

Reinforcement learning for sequential decision-making: a data driven approach for finance



Candidate : Alessio Brini

Advisor: Daniele Tantari

Ph.D. Thesis in Financial Mathematics A.Y. 2021-2022

Abstract

This work presents a variety of reinforcement learning applications to the domain of finance. It composes of two-part. The first one represents a technical overview of the basic concepts in machine learning, which are required to understand and work with the reinforcement learning paradigm and are shared among the domains of applications. Chapter 1 outlines the fundamental principle of machine learning reasoning before introducing the neural network model as a central component of every algorithm presented in this work. Chapter 2 introduces the idea of reinforcement learning from its roots, focusing on the mathematical formalism generally employed in every application. We focus on integrating the reinforcement learning framework with the neural network, and we explain their critical role in the field's development.

After the technical part, we present our original contribution, articulated in three different essays. The narrative line follows the idea of introducing the use of varying reinforcement learning algorithms through a trading application (Brini and Tantari, 2021) in Chapter 3. Then in Chapter 4 we focus on one of the presented reinforcement learning algorithms and aim at improving its performance and scalability in solving the trading problem by leveraging prior knowledge of the setting. In Chapter 5 of the second part, we use the same reinforcement learning algorithm to solve the problem of exchanging liquidity in a system of banks that can borrow and lend money, highlighting the flexibility and the effectiveness of the reinforcement learning paradigm in the broad financial domain. We conclude with some remarks and ideas for further research in reinforcement learning applied to finance.

Acknowledgements

I want to spend some words to thank my supervisor, Daniele Tantari, for his continuous guidance during the last years. My journey as a researcher would have never been possible without meeting him. I am also grateful to professor Maria Elvira Mancino for her constant availability and support in my professional development and to all the professors I had the opportunity to meet at the Scuola, among which Giulia Livieri and my Ph.D. coordinator Stefano Marmi. A deserved mention also for professor Gabriele Tedeschi, an amazing coauthor and friend, and professor Petter Kolm, an experienced guidance in the world of finance.

A special thank also goes to Giacomo Toscano, who has been my mentor in every part of this journey.

On the personal side, I thank infinitely my family, Deborah, Enrico and Lorenzo, for their continuous support in every moment of my life.

To conclude, many things would have been different if I hadn't met my once girlfriend and now wife, Barbara, who made me believe in myself and convinced me to pursue this experience. I will be thankful to her for life.

Contents

I	Machine learning paradigms: from the basics to reinforcement learning	18
1	Machine learning: basics concepts and neural networks	19
1.1	Learning paradigms	21
1.2	Generalization, fitting and estimation	23
1.3	Background Theory	28
1.4	Early models for artificial neural network	29
1.5	Feedforward neural networks	31
1.6	Design choices	34
1.7	Architecture of a neural network	37
1.8	Universal approximation theorem	38
1.9	Optimizing the neural network parameters	40
1.10	Gradient issues during training	43
1.11	Summary	48
2	Reinforcement learning: theory and background	50
2.1	Basic notions	51
2.2	Markov Decision Processes	55
2.3	Tabular Reinforcement Learning	59
2.3.1	Dynamic Programming	59
2.3.2	Model-free Reinforcement Learning	62
2.4	Approximate Reinforcement Learning	68
2.4.1	Value-based methods	73
2.4.2	Policy-based methods	78
2.5	Summary	87

II	Reinforcement learning: simulation analysis with applications to finance	88
3	Deep Reinforcement Trading with Predictable Returns	90
3.1	Introduction	90
3.2	Financial Market Environment	92
3.3	Deep Reinforcement Learning Methods	94
3.3.1	Tabular Reinforcement Learning	95
3.3.2	Approximate Reinforcement Learning	96
3.4	Numerical Experiments	98
3.4.1	Tracking the Benchmark	100
3.4.2	Outperforming the Benchmark	104
3.5	Conclusions	112
3.6	Appendix : Algorithms and Hyperparameters	112
3.6.1	Q-learning	113
3.6.2	DQN	114
3.7	PPO	116
3.7.1	Environment choices	118
4	Reinforcement learning for trading in a market with frictions: a residual approach	120
4.1	Introduction	120
4.2	Multiperiod problem	123
4.3	Financial model for the excess return	126
4.4	The residual RL approach: empirical results	128
4.4.1	Residual RL approach	128
4.4.2	Simulation results	130
4.5	Further directions	135
4.5.1	Residual approaches	136
4.5.2	Regularization approaches	136
4.6	Conclusions	137
5	Reinforcement Learning Policy Recommendation for Inter-bank Network Stability	139
5.1	Introduction	139
5.2	Model	144
5.2.1	The interbank market microstructure	144

CONTENTS

5.2.2	Banks microfoundations: the dynamics of lending agreements and trading strategies	146
5.2.3	The optimization algorithm: Proximal Policy Optimization	151
5.3	Simulation Results	153
5.3.1	Training the PPO algorithm	154
5.3.2	Micro and macro consequences of the policy recommendation	158
5.3.3	The reinforcement learning based recommendation for taming systemic risk	167
5.4	Concluding remarks	171
5.5	Appendix: A sensitivity analysis on model parameters	172
5.6	Appendix: Algorithms and Hyperparameters	177
6	Conclusions: main results	180

List of publications

Articles under review

- Brini, A., & Tantari, D. (2021). Deep Reinforcement Trading with Predictable Returns. arXiv preprint arXiv:2104.14683, *Submitted, Annals of Operation Research*.
- Brini, A., Tedeschi, G., & Tantari, D. (2022). Reinforcement Learning Policy Recommendation for Interbank Network Stability. arXiv preprint arXiv:2204.07134 *Submitted, Journal of Financial Stability*.

Articles in progress

- Brini, A., & Kolm, P., Reinforcement learning for trading in a market with frictions: a residual approach. *In progress*.

Introduction

Since the advent of the Black-Scholes (Black and Scholes, 1973) formula for pricing complex derivatives, mathematical techniques have helped to solve financial problems, aiming to capture different aspects of this world. *Quant* is a word that informally indicates a quantitative financial analyst that uses state-of-the-art mathematical and statistical tools to generate predictive signals and create trading strategies, i.e., make money in the financial market. From the end of the 90s to the last decade, the realized technological advancement sped up the research in the mathematical and computational domains, giving rise to new effective financial analysis methods. Machine learning represents one of the research domains that has been affected by the surges in the availability of computational power and therefore attracted the interest of the financial sector. This field of research takes its roots from the computer science domain but embraces concepts from mathematics and statistics that make it truly interdisciplinary.

The last two decades have marked a shift in the various research strands in finance because of the introduction of advanced mathematical techniques to improve the modeling capabilities and beat the market. Nowadays, machine learning tools have become more ubiquitous in finance research to a significant extent than when the various projects underlying this thesis started. As the authors, we had the priceless privilege to tackle new research questions in an almost new set of applications and use novel techniques for the field of finance, each with its advantages and disadvantages. We also assisted in a massive shift towards data-driven modeling techniques by academics and practitioners in the industry. The Organization for Economic Co-operation and Development (OECD, 2021) has recently recognized the importance of this emerging trend for the financial sector to help the policymakers in supporting the introduction of this new strand of innovation in the industry. Goodell et al., 2021 identifies different directions of research for machine

CONTENTS

learning technologies applied to finance. Big international banks such as JP Morgan have already set up their research lab ¹, and they are working on solutions for the business with the use of machine learning and artificial intelligence. At Fidelity Labs, one of the major asset management companies in the United States, they are trying to put in production a reinforcement learning approach for the optimal solution of the best execution and portfolio management problems (Halperin et al., 2022).

