{S}$
- Discount factor $\gamma \in [0, 1]$

The state space \mathcal{S} is a vector representation of the financial environment. For a portfolio of D assets, the features that describe the state include asset prices, technical indicators and macroeconomic indicators:

- Close price $\mathbf{p}_t \in \mathbb{R}^D$: Adjusted close prices of the assets at time t .
- Open price $\mathbf{o}_t \in \mathbb{R}^D$: Opening prices of the assets at time t .
- High price $\mathbf{h}_t \in \mathbb{R}^D$: Highest prices of the assets at time t .
- Low price $\mathbf{l}_t \in \mathbb{R}^D$: Lowest prices of the assets at time t .
- Volume $\mathbf{v}_t \in \mathbb{R}^D$: Trading volume of the assets at time t .

- Technical indicators $\mathbf{I}_t \in \mathbb{R}^{D \times I}$: For each of the D assets, a vector \mathbf{i}_t of I technical indicators, such as moving averages, relative strength index (RSI), and Bollinger Bands, calculated from the asset prices.
- Macroeconomic indicators $\mathbf{M}_t \in \mathbb{R}^{D \times M}$: For each of the D assets, a vector \mathbf{m}_t of M macroeconomic indicators, such as volatility index and interest rates, which provide additional context about the financial environment.
- Covariance matrix $\mathbf{C}_t \in \mathbb{R}^{D \times D}$: For each of the D assets, a vector of D values representing the covariance between the assets in the portfolio at time t .

The description of technical and macroeconomic indicators is provided in more detail in Appendix B.

The action space \mathcal{A} is the set of possible actions that the agent can take at each time step. For the portfolio optimisation problem, the actions correspond to portfolio weights and are defined as follows:

$$\mathbf{a}_t = \mathbf{w}_t : w_{t,d} \in [0, 1] \quad \forall d \in \{1, \dots, D\}, \quad (3.1)$$

where \mathbf{w}_t is a vector of portfolio weights at time t , representing the allocation of the portfolio to each asset. The weights are constrained to be non-negative and sum to one:

$$\sum_{d=1}^D w_{t,d} = 1, \quad w_{t,d} \geq 0 \quad \forall d \in \{1, \dots, D\}. \quad (3.2)$$

Moreover, they are initialised to be equal for all assets, meaning that the agent starts with an equal allocation to each asset in the portfolio. The reason behind this is to avoid an initial bias and allow the agent to learn an allocation from the environment rather than favour any particular asset.

The transition function T describes how the state of the environment changes in response

to the action taken. It is defined as:

$$s_{t+1} = T(s_t, a_t), \quad (3.3)$$

where s_{t+1} is the new state of the environment after taking action a_t in state s_t . The transition function is determined by the dynamics of the financial market, which are influenced by the asset prices, trading volume and other factors.

The reward function R models the direct reward of taking an action a_t in state s_t and transitioning to a new state s_{t+1} . It is defined as the change in the portfolio value from time t to time $t + 1$:

$$R_{t+1} = R(s_t, a_t, s_{t+1}) = V_{t+1} - V_t, \quad (3.4)$$

where the value of the portfolio at time t is given by the dot product of the portfolio weights and the asset close prices:

$$V_t = \mathbf{w}_t \cdot \mathbf{p}_t. \quad (3.5)$$

An alternative formulation of the reward function is to use the Sharpe ratio [69], which is defined as the ratio of the expected return to the standard deviation of the returns:

$$R(s_t, a_t, s_{t+1}) = \frac{E[r_{t+1}]}{\sigma[r_{t+1}]}, \quad (3.6)$$

where $E[r_{t+1}]$ is the expected return of the portfolio at time $t + 1$ and $\sigma[r_{t+1}]$ is the standard deviation of the returns at time $t + 1$. This formulation encourages the agent to maximise the expected return while minimising the risk of the portfolio.

Regardless of the choice of reward function, the goal of the agent is to learn a policy $\pi : \mathcal{S} \rightarrow \mathcal{A}$ that maximises the expected cumulative reward over time, which can be

expressed as:

$$J(\pi) = \mathbb{E} \left[\sum_{t=0}^T \gamma^t R(s_t, a_t, s_{t+1}) \right], \quad (3.7)$$

where T is the time horizon and γ is the discount factor that determines the importance of future rewards.

The environment is implemented in Python¹ using the `Gym` library [70], which provides a standard interface for reinforcement learning environments. The environment is defined as a class that inherits from the `gym.Env` class and implements the required methods: `reset`, `step` and `render`.

3.3 Deep Reinforcement Learning Algorithms

The proposed solution is based on the DRL framework, which allows the agent to learn the optimal policy directly from the environment by taking actions and receiving rewards. The algorithms used are:

- Advantage Actor-Critic (A2C)
- Proximal Policy Optimisation (PPO)
- Deep Deterministic Policy Gradient (DDPG)
- Twin Delayed Deep Deterministic Policy Gradient (TD3)
- Soft Actor-Critic (SAC)

The implementation is done using the `Stable Baselines3` library [71], which provides a set of state-of-the-art DRL algorithms with a consistent interface and easy-to-use API. The pseudo-code for each algorithm is provided in Appendix A.1.

¹<https://www.python.org/>

3.3.1 Hyper-parameter tuning

Hyper-parameter tuning is a crucial step in the training process of DRL models, as the right choice of hyper-parameters can significantly impact the model’s performance. In this thesis, hyper-parameter tuning refers to the process of optimising the training parameters of the DRL algorithms to maximise their performance in the task of portfolio optimisation. The hyper-parameters tuned include the learning rate, batch size, number of training steps, and other algorithm-specific parameters, which are summarised in Appendix C.

To implement hyper-parameter tuning in Python, the **wandb** [72] library is used, which provides a simple and efficient way to track experiments, visualise results, and manage hyper-parameter sweeps. A sweep is defined as a search for hyper-parameters that optimises a cost function, in our case, the Sharpe ratio. Given that the models were implemented using the **Stable Baselines3** library, the integration with **wandb** allows for seamless tracking of hyper-parameter configurations and their corresponding performance metrics [73].

As mentioned above, sweeps can optimise a cost function to avoid naively testing every possible combination. Using **wandb**, a Bayesian optimisation approach is taken [74], which uses a probabilistic model to estimate the performance of different hyper-parameter configurations and selects the next configuration to test based on the expected improvement over the current best configuration. This allows for a more efficient search of the hyper-parameter space and reduces the number of configurations that need to be tested. Another option to reduce the time taken to find the optimal hyper-parameters is to use early termination. This method will stop a poorly performing run before it has fully completed, saving computational resources.

3.4 Post-hoc Explainability

Given the goal of improving the explainability of the DRL models, this thesis adopts explainability techniques to interpret the model’s decision-making process in a transparent manner. By using post-hoc methods, rather than modifying each model’s architecture to enhance their transparency, the proposal is model-agnostic and can be applied to any DRL model. Consequently, it combines the ability to find the most suitable architecture while maintaining the interpretability.

The explainability techniques implemented are:

- Feature Importance
- Local Interpretable Model-agnostic Explanations (LIME)
- SHapley Additive exPlanations (SHAP)

Following the work from de-la-Rica-Escudero et al. (2025) [68], the implementation of these techniques follows two directions. First, as in their paper, a surrogate model maps the state space to the action space as a proxy for the model’s decisions. The second direction is to use the LIME and SHAP techniques directly on the DRL model to interpret its decisions.

3.4.1 Surrogate Model Explainability

The surrogate model is trained to approximate the behaviour of the DRL model by learning the mapping from its inputs, the environment representation, to its outputs, the portfolio weights. In the paper [68], the authors do not explicitly acknowledge the use of a surrogate model, even though their code implementation does so. A potential reason behind not explaining the model’s actions directly could be the code complexity in using SHAP with a DL model.