Before delving into specific use cases of machine learning for our domain of interest, it is essential to consider its origin as a field of study and its underlying idea. Machine learning is the development of mathematical tools and algorithms that exploit data to learn a task, attempting to imitate human learning behavior. These algorithms repeat iterative procedures that allow for improvement in performing a specific task gradually. For these reasons, it is common to say that a machine learning algorithm is *trained* to accomplish an objective, which the researcher sets to indicate what the computer program should learn. Even though machine learning gained massive popularity in recent years, its long history started before the 50s when McCulloch and Pitts, 1943 presented a mathematical model to imitate the functioning of the human brain, called artificial neuron. The model was refined later by the perceptron of Rosenblatt, 1958, which is one of the first examples of a mathematical model able to classify a set of points into binary classes. The artificial neuron represents the ancestor of modern neural networks. Despite the excitement caused by the novelty of their research at that time, Minsky and Papert, 1969 severely criticized the perceptron model for its limits because it was able to learn only linearly separable functions of the data and struggled with more complex functions. Notwithstanding the encouraging results of creating neuron-like models, neural networks remained overlooked due to their lack of capabilities and background theory. A period called *Artificial intelligence winter* followed, referring to about ten years in which the initial promise of artificial intelligence, which is a broader field of research including machine learning, bumped into the lack of sophisticated computer programs to test their effectiveness. The overinflated hype, which has commonly occurred for many other emerging technologies such as the World Wide Web, exacerbates the gap between the promising academic research and the lack of instruments to put it into practice. In the mid-80s, Rumelhart et al., 1986 shed light on the backpropagation algorithm to

¹JPM AI research

CONTENTS

perform automatic differentiation of complex functions. Based on previous works (Linnainmaa, 1970; Werbos, 1982), they were the first to demonstrate that backpropagating the gradients can yield an expressive internal representation of data in the hidden layers of neural networks. Their work helped the subsequent popularity of these mathematical models because it efficiently computes the enormous amount of partial derivatives required to train a neural network. The backpropagation algorithm, together with the universal approximation theorem of Hornik et al., 1989, proved the capability of a multi-layer neural network to learn any continuous function theoretically. Those groundbreaking researches represented a turning point in the history of such mathematical models. LeCun et al., 1989 provided one of the first examples of such capabilities by teaching neural networks to recognize handwritten zip codes of the US Postal Service. In this simple pattern recognition problem, neural networks reached a human-level accuracy.

Not surprisingly, the central moment that marked the way for neural networks to become widely employed in commercial applications was the release of the ImageNet database (Deng et al., 2009). In many cases, increasing data availability means more possibility to exploit machine learning tools and obtain groundbreaking results. ImageNet consists of a vast database of different classes of images at the disposal of researchers to test and develop novel machine learning algorithms. ImageNet was born with the idea of providing an adequately organized dataset for machine learning research and favored a renovated trust in neural networks' capabilities. Krizhevsky et al., 2012 deployed a multi-layer neural network that was able to outperform in classification accuracy all the other competing models in the ImageNet Large Scale Visual Recognition Challenge (ILSVRC). This success defined machine learning as a groundbreaking tool for solving similar computer vision and pattern recognition problems. From the decade following the work of Krizhevsky et al., 2012 up to the recent days, neural networks have incrementally improved their performances thanks to the effort of a growing community of machine learning researchers. Several alternative formulations have taken place in order to solve many of problems in different domains, from the pattern and face recognition (Le, 2013; Taigman et al., 2014) to language translation (Devlin et al., 2018; Vaswani et al., 2017) and data generation (Radford et al., 2015). These applications are generally identified under the umbrella of deep learning applications since they employ complex neural networks with many layers, hence the adjective *deep*. Neural networks allowed for various and challenging tasks thanks to their capability to pro-

CONTENTS

cess a considerable amount of data and discover patterns or more meaningful representation of the available information.

These characteristics of the neural networks have also been central to the development of a specific paradigm of machine learning, called deep reinforcement learning, which merges the mathematical advancement made through time in the reinforcement learning field with the use of these biologically inspired models. Reinforcement learning is a mathematical framework that models the behavior of an agent that needs to learn how to move in a dynamical system. Within the broad strand of research in artificial intelligence and machine learning of the second part of the last century, Sutton and Barto, 2018 identifies the roots of modern reinforcement learning as the collision of different threads of research that shares common aspects: learning by trial and error and optimal control. The former pursues the idea of animal learning and relies on the concept of events that reinforces the selection of specific actions by the agent. The latter optimizes a measure of the performance of an agent who acts in a system whose dynamics is known. Optimal control has not consistently been recognized as a form of reinforcement learning since the class of mathematical methods used to solve these problems, referred to as dynamic programming (Bellman, 1957), requires full knowledge of the dynamic system. For this reason, optimal control problems resemble more an optimization problem than a learning one. The two strands of research followed different paths for an extended period, possibly due to the nature of their approach, philosophy for learning, and mathematics for optimal control. Werbos, 1987 proposed the first interrelation between learning and dynamic programming approach that culminated with the work of Watkins and Dayan, 1992 which treats the reinforcement learning problem under the mathematical formalism of the optimal control theory. Since then, modern reinforcement learning has been a standalone field of research that embraces both perspectives. The turning point for the research domain has been the introduction of neural networks as instrumental for the learning process. Algorithms and concepts of reinforcement learning (Lin, 1992; Tesauro et al., 1995; Watkins and Dayan, 1992) have been revisited to leverage the data representation capabilities of neural networks as in the work of Mnih et al., 2013, which introduces the Deep Q-Network, one of the first learning algorithm able to play video games at a superhuman level. This novel introduction to the field of research gives rise to the so-called deep reinforcement learning, which nowadays includes a large community of computer scientists and mathematicians. AlphaGo (Silver et al., 2017b) represents a game-changer application

CONTENTS

to highlight the disruptive impact of neural networks for reinforcement learning. It is a computer program devised to play the game of Go and beat the world champion of the homonymous board game by just knowing its rules. Indeed, reinforcement learning is a machine learning paradigm that allows solving a sequential decision-making problem without specific knowledge or information. Even though it is just a board game, Go represents an arduous task to solve for a computer program since it is probably the most complex game ever devised by humans. It requires the player to have deep thinking and reasoning capabilities. The software programmed by DeepMind, nowadays, a subsidiary of Alphabet Inc., was empowered with reinforcement learning techniques and has also been improved since that time (Schrittwieser et al., 2020; Silver et al., 2017a). Games are an appropriate testbed for reinforcement learning because the inherent concept of scores makes it easy to measure incremental signs of progress. This aspect is easily generalized to the financial world, where the concept of profit is ubiquitous, providing an impressive number of opportunities for applying the reinforcement learning paradigm.

Indeed, the astonishing achievements obtained by deep neural networks for applied research in various domains are mostly the results of two important factors:

- The large amount of data produced by modern society;
- The disposal of powerful computational machines equipped with high-performance graphic processing units (GPUs) to efficiently run and test modern machine learning models (Raina et al., 2009).

Combining these two factors is pivotal in the ongoing development of the finance industry and research, opening the possibility to explore novel ideas and test new algorithms. Machine learning creates the opportunity to analyze and exploit the massive data production in the financial world. Nowadays, financial firms, either on the buy or the sell side, are increasingly adopting technology for their daily routine business, and machine learning techniques are crucial in this shift of paradigm (Ryll et al., 2020). Even though neural networks are not the only machine learning models employed for solving financial problems (see (Malliaris and Malliaris, 2015; Wang et al., 2009) for decision trees, (Ananthakumar and Sarkar, 2017; Hua et al., 2007) for logistic regression and (Kim, 2003; Kumar and Thenmozhi, 2006) for support vector machines), they truly embody a novelty rather than standard methods

CONTENTS

of regression analysis commonly used in finance. Neural networks can approximate very complex functions of the data that contain different degrees of nonlinearity. The financial markets and the whole economy in a broader sense are known to be described by highly nonlinear relationships (Amini et al., 2021; Brock, 2018), making neural networks appealing for modeling economic phenomena. Machine learning techniques for finance represent a further development beyond classical econometric analysis, which set the basis for studying finance as a quantitative subject. López de Prado, 2019 highlights the connections and the differences between the new and the classical approach from the perspective of the current availability of more complex and unstructured datasets, which can hardly be analyzed with econometric tools. Despite their promises, the skepticism in the broad adoption of these techniques for financial analysis is due to the underlying complexity that makes it hard to interpret them. For this reason, practitioners often prefer to adopt simpler models which may be less effective for the specific task but whose results are easier to interpret. Therefore, the complexity of machine learning motivates a critical approach before applying these techniques instead of a more simple method of analysis because the benefits of improved data-driven modeling capabilities have to be larger than the loss in explainability. To this end, the field of explainable AI (XAI) emerged during the last decade to guarantee a certain degree of control and awareness of what interactions and relationships complex machine learning is retrieving from the provided data. In this way, complex models as neural networks become interpretable, benefiting the application to sensitive domains like finance.

The use cases for machine learning powered by neural network in finance are multiples and touch different aspects of the subject: forecasting prices or economic indicators (Kaniel et al., 2021; Lim et al., 2019; Ravi et al., 2017; Sehgal and Pandey, 2015), financial trading and portfolio management (Benhamou et al., 2020; Cong et al., 2021; Deng et al., 2016; Guida, 2020; Heaton et al., 2017; Jiang et al., 2017; Noguer i Alonso and Srivastava, 2020), hedging (Du et al., 2020; Kolm and Ritter, 2019b) and credit scoring evaluation (Wang et al., 2018; Zhu et al., 2018).

Among all the machine learning approaches, the reinforcement learning paradigm assumes a central role in adopting new technologies for solving problems in the field of finance. It provides the opportunity to model a sequential decision-making problem with fewer assumptions than standard optimization methods and guarantees flexibility in shaping the structure of the problem. The important novelty introduced by reinforcement learning to

finance concerns the possibility of letting a market participant interact with the market itself, which is a separate entity whose characteristics and dynamics might be partially or fully unknown. Many challenges faced by financial agents are control problems, where not all the variables are observed. Reinforcement learning indeed represents a way to tackle this set of tasks without relying on complex models that try to capture the financial dynamics with restrictive or unrealistic assumptions. To this end, reinforcement learning focuses more on analyzing the feedback produced by the agent-environment interaction. This paradigm shifts from a model-driven to a data-driven approach to solving a sequential control problem. Due to the availability of code repositories maintained by trusted software companies, such as DeepMind² and OpenAI³, researchers can access ready-to-use, and bug-free implemented reinforcement learning algorithms to use for solving control problems in finance. Due to their sequential nature and such a good adaptation to the reinforcement learning paradigm, implementation of code libraries specific for this domain have already emerged (Liu et al., 2020). Hambly et al., 2021 provides a comprehensive review of financial problems tractable under the light of the reinforcement learning paradigms. Some examples are optimal execution (Dabérius et al., 2019; Hendricks and Wilcox, 2014), portfolio optimization (Du et al., 2016; Park et al., 2020; Xiong et al., 2018), market-making (Abernethy and Kale, 2013; Spooner et al., 2018), roboadvising (Alsabah et al., 2021; Wang and Yu, 2021), option pricing (Cao et al., 2021; Halperin, 2020) and order routing (Agarwal et al., 2010; Ganchev et al., 2010). However, most of the research contribution relies on a case-by-case constructed dataset, posing the issue of reproducibility and generalization of the obtained results on a different dataset and problem formulation.

Contribution

The main objective of this thesis is to offer a thoughtful application of reinforcement learning to a set of financial problems, such as learning how to trade financial assets with simulations in a controlled environment to highlight the role of the paradigm in reaching the solution. It is not an application of ready-to-use implemented algorithms on gathered financial data. We provide instead results obtained when we know the nature of the data generating

²DeepMind

³OpenAI

CONTENTS

process, and we can modify it to test the robustness of the approach. In this way, reinforcement learning becomes a valid alternative to the optimal control framework without being used as a black box. Understanding how these algorithms work allows us to intervene and adapt them to the economic context, which is highly different from other domains where machine learning has already been applied and brings challenges. Our attempt is also to leverage prior knowledge of the financial market to avoid an entirely data-driven application of reinforcement learning in a noisy financial environment. The last contribution is a novel application of reinforcement learning to an interbank lending problem where we consider an agent to be a public decision-maker that needs to direct the system of banks toward an exchange of liquidity that should be profitable and guarantees systemic stability.

This thesis focuses on reinforcement learning as a framework to represent and mimic many interactions in the financial sector among the different market participants. Almost any problem in finance has a temporal structure with delayed feedback for the actions that the market participants take. For instance, consider a trader that needs to enter a position in the market and then exit gradually from it as the signal decays over time. These temporal choices are ubiquitous in finance and characterize many tasks that market participants need to carry on.

We study reinforcement learning agents of two types among the different sequential financial tasks. On the one hand, we focus on trading and portfolio optimization at a microscopic level and investigate the effectiveness of a single reinforcement learning agent that trades in a financial market. On the other hand, at a macroscopic level, we study the influence of a regulatory agent, which uses a reinforcement learning approach to guide an entire financial system of banks. We believe that evaluating the reinforcement learning approach on different tasks at multiple scales is crucial to show the importance that it can have for finance as a field of research. We focus on those set of problems that naturally reminds the importance of timing the choices. We also firmly believe that this class of problems represents a good testbed for reinforcement learning application. Intertemporal optimization and control methods already have a substantial presence in the financial literature that solves problems. Hence, they serve as a proper comparison with the reinforcement learning approach. Differently from the majority of the approaches in this regard, our first contribution studies the capabilities of the reinforcement learning framework to solve a financial trading problem in a known environment, where an exact solution already exists (Brini and Tan-

CONTENTS

tari, 2021). Simulating synthetic financial returns that move according to a specific factor-model (Gârleanu and Pedersen, 2013) and trading them in a financial market that admits an optimal solution is relevant for the scope of our research. We want to decouple the performances of the tested algorithms from common sources of noise in real financial applications, such as data quality and preprocessing. The main difference from the existing literature is using a controlled environment that allows comparing different reinforcement learning algorithms to a dynamic programming approach, which is standard for this optimal control problem. By performing a considerable number of synthetic simulations, we test reinforcement learning algorithms of three different classes by selecting the most representative for each to enlighten their strengths and weaknesses concerning the exact benchmark solution and each other. We aim to avoid the issues of modern reinforcement learning algorithms by disentangling the effect of the goodness of the signal used for trading and portfolio application from the actual increase in performance brought by the capability of the reinforcement learning algorithms. We investigate the performance of deep reinforcement learning traders in a market environment with different known mean-reverting factors driving the dynamics. When the framework admits an exact dynamic programming solution, we can assess the limits and capabilities of different value-based algorithms to retrieve meaningful trading signals in a data-driven manner. This approach guarantees flexibility and outperforms the benchmark strategy when the price dynamics are misspecified, i.e., some original assumptions on the market environment do not hold because of extreme returns and volatility clustering effects. In the latter case, we discover that the policy-based reinforcement learning algorithm works better than the value-based one since it is more robust to extreme events and heteroskedastic volatility. This work also introduces and tests a simple tabular reinforcement learning method that helps to motivate the need for deep reinforcement learning algorithms for a financial trading problem.

The insights obtained by such simulated analysis are further elaborated in a currently work in progress project to leverage prior knowledge of the trading problem and use it to improve the reinforcement learning performances. Model-free reinforcement learning is indeed fallacious in modeling complex trading signals, as in a multi-asset trading approach. For these reasons, we modify the training setup of a reinforcement learning algorithm to consider known solutions to the trading problem. We leverage the neural networks' approximation capability to learn how to trade by avoiding high transaction

costs, which are highly nonlinear in the trades. We refer to this approach as residual reinforcement learning, which is scalable to a large enough cross-section of assets and applies to real financial data.

The third contribution of this work focuses on reinforcement learning in an agent-based model setting. (Brini et al., 2022) represent a still new application of reinforcement learning, where the agent learns a strategy to improve the flow of liquidity through interbank lending agreements. We modeled a public policy recommendation through a reinforcement learning approach. We then analyze the effect of such policy on an artificial interbank market, where financial institutions can stipulate lending agreements. In this context, the reinforcement learning paradigm maximizes the system’s long-term fitness by gathering information on the economic environment and directing economic actors to create credit relationships based on the optimal choice between a low-interest rate or high liquidity supply. Financial institutions create or cut their credit connections through time via a preferential attachment evolving procedure that generates a dynamic network based on the combination of the public signal and private information.

The reinforcement learning optimal policy recommendation plays a crucial role in mitigating systemic risk with respect to alternative policy instruments. It provides a certain amount of flexibility to solve the problem efficiently when the underlying environment is as complex as an agent-based model and the underlying dynamics are not fully specified. To interpret the choices carried out by the learned public policy, we also employ a state-of-the-art tool for enhancing the explainability of the retrieved solution and devising insights on what has been learned by the machine learning tool regarding the systems of banks. Moreover, our results show that the emergence of a core-periphery interbank network, combined with a certain level of homogeneity in the size of lenders and borrowers, is an essential feature to ensure the system’s resilience.

Thesis structure

Based on the outlined contributions, this work composes of two parts. The first part is a technical overview of the basic concepts in machine learning, which are required to understand and work with the reinforcement learning paradigm. These concepts are not specific to finance and are shared among the application domains of machine learning. Chapter 1 outlines the funda-

CONTENTS

mental principle of machine learning before introducing the neural network model as a central component of every algorithm presented in this work. Then some theory regarding neural networks introduces how they work and are trained together with the different use cases. Chapter 2 introduces the idea of reinforcement learning from its roots, focusing on the mathematical formalism generally employed in every application. Once we have outlined the framework, we extend this part by introducing neural networks in a reinforcement learning context and explaining their critical role in the field's development. For this purpose, we enter into the details of the reinforcement learning algorithms employed in the practical applications of our work mentioned in the contribution section. Therefore, we provide a detailed explanation of the reasoning and motivation of the two broad families of deep reinforcement learning algorithms, i.e., value-based and policy-based.