Given that the action space is continuous, the surrogate model is implemented using a `RandomForestRegressor` [75], which is a non-parametric model that can capture complex relationships between the inputs and outputs. The model is trained on the state-actions pairs of the test data, which is the object of the explanations. However, since the model has a number of hyper-parameters, it requires careful tuning to achieve optimal performance. Consequently, hyperparameter tuning was used to find the optimal architecture using `HalvingGridSearchCV` [76], which is a method that iteratively narrows down the search space by evaluating a subset of hyper-parameters and discarding the less promising ones. The use of grid search rather than Bayesian Optimisation, as was done for DRL hyper-parameter tuning 3.3.1, is to replicate the approach taken in [68]. Once the optimal hyper-parameters have been found and the surrogate model has been trained, its prediction function is used as a proxy to interpret the original model’s decisions.

Firstly, feature importance is built-in for Random Forest Regressors, and can be easily accessed with the built-in property `feature_importances_` [77]. This method provides the importance of each feature in the state representation by using a combination of the fraction of the samples a feature contributes to and the mean decrease in impurity.

Secondly, LIME and SHAP are applied to the surrogate model via the `predict` function to provide local explanations for individual predictions. Both of these techniques provide insights into the model’s decisions by perturbing the input data and observing the changes in the output.

The LIME implementation is done using the `LimeTabularExplainer` [78], which is designed to work with tabular data and provides a way to explain individual predictions by approximating the model’s behaviour locally with a linear model. Similarly, for SHAP, the framework provides a particular implementation for tree-based models, which is used to compute the SHAP values efficiently by exploiting the structure of the trees. Conse-

quently, the `TreeExplainer` [79] is used to compute the SHAP values for the surrogate model.

3.4.2 Direct Model Explainability

Undoubtedly, a surrogate model adds an additional layer of complexity and may obscure the understanding of the original model’s decisions. Therefore, the LIME and SHAP techniques are also applied directly to the prediction function of DRL model. This approach allows for a more direct interpretation of the model’s decision-making process, without the need of a supplementary level.

For LIME, the implementation is again done using `LimeTabularExplainer`, but the prediction function is now obtained from the relevant DRL algorithm. For SHAP, the `KernelExplainer` is used, which is a model-agnostic method that randomly samples feature coalitions to approximate SHAP values to reduce computation [80].

Chapter 4

Results

This chapter presents the results of conducting experiments under the methodology proposed in Chapter 3. The experiments were designed to evaluate the performance of the implemented DRL models for portfolio optimisation in changing environment representations and market conditions. Moreover, to analyse the interpretability of the model’s decisions, a framework using post-hoc explainability techniques is explored.

4.1 Dataset

Given the general difficulty in finding the appropriate DRL algorithm with a suitable reward function for portfolio optimisation, the five implemented algorithms were tested on five different datasets. Each dataset consists of a different set of financial assets, ranging from three different asset classes. First, three datasets were constructed using the stock constituents of three renowned indexes:

- Dow Jones Industrial Average (DJIA) with 30 stocks,
- Euro Stoxx 50 Index (Euro Stoxx 50) with 50 stocks,

- Financial Times Stock Exchange 100 Index (FTSE 100) with 100 stocks.

The constituents of each of the indexes were retrieved in April 2025 and can be found in Appendix D.1. It is important to note that the datasets were chosen to illustrate different currencies, as this introduces another factor of changing market conditions.

Additionally, two datasets were constructed using commodities and currencies, respectively. The commodities dataset includes six different commodities, which are listed in Appendix D.2. These are a sample of the most traded commodities in the market and were chosen by their availability in the `Yahoo! Finance API`¹. With regard to the currencies dataset, it includes ten different currency pairs, listed in Appendix D.3. These were selected based on their trading volume and liquidity, with all pairs quoted in United States Dollar (USD).

The datasets are constructed using daily data from January 2016 to July 2025 downloaded using the Python `yfinance` library [81]. The dataset is partitioned into two disjoint sets: training and testing, with the training set containing data from January 2016 to December 2023, and the testing set starting on January 2024 until July 2025. The training set is used to train the DRL models, while the testing set is used to evaluate their performance. For hyper-parameter tuning, the training set is further split into a training and validation set, with the validation set corresponding to the period between January 2023 and December 2023. The validation set is used to evaluate the performance of the models for each hyper-parameter combination. The train-validation-test split is summarised in Table 4.1.

Dataset	Training Period	Validation Period	Testing Period
Dow Jones 30	Jan 2016 - Dec 2022	Jan 2023 - Dec 2023	Jan 2024 - Jul 2025
Euro Stoxx 50	Jan 2016 - Dec 2023	Not applicable	Jan 2024 - Jul 2025

¹<https://uk.finance.yahoo.com>

Dataset	Training Period	Validation Period	Testing Period
FTSE 100	Jan 2016 - Dec 2023	Not applicable	Jan 2024 - Jul 2025
Commodities	Jan 2016 - Dec 2022	Jan 2023 - Dec 2023	Jan 2024 - Jul 2025
Currencies	Jan 2016 - Dec 2022	Jan 2023 - Dec 2023	Jan 2024 - Jul 2025

Table 4.1: Train-Validation-Test Split for each dataset

4.2 Experiment Design

To address the challenge of finding a suitable algorithm for portfolio optimisation, the five implemented DRL algorithms were tested on the five datasets described in Section 4.1, with the goal of evaluating the performance of each algorithm in different scenarios and market conditions. Moreover, the environment representation will also be varied to assess the impact of more information on the model’s performance. Four environment representations were considered, each with a different number of features:

- Simple dataset: Open, High, Low, Close, Volume (OHLCV) prices of the assets.
- Covariance dataset: To the simple dataset, the covariance matrix of the assets is added to explicitly model the relationships between the assets.
- Indicators dataset: Technical and macroeconomic indicators are added to the simple dataset.
- Complete dataset: The complete dataset includes the simple dataset, the covariance matrix and the technical and macroeconomic indicators.

The strength of DRL algorithms lies in their ability to learn from high-dimensional data, which is why the goal is to evaluate whether a more exhaustive environment

representation leads to better performance. However, with higher dimensionality comes a higher computational cost.

Another particular challenge is the choice of suitable reward function, which is crucial for the success of the algorithms in any RL setting. The reward function should ideally be designed to encourage the model to learn an investment strategy that maximises returns while minimising risk. As a result, two choices of reward function were considered:

- Change in portfolio value: The reward is the change in portfolio value at each time step, which encourages the model to maximise returns.
- Sharpe ratio: The reward is the Sharpe ratio of the portfolio, which measures the risk-adjusted return.

Finally, the performance of the algorithms is closely related to the choice of hyper-parameters. As such, the hyper-parameters were tuned using Bayesian optimisation, to find the optimal configuration for each algorithm and dataset combination.

Overall, the experiments were designed to evaluate the performance of the implemented DRL algorithms in different scenarios, with the goal of finding the most suitable algorithm for portfolio optimisation. However, testing 5 algorithms on 5 distinct datasets with 4 possible environment representations and 2 potential reward functions would result in a total of 40 different experiments. Additionally, optimising the parameters for each experiment further expands the experimental space and significantly increases the computational time required. Due to limited computational resources², the scope of experiments was adjusted as follows:

- Hyper-parameter tuning was performed only for the Dow Jones 30 dataset with simple and indicators environment representation, as it is the smallest equities

²The university did not provide access to a computing cluster; therefore, all experiments were conducted on a personal computer.

dataset and requires less computational time.

- Since the covariance matrix, significantly increases the dimensionality of the environment representation, it was only included in the experiments with the Dow Jones 30, the currencies and the commodities datasets.
- The reward function was set to the change in portfolio value for all experiments, as it is the most straightforward and intuitive choice for portfolio optimisation. However, the Sharpe ratio performance was evaluated in the Dow Jones 30 dataset, with environment representations that exclude the covariance.

The final experimental design consists of 18 experiments, which are summarised in Table 4.2, where each row represents a unique combination of dataset, environment representation, reward function and whether hyper-parameter tuning is performed for this combination.