Once we have reviewed the mathematical framework and all the concepts underlying machine learning and reinforcement learning theory, the second part of the work presents our work's original contribution, which is articulated in three different essays. The narrative line follows the idea of introducing the use of varying reinforcement learning algorithms through a trading application (Brini and Tantari, 2021) in Chapter 3. Then in Chapter 4 we focus on one of the presented reinforcement learning algorithms and aim at improving its performances and scalability in solving the trading problem by leveraging prior knowledge of the setting. In Chapter 5 of the second part, we use the same reinforcement learning algorithm to solve the problem of exchanging liquidity in a system of banks that can borrow and lend money, highlighting the flexibility and the effectiveness of the reinforcement learning paradigm in the broad financial domain. We conclude with some remarks and ideas for further research in reinforcement learning applied to finance.

Part I

Machine learning paradigms: from the basics to reinforcement learning

Chapter 1

Machine learning: basics concepts and neural networks

Machine learning is a broad field of research that embraces tools from different subjects such as mathematics, statistics, and computer science and focuses on developing efficient numerical approaches to solve complex problems in various domains. It is interdisciplinary since elements of linear algebra, probability theory, statistics, and computer programming are required to be fully understood and tackled. Machine learning is divided into various fields of study, each with specific scopes, such as prediction, object detection, or data generation. At the same time, it is also a subfield of artificial intelligence, commonly defined as the capability of a computer program to imitate intelligent human behavior. Although the difference between machine learning and artificial intelligence can be subtle and sometimes confusing, the latter is a more general concept that aims at implementing a human-like general intelligence by using sophisticated mathematical techniques, among which those included in machine learning.

It is important to stress that learning refers to the known ability of humans and animals to acquire the capability to carry out a task and to be able to generalize when the context changes. The peculiar aspect differentiating machine learning from other existent numerical techniques is the opportunity to program an algorithm that learns a model from the data and hence limits the number of required assumptions. The first step to understanding the contribution of the machine learning approach for financial applications, which will be the focus of our work, is to compare it to econometrics, which is the established approach in financial academic and industry research. Both

approaches want to build predictive models using explanatory variables, usually referred to as features in the literature. Since the development of those research fields occurred in parallel, they often accomplish the same tasks with different methodologies derived from the different philosophical approaches to the problems.

On the one hand, econometrics build models based on assumptions to describe the dynamics of the underlying variables and make inferences from the available data. On the other hand, machine learning use algorithms that require fewer assumptions and are data-driven without necessarily injecting prior knowledge to the model (Charpentier et al., 2018). It is important to remark that learning should be considered not as the task that machine learning aims to achieve but instead as the means of reaching the proposed task. Notwithstanding the difference in approaching the modeling problem, machine learning and econometrics share some basic concepts that we outline in this chapter. Using machine learning instead of econometrics generally improves the modeling capability, allowing to deal with complex datasets and discovering hidden patterns that would be hard to model and incorporate into the latter approach. Even though econometrics can capture complex nonlinear interactions with the variables, this modeling effort requires assumptions based on prior knowledge of the problem. On the contrary, machine learning can discover relationships in the data that have not been observed yet, relaxing the modeling effort and possibly providing more insightful analysis driven by the data.

Whenever we need to define a machine learning algorithm, we should first ask ourselves three questions:

1. What do we want to learn?
2. How do we learn it?
3. How do we measure the progress of learning?

A learning algorithm is a set of iterative procedures that has an answer to all these questions. Usually, an algorithm aims at learning a general task T through some data D according to a performance measure M . These are the core elements that provide the structure of a learning algorithm. Excluding one of those elements conflicts with the purpose of such algorithmic procedures. For instance, an unclear definition of the task to perform does not allow defining the aim of the learning algorithm. The absence, even partial,

of the dataset does not allow having enough examples to learn. A wrong definition of the performance measure can lead to possible misjudgments of the learning procedure. For this reason, in the next section, we define each of these core elements properly and connect them to present the machine learning modeling setup.

Therefore, in this chapter, we want to outline the structure of a learning algorithm and its main components. Then we introduce neural networks, which are the machine learning tool employed in the second part of this work, and provide a detailed overview of the basics, the theory, and the optimization process underlying this family of models.

1.1 Learning paradigms

Different machine learning paradigms are defined depending on the type of dataset that an algorithm is allowed to experience. In what follows, we assume a dataset \mathbb{Z} to be composed by a number M of vectors $\mathbf{z} \in \mathbb{R}^K$ where z_i represents the i -th feature. These vectors are generally referred to as examples or data points. We assume the dataset to be a two-dimensional matrix with M rows and K columns. We present the different machine learning paradigms in a general bidimensional case, even though machine learning algorithms can also work with higher-dimensional datasets, as happens for image recognition tasks.

Supervised learning is the prevailing machine learning paradigm at the core of many commercial applications. These kinds of algorithms experience a dataset composed of some data points \mathbf{x} , which includes the values z_i up to $K - 1$ and the labels or targets y , which are the values z_K . These algorithms aim at learning a function of the data that returns the label given the corresponding example, i.e. $f(\mathbf{x}) = y$ where $\mathbf{x} \in \mathbb{R}^{K-1}$. A single label can be either any real value $y \in \mathbb{R}$ or a value from a discrete set $y \in \{1 \dots k\}$. The former output type refers to a regression problem, where the learned function returns a real-valued output. In contrast, the latter refers to a classification problem, where a data point maps to some category identified by numbers in the finite discrete set. The term supervised refers to the researchers manually labeling the dataset to specify its context for the models, e.g., stating which objects are included in an image so that they oversee and direct the learning process. The choice of the label is indeed the most important part in the setting of a supervised problem because it affects the type of mapping that the

algorithm learns from the data, i.e., the objective of the learning itself. Some examples of supervised learning tasks are object recognition (Hu et al., 2015; Mehdipour Ghazi and Kemal Ekenel, 2016; Wang, 2016), recommendation systems for targeted commercials (Zhou, 2020), and prediction of housing prices given the area’s characteristics (Truong et al., 2020).

In contrast to that, *unsupervised learning* algorithms want to learn valuable properties of the dataset, which generally resort to learning the probability distribution $p(\mathbf{z})$ that produces the data points. The algorithm does not receive any label because it should infer information from the data itself instead of learning a mapping to some output. Unsupervised learning can help find a customer segment, reduce the problem’s dimensionality, compress the data into a new representation (Zhang and Saniie, 2021), or make feature selections before applying supervised algorithms (Sharang and Rao, 2015; Taherkhani et al., 2018).

The main difference is that supervised learning needs to find a mapping from \mathbf{x} to y , while unsupervised learning searches for patterns in the data to provide a different representation and learn some data characteristics.

It is important to remark that there is a subtle line between the two paradigms, and there are machine learning models that perform both tasks. However, a distinct classification helps classify different applications of the same model and maintain an organic view of the field of machine learning. Notwithstanding, the major shared characteristic of these two paradigms is using a dataset that is provided in advance and does not change during the learning process.

Reinforcement learning, the third machine learning paradigm, is different from the previous two in its purpose and the way it exploits the dataset. Reinforcement learning does not experience a fixed dataset \mathbb{Z} , and it is not focused on learning a specific relationship between some features and their label or understanding some properties of the features. Instead, it represents a framework in which an algorithm, referred to as the agent, interacts with external factors, referred to and summarized under the environment, to learn how to carry out a control task by trial and error. For this reason, every time the agent interacts with the environment by taking action, the dataset is possibly augmented if not entirely replaced by a new set of observations. In addition, the dataset also presents a specific structure in which every piece of it plays a particular role so that each vector \mathbf{z} must contain information about the state of the environment, i.e., the variables that characterize it, the action carried out by the agent and the reward, i.e., the response of

the environment to the agent’s action. Therefore, reinforcement learning is a machine learning paradigm to perform optimal control tasks, such as teaching a sensorimotor robot how to move in a restricted space (Ghadirzadeh et al., 2016), optimally managing scarce resources (Mao et al., 2016), or managing a portfolio of equity options (Buehler et al., 2019). We postpone the details about this machine learning paradigm and the use of such a peculiar dataset to the Chapter 2, as it represents the core of this work, and we will describe it there accurately.

A common aspect shared by all those learning paradigms is the training process of the selected machine learning algorithm to perform a predefined task. Training a machine learning model consists in providing the algorithm with the data to learn from and repeatedly updating the set of parameters θ of the algorithm based on the information contained in such data. Hence, the training procedure differs depending on the machine learning paradigm, i.e., the kind of dataset experienced. Model training represents the first step of all machine learning approaches, which results in a trained model that can be tested and evaluated through the selected performance metric and eventually deployed to be used in production. The term “training” is used in the machine learning literature as synonymous with fitting, which is more used in other branches of statistics, but it means essentially the same research of the good set of model parameters to properly represent the data.

1.2 Generalization, fitting and estimation

The aim of a machine learning algorithm goes further than finding the model that best fits the data. Generalizing over unseen examples is a critical concept to distinguish between fitting and learning. We will describe these concepts under the supervised learning framework to ease the exposition, although they are also common to the other machine learning paradigms. Recalling the mapping $f(\mathbf{x}) = y$ that a supervised learning algorithm wants to retrieve from the dataset, the function of the input \mathbf{x} should approximate well the corresponding target also when a specific data point is not used for training.

The capability to teach a machine learning algorithm how to generalize depends on organizing the data before training the model. Generally, the dataset is divided in two parts: the *training* set and the *test* set. The evaluation of the training progress goes through a selected performance metric $J(\theta)$, which depends on the model parameters θ and returns a measure of the

difference between the predicted output and the true value. While this error computed on the training set is important to evaluate the progress of the learning, the discrepancy over the test set is ultimately more important in machine learning because it represents the generalization error. The model does not use the data points included in the test set, and therefore these points are not previously observed by the machine learning algorithm. A good performance on this set of examples, i.e., a low test error, signals a good generalization.

Statistical learning theory assumes the data points in the training and the test set to be generated by a so-called data generating process (DGP). Therefore, each data point is considered to be independent and identically distributed (i.i.d) since they are assumed to be drawn separately from the same underlying distribution $p(\mathbf{x}, y)$. Assuming to draw both sets from the same DGP causes the expected training error to be equal to the expected test error for a given model. The only difference so far is the role we assign to each set. However, this is valid only if we fix the model parameters and then evaluate the algorithm over the training and the test set. On the contrary, a machine learning algorithm is optimized over the training set so that the parameters repeatedly change through iterative updates. Then the test set is evaluated, and the expected test error ends up being greater or equal to the expected training error. This difference clarifies the two concurrent aims of a machine learning algorithm: to reduce the training error while minimizing the test and training error gap. The gap between the training and test error determines the different outcomes of the training procedure of a machine learning algorithm. Obtaining a high training error results in *underfitting* since the model is not able to capture the information in the data. The opposite situation happens when the gap between the test and the training error is large, which means that the model fits perfectly the training data, but it is not able to generalize to unseen examples. This situation is called *overfitting*.

One way to control the overfitting-underfitting tradeoff is to change the model capacity, i.e., the possibility to choose the function to represent the model and to fit the data from a broad set of functions. The model capacity can be controlled through the magnitude of the hypothesis space \mathcal{H} , which is the set of possible functions that the algorithm can select as a solution to the machine learning problem (Abu-Mostafa et al., 2012). For example, in a linear model, the set \mathcal{H} includes all the possible linear equations that the algorithm can choose to fit the data. Therefore, a larger number of parameters

in a linear model composes a bigger hypothesis space for the machine learning model under consideration, which in turn can find a more expressive way to represent the function of the data. In the next section, we provide a brief overview of the underlying background theory in machine learning, including the choice of the algorithm capacity. Figure 1.1 visualize the tradeoff between the model complexity and the predictive capability that a machine learning algorithm can achieve. Usually, a range of model complexity represents a good compromise between having a lower test error and good generalization capability.

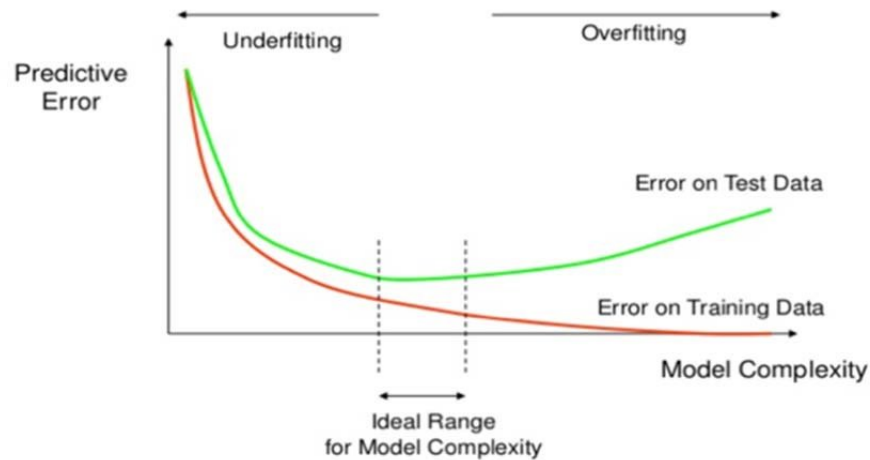


Figure 1.1: Representation of the relationship between train and test error affecting the underfitting-overfitting tradeoff. Source: Al-Behadili et al., 2018

Sticking with the example of supervised learning that focuses on learning a mapping like $f(\mathbf{x}; \theta) = y$, any of those algorithms is considered to perform an estimation of the set of true parameters θ , which are unknown. Assume to have a dataset $D = \{\mathbf{x}_1, \dots, \mathbf{x}_M\}$ of M data points, which are drawn by an unknown DGP and are IID. We define a parameter estimator as the set $\hat{\theta}$ that describes a valid mapping $f(\mathbf{x}; \hat{\theta}) + \epsilon = y$ for each available data points, where ϵ is the part that the model is not able to explain. A large part of the machine learning applications consists in finding the approximation of $f(\mathbf{x}; \theta)$ with a function estimate $f(\mathbf{x}; \hat{\theta})$. In some cases, function estimate or parameter estimate are interchangeable since estimating the parameters is instrumental for estimating the desired mapping from \mathbf{x} to y . Henceforth, we will refer to it simply as the estimator. As anticipated in the previous

section, The estimator is the result of the training process that iteratively searches for suitable values of $\hat{\theta}$ so that it is as close as possible to the true unknown values θ and minimizes the estimation error ϵ .

When considering an estimate, one should look at two important properties: bias and variance. The bias for an estimator $\hat{\theta}_N$ is defined as:

$$\text{Bias} \left(f(\mathbf{x}; \hat{\theta}) \right) = \mathbb{E} \left(f(\mathbf{x}; \hat{\theta}) \right) - f(\mathbf{x}; \theta) \quad (1.1)$$

where the expectation is computed over the set of training samples of length M . We said that an estimator is unbiased if the bias is zero so that $\mathbb{E}(f(\mathbf{x}; \hat{\theta})) = f(\mathbf{x}; \theta)$. Otherwise, an estimator is said to be asymptotically unbiased if the bias goes to zero when the amount of training examples goes to infinity, so that $\lim_{M \rightarrow \infty} \mathbb{E}[f(\mathbf{x}; \hat{\theta})] = f(\mathbf{x}; \theta)$. The bias reflects how well the estimator approximates the actual value on average.

On the other hand, the variance of an estimator is simply the variance of the estimated function,

$$\text{Var}(f(\mathbf{x}; \hat{\theta})) = \mathbb{E}[(f(\mathbf{x}; \hat{\theta}) - \mathbb{E}[f(\mathbf{x}; \theta)])^2], \quad (1.2)$$

The variance measures how much one would expect the estimator to vary if one computes it over another sample of the data from the same DGP. It is of practical use because, generally, estimations come from a finite number of samples, since collecting a huge amount of data points is either infeasible or costly and time-consuming, and the variance provides a measure of how much one would expect the estimate to vary when independently sampling multiple times the dataset from the same DGP. Bias and variance represent two different sources of error when obtaining the desired estimate and therefore provide an evaluation criterion for the variability through a sample. While bias measures the expected deviation from the true value to be estimated, on the contrary, variance measure the deviation from the expected estimation, which is likely to occur if one changes the training examples.

One would expect an estimator to have a low bias and a low variance. However, in statistics and machine learning, this ideal situation is not always achievable, and, as a consequence, there is a known tradeoff between the two sources of error (Kohavi, Wolpert, et al., 1996). It is possible to show that in estimating $f(\mathbf{x}; \hat{\theta})$, the expected mean square error (MSE) with respect to the labels over the test set can be decomposed as a sum of two quantities:

$$\begin{aligned}
 & \mathbb{E} \left[f(\mathbf{x}; \hat{\theta}) - y \right]^2 = \\
 & \mathbb{E} \left[f(\mathbf{x}; \hat{\theta})^2 - 2f(\mathbf{x}; \hat{\theta})y + y^2 \right] = \\
 & \mathbb{E} \left[\hat{f}(x)^2 \right] - 2\mathbb{E} \left[f(\mathbf{x}; \hat{\theta}) \right] f(x) + f(x)^2 = \\
 & \mathbb{E} \left[\hat{f}(x)^2 \right] - \mathbb{E} \left[f(\mathbf{x}; \hat{\theta}) \right]^2 + \mathbb{E} \left[f(\mathbf{x}; \hat{\theta}) \right]^2 - 2\mathbb{E} \left[f(\mathbf{x}; \hat{\theta}) \right] f(x) + f(x)^2 = \\
 & \mathbb{E} \left[\hat{f}(x)^2 \right] - \left(\mathbb{E} \left[f(\mathbf{x}; \hat{\theta}) \right]^2 + \mathbb{E} \left[f(\mathbf{x}; \hat{\theta}) \right]^2 - y \right)^2 = \\
 & \text{Var} \left(f(\mathbf{x}; \hat{\theta}) \right) + \left[\text{Bias} \left(f(\mathbf{x}; \hat{\theta}) \right) \right]^2
 \end{aligned} \tag{1.3}$$