Dataset	Environment Representation	Reward Function	Hyper-parameter Tuning
Dow Jones 30	Simple	Change in Value	Yes
Dow Jones 30	Covariance	Change in Value	No
Dow Jones 30	Indicators	Change in Value	Yes
Dow Jones 30	Complete	Change in Value	No
Dow Jones 30	Simple	Sharpe Ratio	Yes
Dow Jones 30	Indicators	Sharpe Ratio	Yes
Euro Stoxx 50	Simple	Change in Value	No
Euro Stoxx 50	Indicators	Change in Value	No
FTSE 100	Simple	Change in Value	No
FTSE 100	Indicators	Change in Value	No

Dataset	Environment Representation	Reward Function	Hyper-parameter Tuning
Commodities	Simple	Change in Value	No
Commodities	Covariance	Change in Value	No
Commodities	Indicators	Change in Value	No
Commodities	Complete	Change in Value	No
Currencies	Simple	Change in Value	No
Currencies	Covariance	Change in Value	No
Currencies	Indicators	Change in Value	No
Currencies	Complete	Change in Value	No

Table 4.2: Summary of experiments conducted.

4.3 Evaluation

As outlined in the previous section, the experiments are designed to evaluate the performance of the implemented DRL algorithms in different scenarios and market conditions. The evaluation will focus on key performance metrics, as well as benchmarking against traditional portfolio optimisation techniques.

4.3.1 Performance Metrics

The performance metrics are provided through the `pyfolio` library [82], which includes a `perf_stats` method to calculate various performance metrics of a strategy.

The main metrics for comparison are:

- The cumulative return is the total change in investment price over a period of time,

representing the overall percentage gain or loss from the initial investment value. The formula is given by:

$$\text{Cumulative return} = \frac{\text{Final portfolio value} - \text{Initial portfolio value}}{\text{Initial portfolio value}}. \quad (4.1)$$

- The annualised return is the geometric average of the amount of money earned by an investment each year over a given period of time, providing a standardised measure of annual performance. It is calculated as follows:

$$\text{Annualised Return} = \left(\frac{\text{Final portfolio value}}{\text{Initial portfolio value}} \right)^{\frac{1}{\text{Number of years}}} - 1. \quad (4.2)$$

- The annualised volatility is the standard deviation of returns annualised to provide a measure of investment risk on a yearly basis and can be computed with the following formula:

$$\text{Annualised Volatility} = \text{Standard Deviation of Returns} \times \sqrt{\text{Yearly trading days}}, \quad (4.3)$$

where the number of trading days per year is typically assumed to be 252.

- The Sharpe ratio is a measure of risk-adjusted performance that compares the excess return of an investment to a risk-free asset against its volatility. The ratio is given by:

$$\text{Sharpe Ratio} = \frac{R_p - R_f}{\sigma_p}, \quad (4.4)$$

where R_p is the annualised return of the portfolio, R_f is the annualised risk-free rate, and σ_p is the annualised volatility of the portfolio.

- The max drawdown is the maximum percentage loss from a peak to a trough during a specified period, indicating the worst-case scenario for portfolio decline.

Its formula is:

$$\text{Max Drawdown} = \frac{\text{Peak Value} - \text{Trough Value}}{\text{Peak Value}}. \quad (4.5)$$

Other metrics available through the `pyfolio` library are outlined in Appendix E.

4.3.2 Benchmark Strategies

Aside from computing relevant performance metrics, the algorithms will be benchmarked against traditional portfolio optimisation methods. The benchmarks are designed to provide a baseline for comparison and to evaluate the performance of the DRL algorithms in relation to established methods. The following benchmark strategies were considered.

- Equal-weighted portfolio: A simple strategy that allocates an equal weight to each asset in the portfolio.
- Mean-variance optimisation: A classic portfolio optimisation method that aims to maximise Sharpe ratio.
- Min-variance portfolio: A classic portfolio optimisation method that seeks to minimise the portfolio's volatility.
- Momentum portfolio: A strategy that invests in assets with positive momentum, i.e. those that have performed well in the previous time step, and avoids those with negative momentum.

The implementation of the mean-variance and the min-variance portfolio allocation strategies has been done using the `PyPortfolioOpt` Python library [83], whereas the equal-weighted and momentum strategies have been implemented using custom code.

Finally, if the portfolio is made up of equities of a relevant index, the benchmark will

also include the index itself, which serves as a reference point for the performance of the portfolio.

Chapter 5

Legal, Social, Ethical and Professional Issues

The development and deployment of Deep Reinforcement Learning (DRL) algorithms to perform profitable portfolio allocation raises several legal, social, ethical and professional issues that must be considered. The incorporation of post-hoc explainable techniques to understand, interpret and clarify the decision-making process of these algorithms is a crucial step towards addressing these concerns. This chapter explores them in detail, referencing relevant regulatory frameworks and professional codes of conducts, and highlights potential mitigation strategies.

5.1 Legal Issues

When deploying Machine Learning (ML) algorithms in finance, the European Union (EU)'s Artificial Intelligence Act (AI Act) [65] classifies them as high-risk under Article 6 due to their potential impact on an individual's economic well-being. Consequently, high-risk systems must comply with Chapter 3 - Section 2, Articles 9-15, which include

the obligations for risk management, data governance, documentation, transparency and human supervision. In particular, Article 13 requires that high-risk AI systems's outputs are interpretable. The use of explainability techniques ensures that the decision-making process of the algorithm can be understood and justified, thus adhering to the legal requirements set forth by the AI Act.

In addition, algorithmic trading systems must adhere to the European Union directive on Markets in Financial Instruments Directive II (MiFID II) [84] and the United Kingdom (UK) Financial Conduct Authority (FCA) guidance on *Algorithmic Trading Compliance in Wholesale Markets*. These regulations require robust governance and oversight frameworks, risk controls and thorough testing infrastructure. In this project, back-testing has been carried out to assess the algorithm's performance under varying market conditions followed by explainability techniques to enable the auditability of the system.

5.2 Social Issues

Despite the growing exposure to Artificial Intelligence (AI) systems with the rise of Large Language Models (LLMs), there is still a significant lack of understanding and trust in these systems. Some of the reasons behind this mistrust include:

- the black box nature of models, making outputs hard to interpret,
- lack of human-like qualities in the models, and
- perceived limitations to adapt to new situations and learn from previous mistakes.

In the context of this work, the integration of explainability techniques addresses the black box behaviour of the implemented technology by offering clear, data-driven justifications. However, transparency must also be accessible. It is essential to provide

a user-friendly interface to easily access the explanations and generate plain language descriptions for non-technical audiences.

5.3 Ethical Issues

The ethical implications of DRL in portfolio optimisation are multifaceted. The primary concern is the potential for these systems to make decisions that may not align with human values or ethical standards. For instance, if the algorithm prioritises profit maximisation without considering the social or environmental impact of its investment choices, it could lead to unethical outcomes, such as supporting companies with poor labour practices or those contributing to environmental degradation. To address this concern, the following practices can be put in place:

1. the end-user should have full control over the assets included in their portfolio, allowing them to exclude companies that do not meet their ethical standards; or
2. the environment’s representation can be expanded to include ethical dimensions, like social impact or environmental sustainability metrics.

Regarding the environment representation, it was not possible to incorporate such data as it is not readily available nor in a standardised format. This is why, although there was an interest to broaden the environment representation to include other sources outside of the financial domain, it was not feasible to do so given the available resources.

A critical point in portfolio allocation is risk management. One of the reward functions used in this research balances return maximisation with risk minimisation, ensuring that the agent learns to avoid overly risky investments. Additionally, the inclusion of a risk-free asset enables a safe-guard mechanism. The agent can be implemented to allocate all the resources to cash when market volatility exceeds a pre-defined threshold and resume

only when said volatility decreases.

Finally, potential biases in the training data or model structure could result in allocations favouring certain industries or assets classes. This thesis mitigates such biases by constructing a portfolio from diversified indices (Appendix D.1) and using explainability techniques to audit the decisions for systematic biases.

5.4 Professional Issues

This work has been conducted in accordance with the British Computer Society (BCS) Code of Conduct [85] and the The Institution of Engineering and Technology (IET) Rules of Conduct [86]. All the work in this report is original unless stated otherwise, and all external contributions are explicitly acknowledged, following the principles of the BCS Code of Conduct [85]. Moreover, all third-party libraries, datasets and code have been explicitly acknowledged and their use complies with the respective licences and guidelines.

Most importantly, the project has been conducted within my own area of expertise: computational finance, Reinforcement Learning, and Explainable Artificial Intelligence. If any aspects of the project were beyond these, proper academic sources have been consulted to ensure the integrity and quality of the work.

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] C. Molnar, *LIME*, ch. 14. Christoph Molnar, 3 ed., 2025.
- [4] Wikipedia, “Dow jones industrial average - wikipedia,” 7 2025.
- [5] Wikipedia, “Euro stoxx 50 - wikipedia,” 7 2025.
- [6] Wikipedia, “Ftse 100 index - wikipedia,” 8 2025.
- [7] 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.
- [8] 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.
- [9] 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.

- [10] 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.
- [11] 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.
- [12] 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.
- [13] J. Moody and M. Saffell, “Learning to trade via direct reinforcement,” *IEEE Transactions on Neural Networks*, vol. 12, pp. 875–889, 7 2001.
- [14] 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.
- [15] 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.
- [16] 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.
- [17] 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.

- [18] 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.
- [19] D. Gunning, M. Stefik, J. Choi, T. Miller, S. Stumpf, and G. Z. Yang, “Xai—explainable artificial intelligence,” *Science Robotics*, vol. 4, 12 2019.
- [20] 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.
- [21] 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.
- [22] Y. Sato, “Model-free reinforcement learning for financial portfolios: A brief survey,” *arXiv preprint arXiv:1904.04973*, 4 2019.
- [23] B. Bruce and J. Greene, *Chapter 4 - Portfolio Construction*, pp. 133–178. Academic Press, 2014.
- [24] 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.
- [25] H. Markowitz, “Portfolio selection,” *The Journal of Finance*, vol. 7, pp. 77–91, 3 1952.
- [26] F. M. Clinic, “Mean variance portfolio theory.”
- [27] IBM, “What is machine learning (ml)?,” 9 2021.

- [28] 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.
- [29] I. Goodfellow, Y. Bengio, and A. Courville, *Deep Learning*. MIT Press, 2016. <http://www.deeplearningbook.org>.
- [30] 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.
- [31] 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.
- [32] 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.
- [33] R. Bellman, “A markovian decision process,” *Journal of Mathematics and Mechanics*, vol. 6, pp. 679–684, 1957.
- [34] R. Bellman, *Dynamic Programming*. Princeton University Press, 1 1957.
- [35] S. Thrun, “Efficient exploration in reinforcement learning,” tech. rep., Carnegie Mellon University, 1 1992.
- [36] 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.

- [37] 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.
- [38] 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.
- [39] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. K. Openai, “Proximal policy optimization algorithms,” *arXiv preprint arXiv:1707.06347*, 7 2017.
- [40] 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.
- [41] 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.
- [42] 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.
- [43] D. H. G. B. Tokyo and J. Schmidhuber, “Recurrent world models facilitate policy evolution,” *Advances in Neural Information Processing Systems*, vol. 31, 2018.
- [44] 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.
- [45] 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.
- [46] 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.
- [47] G. E. Uhlenbeck and L. S. Ornstein, “On the theory of the brownian motion,” *Physical Review*, vol. 36, p. 823, 9 1930.
- [48] P. J. Phillips, C. A. Hahn, P. C. Fontana, A. N. Yates, K. Greene, D. A. Broniatowski, and M. A. Przybocki, “Four principles of explainable artificial intelligence,” *National Institute of Standards and Technology*, vol. Internal Report 8312, 9 2021.
- [49] P. Sequeira and M. Gervasio, “Interestingness elements for explainable reinforcement learning: Understanding agents’ capabilities and limitations,” *Artificial Intelligence*, vol. 288, p. 103367, 11 2020.
- [50] S. Greydanus, A. Koul, J. Dodge, and A. Fern, “Visualizing and understanding atari agents,” *Proceedings of the 35th International Conference on Machine Learning*, vol. 80, pp. 1792–1801, 7 2018.
- [51] M. T. Ribeiro, S. Singh, and C. Guestrin, ““why should i trust you?”: Explaining the predictions of any classifier,” *Proceedings of the ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pp. 1135–1144, 8 2016.

- [52] S. M. Lundberg and S. I. Lee, “A unified approach to interpreting model predictions,” *Advances in Neural Information Processing Systems*, vol. 2017-December, pp. 4766–4775, 5 2017.
- [53] B. Amirshahi and S. Lahmire, “Investigating the effectiveness of twitter sentiment in cryptocurrency close price prediction by using deep learning,” *Expert Systems*, vol. 42, p. e13428, 8 2023.
- [54] L. Breiman, “Random forests,” *Machine Learning*, vol. 45, pp. 5–32, 10 2001.
- [55] S. Nembrini, I. R. König, and M. N. Wright, “The revival of the gini importance?,” *Bioinformatics (Oxford, England)*, vol. 34, pp. 3711–3718, 11 2018.
- [56] L. Shapley, *A value for n-person games*, pp. 307–317. Princeton University Press, 1953.
- [57] E. Štrumbelj and I. Kononenko, “A general method for visualizing and explaining black-box regression models,” *Lecture Notes in Computer Science (including sub-series Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*, vol. 6594 LNCS, pp. 21–30, 2011.
- [58] S. M. Lundberg, G. G. Erion, and S.-I. Lee, “Consistent individualized feature attribution for tree ensembles,” *Advances in Neural Information Processing Systems*, 2019.
- [59] C. Zhang, N. N. A. Sjarif, and R. Ibrahim, “Deep learning models for price forecasting of financial time series: A review of recent advancements: 2020–2022,” *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery*, vol. 14, p. e1519, 9 2023.
- [60] A. Millea and F. Guijarro, “Deep reinforcement learning for trading—a critical survey,” *Data 2021, Vol. 6, Page 119*, vol. 6, p. 119, 11 2021.

- [61] R. Chaudhary, “Advanced stock market prediction using long short-term memory networks: A comprehensive deep learning framework,” *Journal of Financial Data Science*, 5 2025.
- [62] X.-Y. Liu, “Finrl: A deep reinforcement learning library for automated stock trading in quantitative finance,” *SSRN Electronic Journal*, 11 2020.
- [63] Z. Jiang and J. Liang, “Cryptocurrency portfolio management with deep reinforcement learning,” *2017 Intelligent Systems Conference, IntelliSys 2017*, vol. 2018-January, pp. 905–913, 12 2016.
- [64] Z. Zhang, S. Zohren, and S. Roberts, “Deep reinforcement learning for trading,” *Papers*, 2019.
- [65] E. Parliament and C. of the European Union, “Regulation (eu) 2024/1689 of the european parliament and of the council of 13 june 2024 laying down harmonised rules on artificial intelligence and amending certain union legislative acts (artificial intelligence act),” 2024. Official Journal of the European Union, L 2024/1689, 12.7.2024.
- [66] Mandeep, A. Agarwal, A. Bhatia, A. Malhi, P. Kaler, and H. S. Pannu, “Machine learning based explainable financial forecasting,” *2022 4th International Conference on Computer Communication and the Internet, ICCCI 2022*, pp. 34–38, 2022.
- [67] C. A. Zhang, S. Cho, and M. Vasarhelyi, “Explainable artificial intelligence (xai) in auditing,” *International Journal of Accounting Information Systems*, vol. 46, p. 100572, 9 2022.
- [68] A. de La-Rica-Escudero, E. C. Garrido-Merchán, and M. Coronado-Vaca, “Explainable post hoc portfolio management financial policy of a deep reinforcement learning agent,” *PLOS ONE*, vol. 20, p. e0315528, 1 2025.