which are respectively the variance of the function estimate and its squared bias. The expected MSE on the test set refers to the average MSE that one would obtain repeatedly estimating the function f using numerous training sets and testing every time at the same test example x_0 . The overall expected test MSE can be computed by averaging the expected MSE over all possible data points in the test set. In general, a model with more parameters, hence fewer assumptions, can obtain estimates with lower bias, i.e., to approximate well the underlying phenomena, although with higher variance at the same time. The relative rate of change of these two quantities, when the model varies, produces either an increase or decrease of the MSE in Eq. 1.3. Friedman et al., 2001 provide a detailed explanation of this trade-off, which is visualized in Figure 1.2 for a two dimensional case. On the one hand, models with low bias and high variance can be easily obtained, for instance, a curve passing through each data point as in the left panel of Figure 1.2, where the function parameters are close to the true values, but the estimator is too attached to a specific set of data points.

On the other hand, it is also easy to come up with a model with low variance and high bias, i.e., a constantly horizontal line as in the central panel of Figure 1.2, for which the function parameters are far from the true values, although the variance of the estimate is expected to be low because the function is not tied specifically to any sampled data points. The former case shows a lack of generalization capability, while the latter exhibits a lack of fitting capability. The challenge in machine learning is to find the model which keeps both measures low, as in the right panel of Figure 1.2, where the

estimated function does not pass precisely through each data point, but it is close enough to each of them to represent the underlying DPG, resulting in a model that has both modeling and generalization capabilities. In what follows, we will present the neural networks type of models, which are highly flexible, hence able to eliminate the source of bias, but which often exhibit significant variance in their estimate.

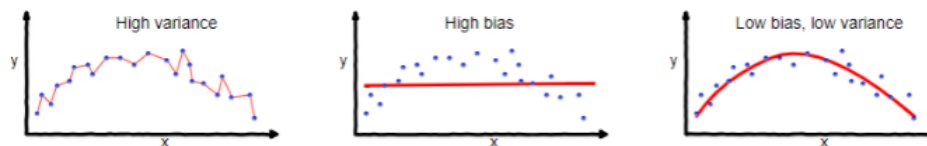


Figure 1.2: Visual explanation of bias and variance concepts in a two dimensional case.

1.3 Background Theory

The work of Vapnik and Chervonenkis, 1971 represents the foundation of the statistical learning theory and provides a way to quantify the capacity of a machine learning algorithm (see Abu-Mostafa et al., 2012 for extensive coverage of these concepts). They introduce the Vapnik-Chervonenkis (VC) dimension, which measures the capacity of a binary classifier. It is defined as the largest possible size of the training set so that the model can classify any of these data points. The VC dimension provides an upper bound for the gap between training and test error directly proportional to the model capacity and inversely proportional to the number of training examples. However, such bound holds in the simple case of a binary classifier, while it is challenging to provide similar results with more complex models. There is a discrepancy between representational, i.e., theoretical, and practical capacities in many machine learning models, mainly due to the underlying optimization algorithms, which often are not guaranteed to converge and have very little theory supporting them. Hence, some models would, in principle, be able to solve the proposed task, but the optimization procedure fails in reaching the goal. Neural networks are an excellent example of this misalignment in capacity (see Section 1.4).

Wolpert and Macready, 1997 introduces the *No Free Lunch Theorem*, which states that if it would be possible to collect an infinite amount of data, and hence average over all data generating distributions, the performance metric of any machine learning model over previously unseen data points would be the same. This statement implies that a universal outperforming model is impossible. A model outperforms the others only because it is particularly suited to solve the structure of the specific problem under consideration. Therefore, in principle, no model is universally better than the other. In a more practical approach, where one can make assumptions about the DGP and restrict the set of distributions where the training samples are drawn, some models become better than others. The need for a model that suits the task to perform is a motivation for the dominance of certain types of models, such as neural networks, in the context of all the different machine learning paradigms.

1.4 Early models for artificial neural network

In this section, we present the neural network paradigm, which has been central in the development of modern machine learning. We begin with a brief introduction of pioneering models that set the basis for the existence of neural networks. Then we introduce the feedforward neural network as the basic structure of the modern field of deep learning, and we outline all the design decisions that one should take in dealing with this type of model. Additionally, we briefly recap different neural network architectures, each with their advantages and disadvantages, and the main algorithm used to train them.

The idea of a biologically inspired model comes from McCulloch and Pitts, 1943, which presents a simple mathematical model that tries to mimic the work of the neuron in the complex structure of the human brain. Figure 1.3 shows the method of operation of the artificial neuron, which accepts binary values as input and produces a single binary output according to a certain threshold. Two different functions operate on the binary inputs: a function $g(x_1, x_2, x_3, \dots, x_n) = g(\mathbf{x}) = \sum_{i=1}^n x_i$, which sums all the binary values of the input, and a function $f(g(\mathbf{x}))$ which returns one if the aggregated sum is greater than θ , otherwise zero.

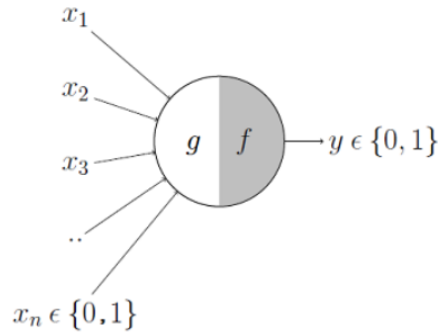


Figure 1.3: Artificial neuron representation of McCulloch and Pitts, 1943

An extension of this model is provided by the perceptron of Rosenblatt, 1958, shown in Figure 1.4. Such model can work with either a binary or a non-binary vector of input values \mathbf{x} and computes a weighted sum $g(\mathbf{x}, \mathbf{w}) = \sum_{i=1}^n w_i x_i$ using the vector of weights \mathbf{w} , before applying a similar threshold θ as in the artificial neuron model.

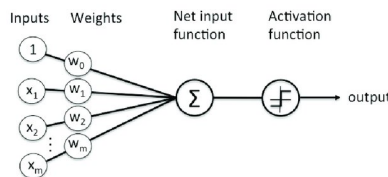


Figure 1.4: Perceptron representation of Rosenblatt, 1958

The scalar parameter θ controls the neuron's activation and allows representing a linearly separable function of the data. Being linearly separable means that a hyperplane exists that splits the input space into two parts, containing all the points of a specific class. It is said in this case that there is a line that can perfectly separate all the data points. Minsky and Papert, 1969 criticizes this type of model because they are not able to classify data points that are not linearly separable. In this regard, the artificial neuron and the perceptron can model several functions of binary values, also called Boolean functions, but they fail, for example, with the exclusive-or (XOR) function. This operation on binary values x_1 and x_2 returns one when either $x_1 = 1$ or $x_2 = 1$, otherwise returns zero. This set of data is not linearly separable, as clarified from Figure 1.5, since there is no line in the two-dimensional space that can separate the two classes of outputs.

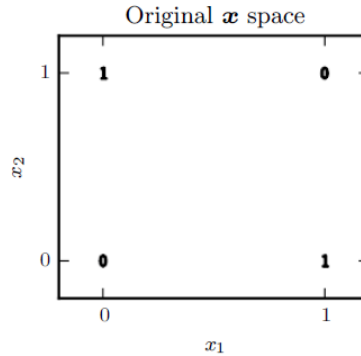


Figure 1.5: Representation of the XOR function in the two-dimensional space (Goodfellow et al., 2016). The number in the figure represents the expected model output for the pair of points x_1, x_2

The limitation of these models is the possibility to assume just binary output values (0 or 1) due to the threshold θ . In the case of $x_1 = 0$ ($x_1 = 1$), we should expect the model output to increase (decrease) with the other data point x_2 . This is not possible with the perceptron because the model assumes the weight w_2 to be fixed independently of the value of x_1 . The XOR case is often representative of a class of more complex problems that a linear classifier is not able to solve, such as nonlinear pattern recognition (Friedman et al., 2001), combinatorial game theory (Albert et al., 2019) and cryptography (Tuyls et al., 2005). For this purpose, a feedforward neural network comes in help, as outlined in the next section.