- [69] W. Sharpe, “The sharpe ratio,” *The Journal of Portfolio Management*, vol. Fall, 1994.
- [70] G. Brockman, V. Cheung, L. Pettersson, J. Schneider, J. Schulman, J. Tang, and W. Z. Openai, “Openai gym,” *arXiv preprint arXiv:1606.01540*, 6 2016.
- [71] A. Raffin, A. Hill, A. Gleave, A. Kanervisto, M. Ernestus, and N. Dormann, “Stable-baselines3: Reliable reinforcement learning implementations,” *Journal of Machine Learning Research*, vol. 22, pp. 1–8, 2021.
- [72] L. Biewald, “Experiment tracking with weights and biases,” 2020. Software available from wandb.com.
- [73] W. . Biases, “Stable baselines 3 — weights & biases documentation,” 8 2025.
- [74] S. Falkner, A. Klein, and F. Hutter, “Bohb: Robust and efficient hyperparameter optimization at scale,” *35th International Conference on Machine Learning, ICML 2018*, vol. 4, pp. 2323–2341, 7 2018.
- [75] S. Learn, “Randomforestregressor — scikit-learn documentation.”
- [76] S. Learn, “Halvinggridsearchcv — scikit-learn 1.7.1 documentation.”
- [77] S. Learn, “1.11.2.5. feature importance evaluation — scikit-learn 1.7.1 documentation.”
- [78] Lime, “Lime tabular explainer — lime 0.1 documentation.”
- [79] SHAP, “shap.treeexplainer — shap documentation.”
- [80] SHAP, “shap.kernelexplainer — shap documentation.”
- [81] ranaroussi, “finance - python library.”
- [82] Q. Inc, “pyfolio.”

- [83] R. A. Martin, “Pyportfoliopt: portfolio optimization in python,” *Journal of Open Source Software*, vol. 6, no. 61, p. 3066, 2021.
- [84] European Parliament and Council of the European Union, “Directive 2014/65/eu of the european parliament and of the council of 15 may 2014 on markets in financial instruments and amending directive 2002/92/ec and directive 2011/61/eu (recast),” 2014. Official Journal of the European Union, L 173, 12.6.2014, p. 349-496.
- [85] B. C. Society, “Bcs code of conduct,” 2022.
- [86] I. of Engineering and Technology, “Rules of conduct,” 2024.

Appendix A

Algorithms

A.1 Deep Reinforcement Learning Algorithms

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

Initialise:Global shared policy parameters θ and value parameters θ_v Number of parallel workers N Global step counter $T \leftarrow 0$ Hyper-parameters: discount γ , max steps per update t_{\max} , max total steps T_{\max} , learning rates α_π, α_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. θ :

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

Accumulate gradients w.r.t. θ_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 Optimisation (PPO) Pseudo-code

Initialise:Policy parameters θ_0 and value function parameters ϕ_0 Global step counter $T \leftarrow 0$ Hyper-parameters: discount γ , GAE parameter λ , clipping parameter ϵ , learning rates α_π, α_v , epochs per update K_{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 π_{θ_k} in environmentStore transitions: $\{(s_t, a_t, r_t, s_{t+1}, \text{done}_t)\}$ **for** each trajectory τ 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_{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 θ^Q, θ^μ

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

Replay buffer \mathcal{B}

Hyper-parameters: discount γ , soft update rate τ , batch size N , exploration noise process \mathcal{N} , learning rates α_Q, α_μ , 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 θ_1, θ_2

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

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

Replay buffer \mathcal{B}

Hyper-parameters: discount γ , target policy noise $\tilde{\sigma}$, noise clip c , policy delay d , exploration noise σ , update rate τ , 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 θ Two critic networks $Q_{\phi_1}(s, a)$ and $Q_{\phi_2}(s, a)$ with parameters ϕ_1, ϕ_2 Target critic networks: $\phi'_1 \leftarrow \phi_1, \phi'_2 \leftarrow \phi_2$ Replay buffer \mathcal{D} Hyper-parameters: discount γ , temperature α (fixed or learnable), target entropy \mathcal{H} (if α is learnable), batch size N , learning rates $\lambda_Q, \lambda_\pi, \lambda_\alpha$, soft update rate τ **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** α 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**

A.2 Explainability Algorithms

Algorithm 6 Shapley Value Approximation

Require: Number of iterations M , instance of interest \mathbf{x} , feature index j , data matrix \mathbf{X} , model \hat{f}

Ensure: Estimated Shapley value $\phi_j(\mathbf{x})$

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

 Draw a random instance \mathbf{z} from data matrix \mathbf{X}

 Choose a random permutation \mathbf{o} of the feature indices

 Order \mathbf{x} according to \mathbf{o} : $\mathbf{x}_{\mathbf{o}} = (x_{(1)}, \dots, x_{(j)}, \dots, x_{(p)})$

 Order \mathbf{z} according to \mathbf{o} : $\mathbf{z}_{\mathbf{o}} = (z_{(1)}, \dots, z_{(j)}, \dots, z_{(p)})$

 Construct two new instances:

$\mathbf{x}_{+j} = (x_{(1)}, \dots, x_{(j-1)}, x_{(j)}, z_{(j+1)}, \dots, z_{(p)})$

$\mathbf{x}_{-j} = (x_{(1)}, \dots, x_{(j-1)}, z_{(j)}, z_{(j+1)}, \dots, z_{(p)})$

 Compute marginal contribution:

$$\phi_j^{(m)} = \hat{f}(\mathbf{x}_{+j}) - \hat{f}(\mathbf{x}_{-j})$$

end for

Compute average Shapley value:

$$\phi_j(\mathbf{x}) = \frac{1}{M} \sum_{m=1}^M \phi_j^{(m)}$$

Appendix B

State Representation

B.1 Technical Indicators

The following technical indicators are used to represent the state of the financial environment. They are calculated based on the historical price data of the assets in the portfolio.

- Simple Moving Average (SMA): Lagging indicator that smooths the price over a period of time. It is the unweighted mean of the previous k data points and is computed following:

$$\text{SMA}_t^k = \frac{1}{k} \sum_{i=t-k+1}^t p_i, \quad (\text{B.1})$$

where t is the current time step, k is the look-back period and p_i is the close price at time i . In this thesis, the SMA is calculated for 5, 10 and 20 days.

- Exponential Moving Average (EMA): Weighted moving average that gives more importance to recent prices with the goal of making it more responsive to new

information. It is calculated using the formula:

$$\text{EMA}_t^k = \begin{cases} p_t & \text{if } t = 0 \\ \frac{2}{k+1}p_t + \left(1 - \frac{2}{k+1}\right) \text{EMA}_{t-1}^k & \text{if } t > 0 \end{cases}, \quad (\text{B.2})$$

where k is the look-back period, p_t is the close price at time t and $\frac{2}{k+1}$ is the smoothing factor. In this thesis, the EMA is calculated for 5 and 10 days.

- Moving Average Convergence Divergence (MACD): Momentum indicator that shows the relationship between two moving average of an asset's price. The formula is:

$$\text{MACD}_t = \text{EMA}_t^{k_1} - \text{EMA}_t^{k_2}, \quad (\text{B.3})$$

where k_1 and k_2 are the look-back periods for the short-term (12 periods) and long-term (26 periods) EMAs, respectively, and p_t is the close price at time t .