1.5 Feedforward neural networks

Feedforward neural networks represent a turning point for the development of modern machine learning models. All the three machine learning paradigms described in Section 1.1 widely employ these models. To provide a clear explanation, we introduce the neural network models in the supervised learning framework. Then, in the next chapter, we will describe their use in a reinforcement learning context.

Recalling the aim of a supervised learning model to learn the mapping $y = f(\mathbf{x})$, i.e. learning the output y by exploiting the data points contained in \mathbf{x} , neural networks are able to learn a parametrization of the unknown function $f(\mathbf{x})$ through a set of parameter θ , so that the mapping becomes

$$y = f(\mathbf{x}; \theta).$$

Figure 1.6 provides a visualization of the basic structure of a multilayer feedforward neural network. It partially resembles the structure of a perceptron, hence the name multilayer perceptron, which is sometimes used in place of a feedforward neural network. As the name suggests, one of the differences from the model of Rosenblatt, 1958 is the presence of an intermediate layer, called the hidden layer. The architecture of a neural network includes an input layer, which is the first set of neurons on the left of the Figure 1.6, already present in the original perceptron. Then there is a hidden layer, whose elements are connected by directed edges to all the elements of the input layer and all the elements of the subsequent output layer. These edges represent the parameters θ to approximate the desired function.

We refer to this model as *feedforward* because of the way that information flows from input to output since there is no recursion at any point and no feedback loop that connects the previous outputs to the current inputs. Since each unit, also called a neuron, receives information from all the neurons in the previous layer, one can think of a layer as a function that takes a vector as input and returns another vector as output. The former vector has a dimension equal to the number of neurons in the previous layer and the latter equal to the number of neurons in the current layer.

The formulation in terms of neurons that elaborate and aggregate information motivates the adjective *neural*, because of the loose biological inspiration already present in the first formulation of these models. The number of units per layer varies depending on the task and assume different meaning depending on the layer. Increasing the units in the input layer means that the model can access more information, i.e., each vector \mathbf{x} in the set of training examples $\{\mathbf{x}_1, \dots, \mathbf{x}_M\}$ contains more information. The number of units in the output layer strictly depends on what the neural network wants to achieve. The dimension of the corresponding labels y , which can be either a scalar or a vector, already provides information about the size of these two layers, defining respectively a single output or a multi-output neural network. On the contrary, the number of hidden layer units may increase the model capacity, but it is not directly related to the dataset that the model experiences.

We remark that Figure 1.6 represents a composition of different functions like $f(\mathbf{x}) = (f_3 \circ f_2 \circ f_1)(\mathbf{x})$, where f_1, f_2 and f_3 are respectively the function describing the input, the hidden and the output layer. The name *network* comes from the opportunity to represent this expression as a directed graph.

Many other hidden layers can be stacked to compose the function, hence defining the depth of the model. Therefore, feedforward neural networks with many hidden layers, and possibly wide, i.e., with many neurons, are generally considered as *deep learning* models.

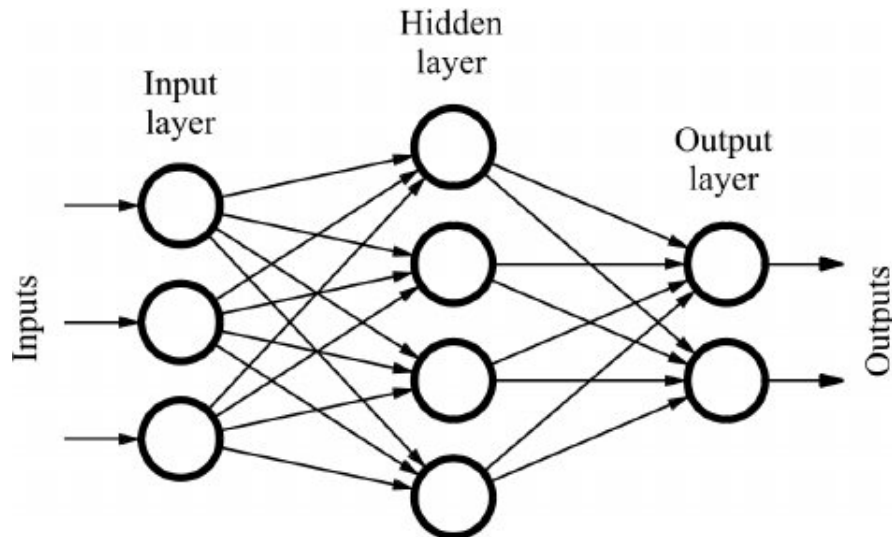


Figure 1.6: Representation of a feedforward neural network with three inputs, two outputs and a hidden layer (From Fadja et al., 2018).

A feedforward neural network overcomes the perceptron's limits, which can only learn a linearly separable function of the data by returning only one of two values (0 or 1) because of the hard-limit threshold function. Instead, a neural network not only produces a different representation of the input space through hidden layers but can also output real values. The former characteristics allow learning a broader set of functions of the inputs that are nonlinear in the neural network parameters.

We can think about the set of these parameters used to approximate the function as composed of two distinct subsets: ψ that helps find the new representation on the feature space and ω for the output layer activation. The approximation becomes $y = f(\mathbf{x}; \theta) = f(\mathbf{x}; \psi, \omega) = \theta(\phi(\mathbf{x}; \psi); \omega)$, where the function ϕ is the learned mapping from the input space to the new feature space and the function θ is the activation function applied in the output layer. Given this general form, we can distinguish two different tasks that a multilayer perceptron can perform. The first is the regression problem when the function θ is linear in the parameters ω . This type of task consists of a

generalization of a classical linear regression problem, where the linear model is applied to the input space learned by the hidden layers $\phi(\mathbf{x}; \psi)$. The second one is the classification task when the function θ is nonlinear, and we need the neural network to output the probability distribution for all the possible outcomes of the problem.

The intermediate layers are called *hidden* because the training dataset does not suggest what these layers should do to elaborate the information, which is entirely retrieved by the mechanism used to learn the representation function ϕ . Elaborating a new representation of the original input space is known to help the model reach a better generalization on the test set and solve previously impossible problems, such as the XOR one, by applying a linear model to a nonlinear transformation of the original data.

Nonetheless, linear models have limited capacity since they can represent only linear functions. Generally, they are solved through closed-form solutions or convex optimization methods. The introduction of nonlinearity in the neural network makes these methods more effective for solving the problem. For these reasons, neural networks are generally trained through gradient-based algorithms, which are methods employed for various machine learning models, especially when the dataset is quite large. Despite the differences in the optimization procedures, training a neural network requires some design choices similar to those required by a linear model, such as choosing an error function, otherwise called cost function, to optimize.

1.6 Design choices

This section explores possible choices in designing a neural network model. This exposition outlines the flexibility of this family of models that, through minor changes, can solve very different problems.

Recalling the concept of an estimator from Section 1.2, machine learning uses some principles to search for a proper function estimate to avoid the need to make many possible guesses and analyze their bias and variance singularly. We introduce the error function as the performance metric in Section 1.2 to measure the discrepancy between the estimator and the true value of the parameters. We again assume to have a dataset $D = \{\mathbf{x}_1, \dots, \mathbf{x}_M\}$ composed by M data points of which we can estimate a common mapping for each \mathbf{x}_i to its respective label y_i , in a supervised learning setting. In this case, a common

choice for the error function is the mean square error (MSE) function:

$$J^{\text{MSE}}(\theta) = \frac{1}{M} \sum_{i=1}^M \left(f(\mathbf{x}_i; \hat{\theta}) - y_i \right)^2 \quad (1.4)$$

However, sometimes it is useful to learn a conditional distribution $p(y | \mathbf{x})$ of the labels given the input data, instead of learning the best function f that ties \mathbf{x} to y . In this case, the model to learn becomes probabilistic, and it is usually treated under the maximum likelihood estimation (MLE) principle, where the error function is

$$J^{\text{MLE}} = -\mathbb{E}_{\mathbf{x}, y \sim p(\mathbf{x}, y)} \log p(y | \mathbf{x}). \quad (1.5)$$

with data drawn from the joint distribution $p(\mathbf{x}, y)$. It is still possible to learn the distribution $p(y | \mathbf{x})$ with a neural network. If we express the conditional distribution as a Gaussian, centered on the function estimate and with a fixed variance σ^2 , $p(y | \mathbf{x}) = \mathcal{N}\left(y; f(\mathbf{x}; \hat{\theta}), \sigma^2\right)$, one can show that the MLE procedure gives exactly the same result of minimizing the MSE. Given that the sampled data points are assumed to be IID, the conditional log-likelihood is

$$\sum_{i=1}^M \log p(y_i | \mathbf{x}_i; \theta) = -M \log \sigma - \frac{M}{2} \log(2\pi) - \sum_{i=1}^M \frac{\left(f(\mathbf{x}_i; \hat{\theta}) - y_i \right)^2}{2\sigma^2}. \quad (1.6)$$

It is immediate to see that maximizing Eq. 1.5 is equal to minimizing Eq. 1.4 up to a constant. This justifies the use of the MSE as an error metric when the output of the neural network is real-valued statistics of the target distribution.