- Relative Strength Index (RSI): Momentum indicator that measures the magnitude of recent changes to identify overbought or oversold conditions in the market. It is computed using smoothed moving averages for the upward change in closing price p_t , defined as $u_t = \max\{p_t - p_{t-1}, 0\}$:

$$\text{SMMA}_t^k(u_t) = \frac{1}{k}u_t + \left(1 - \frac{1}{k}\right) \text{SMMA}_{t-1}^k(u_{t-1}) \quad (\text{B.4})$$

and for the downward change, defined as $d_t = \max\{p_{t-1} - p_t, 0\}$:

$$\text{SMMA}_t^k(d_t) = \frac{1}{k}d_t + \left(1 - \frac{1}{k}\right) \text{SMMA}_{t-1}^k(d_{t-1}). \quad (\text{B.5})$$

Then, the Relative Strength Index (RSI) is given by:

$$\text{RSI}_t = 100 - \frac{100}{1 + \text{RS}_t}, \quad (\text{B.6})$$

where the relative strength RS_t is defined as:

$$RS_t = \frac{SMMA_t^k(u_t)}{SMMA_t^k(d_t)}. \quad (B.7)$$

- Commodity Channel Index (CCI): Momentum indicator that measures the deviation of the price from its historical average price over a period of time. It is calculated as:

$$CCI_t^k = \frac{TP_t^k - MA_t^k}{0.015 \cdot MD_t}, \quad (B.8)$$

where k is the number of periods (14 days), the typical price TP_t is:

$$TP_t = \sum_{i=t-k+1}^t \frac{p_i + h_i + l_i}{3}, \quad (B.9)$$

with p_t , h_t and l_t being the close, high and low prices at time step t , respectively, and the moving average MA_t^k is:

$$MA_t^k = \frac{1}{k} \sum_{i=t-k+1}^t TP_i \quad (B.10)$$

and the mean deviation MD_t is:

$$MD_t = \frac{1}{k} \sum_{i=t-k+1}^t |TP_i - MA_i^k|. \quad (B.11)$$

- Bollinger Bands: Volatility indicator that defines the trend-line for high and low prices based on the deviation of the asset from the moving average. The upper bands is calculated as:

$$BOLLUB_t^k = MA_t^k + m \cdot SD_t^k, \quad (B.12)$$

where k is the number of periods (20 days), m is the number of standard deviations

away from the moving average (2 standard deviations) and SD_t^k is the standard deviation of the typical price over the same period. Similarly, the lower band is calculated as:

$$\text{BOLLDB}_t^k = \text{MA}_t^k - m \cdot \text{SD}_t^k. \quad (\text{B.13})$$

- Average True Range (ATR): Volatility indicator that measures the average range of price movement over a period of time, usually 14 days. It is calculated as:

$$\text{ATR}_t^k = \frac{1}{k} \sum_{i=t-k+1}^t \text{TR}_i, \quad (\text{B.14})$$

where the true range TR_i is defined as:

$$\text{TR}_i = \max \{h_i - l_i, |h_i - p_{i-1}|, |l_i - p_{i-1}|\}, \quad (\text{B.15})$$

where h_i and l_i are the high and low prices of the asset at time i , respectively, and p_{i-1} is the close price of the asset at time $i - 1$. The true range finds the maximum of the following three:

- most recent period high minus most recent period low,
 - absolute value of the most recent period high minus the previous close, and
 - absolute value of the most recent period low minus the previous close.
- Average Directional Index (ADX): Trend indicator used to measure the strength of a trend by quantifying the price movement. It is calculated as:

$$\text{ADX}_t^k = \frac{1}{k} \sum_{i=t-k+1}^t \text{DX}_i, \quad (\text{B.16})$$

where the Directional Movement Index (DX) is defined as:

$$DX_t = \frac{100 \cdot |PDI_t - MDI_t|}{PDI_t + MDI_t} \quad (B.17)$$

with the Positive Directional Index (PDI) and Negative Directional Index (MDI) calculated as

$$PDI_t = \frac{100 \cdot SMMA_t^k(DM^+)}{ATR_t^k} \quad (B.18)$$

and

$$MDI_t = \frac{100 \cdot SMMA_t^k(DM^-)}{ATR_t^k}, \quad (B.19)$$

where DM^+ and DM^- are the positive and negative directional movements, respectively, calculated as:

$$DM^+ = \max(0, h_t - h_{t-1}) \quad (B.20)$$

and

$$DM^- = \max(0, l_{t-1} - l_t), \quad (B.21)$$

where h_t and l_t are the high and low prices of the asset at time t , respectively.

- Rate of Change (ROC): Momentum indicator that measures the percentage change in price between the current price p_t and the price p_{t-k} periods ago. It is given by the formula:

$$ROC_t^k = \frac{p_t - p_{t-k}}{p_{t-k}} \cdot 100. \quad (B.22)$$

In this thesis, the ROC is calculated for $k = 10$ days.

B.2 Macroeconomic Indicators

The following macroeconomic indicators are used to represent the state of the financial environment. They provide additional context about the market conditions and are calculated based on external data sources.

- Volatility Index (VIX) measures the market's expectation of future volatility based on options prices. This is only available for US markets.
 - VIX is calculated using the implied volatility of Standard and Poor's 500 Index (S&P 500) options.
 - VXD is calculated using the implied volatility of Dow Jones 30 Index (DW30) options.
- Currency index captures the impact from the monetary market on the stock market.
 - U.S. Dollar Index (DXY): United States (U.S.) dollar's value relative to a basket of foreign currencies.
 - Euro Index (EXY): Euro's value relative to a basket of foreign currencies.
 - British Pound Currency Index (BXY): British pound's value relative to a basket of foreign currencies.
- Interest rates reflect the cost of borrowing money and the return on savings. This is only available for US markets.
 - 3-Month Treasury Yield (IRX): Reflects the return on investment for a 3-month government bond.
 - 5-Year Treasury Yield (FVX): Reflects the return on investment for a 5-year government bond.

- 10-Year Treasury Yield (TNX): Reflects the return on investment for a 10-year government bond.

Appendix C

Hyper-parameter tuning

For the five implemented Deep Reinforcement Learning algorithms, the following table summarises the hyper-parameters that were tuned during the training process.

Model	Hyperparameter	Values / Range
A2C	Number of steps	{5, 10, 20, 30, 40}
	Entropy coefficient	Uniform[1×10^{-8} , 1×10^{-3}]
	Learning rate	Uniform[1×10^{-5} , 1×10^{-2}]
PPO	Number of steps	{128, 256, 512, 1024, 2048}
	Entropy coefficient	Uniform[1×10^{-8} , 1×10^{-3}]
	Learning rate	Uniform[1×10^{-5} , 1×10^{-2}]
	Batch size	{32, 64, 128, 256, 512}
DDPG	Batch size	{64, 128, 256}
	Buffer size	{50000, 100000, 200000, 500000}
	Learning rate	Uniform[1×10^{-5} , 1×10^{-2}]
TD3	Batch size	{64, 100, 128, 256}
	Buffer size	{500000, 1000000, 2000000}

Model	Hyperparameter	Values / Range
	Learning rate	Uniform[1×10^{-5} , 1×10^{-2}]
SAC	Batch size	{32, 64, 128}
	Buffer size	{100000, 500000, 1000000, 2000000}
	Learning rate	Uniform[1×10^{-5} , 1×10^{-2}]
	Learning starts	{500, 1000, 2000, 5000}
	Entropy coefficient	{"auto", "auto_0.1", "auto_0.01"}

Table C.1: Hyper-parameter tuning configurations for different RL algorithms.

Appendix D

Datasets

D.1 Equities

D.1.1 Dow Jones 30

The Dow Jones Industrial Average (DJIA) is a stock market index with 30 companies listed on United States (U.S.) stock exchanges. The trading symbol is $\hat{\text{DJI}}$ and as of April 2025, its constituents are:

Symbol	Name	Sector	Industry
AXP	American Express	Financial Services	Credit Services
AMGN	Amgen	Healthcare	Drug Manufacturers - General
AAPL	Apple	Technology	Consumer Electronics
AMZN	Amazon	Consumer Cyclical	Internet Retail
BA	Boeing	Industrials	Aerospace & Defense

Symbol	Name	Sector	Industry
CAT	Caterpillar	Industrials	Farm & Heavy Construction Machinery
CRM	Salesforce	Technology	Software - Application
CSCO	Cisco	Technology	Communication Equipment
CVX	Chevron	Energy	Oil & Gas Integrated
DIS	Walt Disney	Communication Services	Entertainment
GS	Goldman Sachs	Financial Services	Capital Markets
HD	Home Depot	Consumer Cyclical	Home Improvement Retail
HON	Honeywell	Industrials	Conglomerates
IBM	IBM	Technology	Information Technology Services
JNJ	Johnson & Johnson	Healthcare	Drug Manufacturers - General
JPM	JP Morgan Chase	Financial Services	Banks - Diversified
KO	Coca-Cola	Consumer Defensive	Beverages - Non-Alcoholic
MCD	McDonald's	Consumer Cyclical	Restaurants
MMM	3M	Industrials	Conglomerates
MRK	Merck	Healthcare	Drug Manufacturers - General
MSFT	Microsoft	Technology	Software - Infrastructure

Symbol	Name	Sector	Industry
NKE	Nike	Consumer Cyclical	Footwear & Accessories
NVDA	NVIDIA	Technology	Semiconductors
PG	Procter & Gamble	Consumer Defensive	Household & Personal Products
SHW	Sherwin-Williams	Basic Materials	Specialty Chemicals
TRV	Travelers	Financial Services	Insurance - Property & Casualty
UNH	UnitedHealth	Healthcare	Healthcare Plans
V	Visa	Financial Services	Credit Services
VZ	Verizon	Communication Services	Telecom Services
WMT	Walmart	Consumer Defensive	Discount Stores

Table D.1: Constituents of the Dow Jones Industrial Average index as of April 2025. [4]

D.1.2 Euro Stoxx 50

The Euro Stoxx 50 Index (Euro Stoxx 50) is a stock market index that represents 50 of the largest companies in the Eurozone. The trading symbol is $\hat{\text{STOXX50E}}$ and as of April 2025, its constituents are:

Symbol	Name	Sector	Industry
ADS.DE	adidas	Consumer Cyclical	Footwear & Accessories
ADYEN.AS	Adyen	Technology	Software - Infrastructure

Symbol	Name	Sector	Industry
AD.AS	Koninklijke Ahold Delhaize	Consumer Defensive	Grocery Stores
AI.PA	L'Air Liquide	Basic Materials	Specialty Chemicals
AIR.PA	Airbus	Industrials	Aerospace & Defense
ALV.DE	Allianz	Financial Services	Insurance - Diversified
AB1.BR	Anheuser-Busch InBev	Consumer Defensive	Beverages - Brewers
ASML.AS	ASML Holding	Technology	Semiconductor Equip- ment & Materials
CS.PA	AXA	Financial Services	Insurance - Diversified
BAS.DE	BASF	Basic Materials	Chemicals
BAYN.DE	Bayer	Healthcare	Drug Manufacturers - General
BBVA.MC	BBVA	Financial Services	Banks - Diversified
SAN.MC	Banco Santander	Financial Services	Banks - Diversified
BMW.DE	BMW	Consumer Cyclical	Auto Manufacturers
BNP.PA	BNP Paribas	Financial Services	Banks - Regional
BN.PA	Danone	Consumer Defensive	Packaged Foods
DB1.DE	Deutsche Börse	Financial Services	Financial Data & Stock Exchanges
DHL.DE	Deutsche Post	Industrials	Integrated Freight & Logistics
DTE.DE	Deutsche Telekom	Communication Ser- vices	Telecom Services
ENEL.MI	Enel	Utilities	Utilities - Diversified

Symbol	Name	Sector	Industry
ENI.MI	Eni	Energy	Oil & Gas Integrated
EL.PA	EssilorLuxottica	Healthcare	Medical Instruments & Supplies
RACE.MI	Ferrari	Consumer Cyclical	Auto Manufacturers
FLTR.L	Flutter Entertainment	Consumer Cyclical	Gambling
RMS.PA	Hermès International	Consumer Cyclical	Luxury Goods
IBE.MC	Iberdrola	Utilities	Utilities - Diversified
ITX.MC	Industria de Diseño Textil	Consumer Cyclical	Apparel Retail
IFX.DE	Infineon Technologies	Technology	Semiconductors
INGA.AS	ING Groep	Financial Services	Banks - Diversified
ISP.MI	Intesa Sanpaolo	Financial Services	Banks - Regional
KER.PA	Kering	Consumer Cyclical	Luxury Goods
OR.PA	L'Oréal	Consumer Defensive	Household & Personal Products
MC.PA	LVMH Moët Hennessy - Louis Vuitton	Consumer Cyclical	Luxury Goods
MBG.DE	Mercedes-Benz Group AG	Consumer Cyclical	Auto Manufacturers

Symbol	Name	Sector	Industry
MUV2.DE	Münchener Rückversicherungs- Gesellschaft	Financial Services	Insurance - Reinsurance
NOKIA.HE	Nokia	Technology	Communication Equip- ment
NDA- FI.HE	Nordea Bank	Financial Services	Banks - Regional
RI.PA	Pernod Ricard	Consumer Defensive	Beverages - Wineries & Distilleries
PRX.AS	Prosus	Communication Ser- vices	Internet Content & In- formation
SAF.PA	Safran	Industrials	Aerospace & Defense
SGO.PA	Compagnie de Saint-Gobain	Industrials	Building Products & Equipment
SAN.PA	Sanofi	Healthcare	Drug Manufacturers - General
SAP.DE	SAP	Technology	Software - Application
SU.PA	Schneider Electric	Industrials	Specialty Industrial Machinery
SIE.DE	Siemens	Industrials	Specialty Industrial Machinery
STLAM.MI	Stellantis	Consumer Cyclical	Auto Manufacturers
TTE.PA	TotalEnergies	Energy	Oil & Gas Integrated
DG.PA	Vinci	Industrials	Engineering & Con- struction

Symbol	Name	Sector	Industry
UCG.MI	UniCredit	Financial Services	Banks - Regional
VOW.DE	Volkswagen	Consumer Cyclical	Auto Manufacturers

Table D.2: Constituents of the Euro Stoxx 50 index as of April 2025. [5]

D.1.3 FTSE 100

The Financial Times Stock Exchange 100 Index (FTSE 100) is a stock market index that represents 100 of the largest companies listed on the London Stock Exchange. The trading symbol is $\hat{\text{FTSE}}$ and as of April 2025, its constituents are:

Symbol	Name	Industry
III.L	3i Group	Financial Services
ADM.L	Admiral Group	Insurance
AAF.L	Airtel Africa	Telecommunications Services
ALW.L	Alliance Witan	Investment Trusts
AAL.L	Anglo American	Mining
ANTO.L	Antofagasta	Mining
AHT.L	Ashtead Group	Support Services
ABF.L	Associated British Foods	Food & Tobacco
AZN.L	AstraZeneca	Pharmaceuticals & Biotechnology
AUTO.L	Auto Trader Group	Media
AV.L	Aviva	Life Insurance

Symbol	Name	Industry
BAB.L	Babcock International Group	Aerospace & Defence
BA.L	BAE Systems	Aerospace & Defence
BARC.L	Barclays	Banks
BTRW.L	Barratt Redrow	Household Goods & Home Construction
BEZ.L	Beazley	Insurance
BKG.L	The Berkeley Group Holdings	Household Goods & Home Construction
BP.L	BP	Oil & Gas Producers
BATS.L	British American Tobacco	Tobacco
BT-A.L	BT Group	Telecommunications Services
BNZL.L	Bunzl	Support Services
CNA.L	Centrica	Multiline Utilities
CCEP.L	Coca-Cola Europacific Partners	Beverages
CCH.L	Coca-Cola HBC	Beverages
CPG.L	Compass Group	Support Services
CTEC.L	ConvaTec Group	Health Care Equipment & Services
CRDA.L	Croda International	Chemicals
DCC.L	DCC	Support Services
DGE.L	Diageo	Beverages
DPLM.L	Diploma	Industrial Support Services
EDV.L	Endeavour Mining	Precious Metals and Mining

Symbol	Name	Industry
ENT.L	Entain	Travel & Leisure
EZJ.L	easyJet	Travel & Leisure
EXPN.L	Experian	Support Services
FCIT.L	F&C Investment Trust	Collective Investments
FRES.L	Fresnillo	Mining
GAW.L	Games Workshop Group	Leisure Goods
GLEN.L	Glencore	Mining
GSK.L	GSK	Pharmaceuticals & Biotechnology
HLN.L	Haleon	Pharmaceuticals & Biotechnology
HLMA.L	Halma	Electronic Equipment & Parts
HIK.L	Hikma Pharmaceuticals	Pharmaceuticals & Biotechnology
HSX.L	Hiscox	Non-life insurance
HWDN.L	Howden Joinery Group	Homebuilding & Construction Sup- plies
HSBA.L	HSBC Holdings	Banks
IHG.L	InterContinental Hotels Group	Travel & Leisure
IMI.L	IMI	Industrial Engineering
IMB.L	Imperial Brands	Tobacco
INF.L	Informa	Media
ICG.L	ICG	Financial Services
IAG.L	International Consoli- dated Airlines Group	Travel & Leisure
ITRK.L	Intertek Group	Support Services

Symbol	Name	Industry
JD.L	JD Sports Fashion	General Retailers
KGF.L	Kingfisher	Retailers
LAND.L	Land Securities	Real Estate Investment Trusts
LGEN.L	Legal & General	Life Insurance
LLOY.L	Lloyds Banking	Banks
LMP.L	LondonMetric Property	Real Estate Investment Trusts
LSEG.L	London Stock Exchange Group	Financial Services
MNG.L	M&G	Financial Services
MKS.L	Marks and Spencer	Food & Drug Retailing
MRO.L	Melrose Industries	Aerospace & Defence
MNDI.L	Mondi	Containers & Packaging
NG.L	National Grid	Multiline Utilities
NWG.L	NatWest	Banks
NXT.L	NEXT	General Retailers
PERSON.L	Pearson	Media
PSH.L	Pershing Square Hold- ings	Financial Services
PSN.L	Persimmon	Household Goods & Home Con- struction
PHNX.L	Phoenix Group	Life Insurance
PCT.L	Polar Capital Technol- ogy	Investment Trusts
PRU.L	Prudential	Life Insurance

Symbol	Name	Industry
RKT.L	Reckitt Benckiser	Household Goods & Home Construction
REL.L	RELX	Media
RTO.L	Rentokil Initial	Support Services
RMV.L	Rightmove	Media
RIO.L	Rio Tinto	Mining
RR.L	Rolls-Royce Holdings	Aerospace & Defence
SGE.L	The Sage Group	Software & Computer Services
SBRY.L	J Sainsbury	Food & Drug Retailing
SDR.L	Schroders	Financial Services
SMT.L	Scottish Mortgage	Collective Investments
SGRO.L	SEGRO	Real Estate Investment Trusts
SVT.L	Severn Trent	Multiline Utilities
SHEL.L	Shell	Oil & Gas Producers
SMIN.L	Smiths Group	Industrial Engineering
SN.L	Smith & Nephew	Health Care Equipment & Services
SPX.L	Spirax Group	Industrial Engineering
SSE.L	SSE	Electrical Utilities & Independent Power Producers
STAN.L	Standard Chartered	Banks
STJ.L	St. James's Place	Financial Services
TW.L	Taylor Wimpey	Household Goods & Home Construction
TSCO.L	Tesco	Food & Drug Retailing
ULVR.L	Unilever	Personal Goods

Symbol	Name	Industry
UU.L	United Utilities	Multiline Utilities
UTG.L	Unite Group	Real Estate Investment Trusts
VOD.L	Vodafone Group	Mobile Telecommunications
WEIR.L	The Weir Group	Industrial Goods and Services
WTB.L	Whitbread	Retail Hospitality
WPP.L	WPP	Media

Table D.3: Constituents of the FTSE100 index as of April 2025. [6]

D.2 Commodities

The commodities market includes a variety of physical goods that are traded on exchanges. The following table lists a sample of 6 commodities:

Symbol	Name
CL=F	Crude Oil
NG=F	Natural Gas
GC=F	Gold
SI=F	Silver
ALI=F	Aluminum
HG=F	Copper

Table D.4: 6 Commodities Futures Contracts

D.3 Currencies

The currencies market includes a variety of currency pairs that are traded on exchanges.

The following table lists a sample of 10 currency pairs, selected for their trading volume.

Note that the trading symbols all have United States Dollar (USD) as the quote currency.

Symbol	Name
EURUSD=X	Euro (EUR)/USD
GBPUSD=X	British Pound (GBP)/USD
JPYUSD=X	Japanese Yen (JPY)/USD
AUDUSD=X	Australian Dollar (AUD)/USD
CADUSD=X	Canadian Dollar (CAD)/USD
CNYUSD=X	Chinese Yuan (CNY)/USD
CHFUSD=X	Swiss Franc (CHF)/USD
HKDUSD=X	Hong Kong Dollar (HKD)/USD
KRWUSD=X	South Korean Won (KRW)/USD
INRUSD=X	Indian Rupee (INR)/USD

Table D.5: 10 Currency Pairs

Appendix E

Evaluation Metrics

The following list outlines additional evaluation metrics that can be used to assess the performance of trading strategies.

- Calmar ratio: A risk-adjusted performance measure that compares the average annual compounded rate of return to the maximum drawdown, providing insight into how much return an investor receives for each unit of risk taken.

$$\text{Calmar Ratio} = \frac{\text{Average Annualised Return}}{\text{Max Drawdown}} \quad (\text{E.1})$$

- Stability: A measure that determines the R-square of a linear fit to the cumulative log returns, to indicate how well returns follow a linear trend over time.

$$\text{Stability} = R^2 = 1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2} \quad (\text{E.2})$$

where y_i is the observed cumulative log return, \hat{y}_i is the predicted cumulative log return from the linear fit, \bar{y} is the mean of the observed cumulative log returns and n is the number of observations.

- Omega ratio: A probability-weighted ratio that compares the gains versus losses relative to a threshold return target, either set to zero or the risk-free rate. It is a proxy for the likelihood of achieving returns above a certain threshold compared to the likelihood of experiencing losses below that threshold.

$$\text{Omega Ratio} = \Omega(\theta) = \frac{\int_{\theta}^{\infty} (1 - F(r)) dr}{\int_{-\infty}^{\theta} F(r) dr} \quad (\text{E.3})$$

where θ is the threshold return and $F(r)$ is the cumulative distribution function of returns.

- Sortino ratio: A modification of the Sharpe ratio that only considers downside risk, measuring the risk-adjusted return of an investment by comparing the excess return to the downside deviation.

$$\text{Sortino Ratio} = \frac{R_p - R_f}{\sigma_d} \quad (\text{E.4})$$

where R_p is the annualised return of the portfolio, R_f is the annualised risk-free rate, and σ_d is the downside deviation of the portfolio returns, i.e. the standard deviation of only the negative returns.

- Skewness: A measure of the asymmetry of the return distribution, indicating whether returns have a tendency toward positive or negative extremes.

$$\text{Skewness} = \frac{\sum_{i=1}^n (R_i - \mu)^3}{(n - 1) \cdot \sigma^3} \quad (\text{E.5})$$

where R_i is the return of the portfolio at time i , μ is the mean return, σ is the standard deviation of returns, and n is the number of observations.

- Kurtosis: A measure of the "tailedness" of the return distribution, indicating the

frequency of extreme values compared to a normal distribution.

$$\text{Kurtosis} = \frac{\sum_{i=1}^n (R_i - \mu)^4}{(n-1) \cdot \sigma^4} \quad (\text{E.6})$$

where R_i is the return of the portfolio at time i , μ is the mean return, σ is the standard deviation of returns, and n is the number of observations.

- Tail ratio: A measure of the relative magnitude of extreme positive returns compared to extreme negative returns, calculated as the ratio between the right (95th percentile) and the left (5th percentile) tails of the distribution of returns.

$$\text{Tail Ratio} = \frac{\text{Right Tail (95th Percentile)}}{\text{Left Tail (5th Percentile)}} \quad (\text{E.7})$$

- Daily Value at Risk (VaR): A measure of the potential loss that could occur in a portfolio over a single day at a specified confidence level, typically 95%. The computation is done using the variance-covariance calculation:

$$\text{Daily VaR} = (R_p - z \cdot \sigma_p) \times V_p \quad (\text{E.8})$$

where R_p is the expected daily return, z is the z-score of the desired level of confidence, σ_p is the daily standard deviation of portfolio returns, and V_p is the value of the portfolio.