Beyond the case of a Gaussian conditional distribution, the general formulation of the error function in terms of MLE adapts to various problems in machine learning, such as regression but also classification. One can show that maximizing Eq. 1.5 is also equivalent to minimizing the cross-entropy between the labels and the data points in the dataset. This expression of the cost function is helpful when the neural network's output is binary data and we are solving a classification problem. The choice of the error function for a neural network is similar to other machine learning and even econometrics models and needs to be adapted to the type of problem we want to solve. In this context, the MLE formulation allows for the expression of deterministic

or probabilistic models, covering the wide space of different types of problems in the machine learning domain.

Similar to the error function, the choice of the output layer units is connected to the task type performed by the neural network. Selecting the proper type of units in a layer means choosing whether and which type of nonlinearity is used to solve the problem. Recalling Figure 1.6, we interpret the results of the last hidden layer, even though in the proposed example there is just one, as a modified representation of the original data $\mathbf{h} = f(\mathbf{x}, \theta)$. This hidden representation goes to the final layer, which, depending on the type of unit, elaborates it and provides the model's output. Generally, when we want the neural network to output real-valued numbers, linear units are employed. These units simply takes the intermediate output \mathbf{h} and apply an affine function, resulting in an output $\hat{y} = \mathbf{W}^T \mathbf{h} + \mathbf{b}$. Linear units are useful to output the mean of a Gaussian distribution $p(y | \mathbf{x}) = \mathcal{N}(y; \hat{y}, \mathbb{I})$.

If we are solving a classification problem instead of a regression one, the choice of the output units depends on how many classes are included in the task. In the case of a binary choice between two classes, the output layer is composed of sigmoid units. These units still apply the affine function to the intermediate output \mathbf{h} and the result is passed to a nonlinear logistic sigmoid function so that $\hat{y} = \sigma(\mathbf{W}^T \mathbf{h} + \mathbf{b})$. The conditional function to be modeled in this case is a Bernoulli $p(y = 1 | \mathbf{x})$, hence the values must lie in the interval $[0, 1]$ to represent a probability, as guaranteed by the sigmoid function $\sigma(x) = \frac{1}{1 + \exp^{-x}}$. Instead, when the output of a problem can be chosen among n multiple classes, we use the softmax function as a generalization of the sigmoid function for binary classes since, in this case, the output is a vector \mathbf{y} instead of a single value. Such function normalizes the values of each vector component in the interval $[0, 1]$ so that they sum to one. This is achieved by computing $\text{softmax}(\mathbf{x})_i = \frac{\exp^{x_i}}{\sum_j \exp\{x_j\}}$.

The type of units described can be used equivalently in a hidden layer, even though linear and softmax units are rarely used. The former does not introduce nonlinearity, which we usually want to achieve in the intermediate layer. At the same time, the latter is used if we want the model to choose among n possible values for some hidden variable. While the choice of the output units is relevant to defining the type of task we want to perform, the choice of the hidden units describes the type of nonlinearity we want to model. There is no guarantee that any type of nonlinearity works better in any case than the other. Therefore the representational power of the hidden

layers of the neural network is an active field of research based on heuristics and extensive simulation tests.

Similar to the output units, the hidden units differ in the type of nonlinear activation function used after calculating the affine function on the output of the previous layer. Sigmoid functions were largely used in the first successful deep learning implementation (Krizhevsky et al., 2012), together with the hyperbolic tangent function, which is strictly linked to the former since $\tanh(x) = 2\sigma(2x) - 1$. Nowadays, these two type of activation functions are less popular than rectified linear units (ReLU) and their variants that applies the function $g(\mathbf{h}) = \max\{0, \mathbf{h}\}$ on top of the affine transformation. Several generalizations of the ReLU exist, and it is a matter of heuristic and case-by-case choice to find the one that performs better for a specific task. They exist because gradient-based learning fails, for example, where the ReLU activation outputs zero. Hence, they help obtain a gradient everywhere by providing a non-zero slope in that part of the activation function domain. These types of units ease the model's optimization thanks to their similarity to linear units, whose gradient never saturates. We provide a more detailed explanation of their benefits from sigmoid and hyperbolic tangent activation after introducing the optimization techniques of a neural network and its related algorithms in the coming sections of this chapter.

1.7 Architecture of a neural network

The neural network's architecture consists of its high-level configuration: the number of different layers, the type of such layer, and the composition, i.e., the type of units in each of them. The architecture somehow recaps all the design choices carried out to prepare the model before the learning process and deserves a special mention because it may lead to different type of neural networks for multiple purposes.

As explored in Section 1.5, the simplest neural network organizes into a group of units called layers, where each layer is a function of the previous one. Hence, we can express the feedforward neural network as a chain of affine and nonlinear functions that goes through all the network structures up to the output layer. Defining \mathbf{x} as the input, the first and the second layer

can be written as

$$\mathbf{h}_1 = g_1 (\mathbf{W}_1^T \mathbf{x} + \mathbf{b}_1), \quad (1.7)$$

$$\mathbf{h}_2 = g_2 (\mathbf{W}_2^T \mathbf{h}_1 + \mathbf{b}_2) \quad (1.8)$$

where $\mathbf{W}_1, \mathbf{b}_1$ and $\mathbf{W}_2, \mathbf{b}_2$ represent respectively the parameter and the biases associated with those layers. This chain may eventually continue and is long in the case of modern deep neural networks. The depth and the width of the neural network are task-dependent and adapt to the problem complexity, as most of the design choices for this family of models.

Since the beginning of machine learning research, some architecture and network types have been proposed. Feedforward neural networks represent the primary building block and inspiration for another kind of network structure. Delving into the details of these structures is out of the scope of this work since we are going to use just a feedforward neural network in our empirical research. By the way, we refer respectively to Chapters 9 and 10 of Goodfellow et al., 2016 to get more insights about two prevalent neural network structures: convolutional neural networks (LeCun, Bengio, et al., 1995) and recurrent neural networks (Graves, 2012; Hochreiter and Schmidhuber, 1997). The former has been extensively employed to solve computer vision and object recognition tasks. At the same time, the latter played a huge role in the achievement of machine translation and text recognition tasks since they are handy for modeling data sequences.

The need for these different architectures is motivated by different structures in which the data are arranged. However, theorems state that a feedforward neural network with just one hidden layer is potentially sufficient to learn a mapping in the training set. In the next section, we will detail these theoretical statements and hence motivate the need for more complex neural networks.

1.8 Universal approximation theorem

Neural networks exhibit a peculiar approximation capability since they need few assumptions to learn a specific family of nonlinear functions. Hidden layers in their structure lead to a universal approximation framework. This characteristic eases the modeling task since we do not need prior information on the training data's specific nonlinear function. Instead, we can be as

general as possible and train a neural network, giving to the hidden layers the role of learning such nonlinearity.

The Universal Approximation Theorem (Cybenko, 1989; Hornik et al., 1989) states that a feedforward neural network with a linear output layer and just one hidden layer can approximate with a negligible amount of error any continuous function on a closed and bounded subset of \mathbb{R}^n . This theoretical result is guaranteed if the activation function in the hidden layer units is a function that saturates, i.e., becomes very flat, for large positive and negative values, like the sigmoid or the hyperbolic tangent functions. The theorem is also generalized to another type of common nonlinearity, such as the rectified variants (Leshno et al., 1993).

However, the theoretical results only imply that a large enough feedforward neural network can represent the function we would like to approximate. The optimization step, i.e., the model's training, often does not result in the negligible amount of error expected from the theorem. This issue relates to the fact that neural networks are usually trained through gradient-based optimizers (see Section 1.9) and may fail in finding the proper set of parameter values. In addition, there is the risk of overfitting when the model capacity is too high, leading to local minima of the objective function and hence to a suboptimal training result. Therefore, even if the universal approximation theorem states that a powerful universal approximator exists, there is no guarantee to retrieve the exact parameters in a finite amount of time, neither there is specific guidance on how to structure the model (see Section 1.6). We refer to Goodfellow et al., 2016 for a detailed review of the literature on the tradeoff between the neural network width, as stated by the theorem, and the neural network depth, which in contrast, employs many hidden layers. The latter has performed better on practical tasks since a model with a single huge hidden layer becomes increasingly difficult to train.

In the next section, we discuss the optimization process of a neural network, where we shed light on the set of calculations needed and the proper algorithms to achieve that.

1.9 Optimizing the neural network parameters

As anticipated in the previous sections, most neural networks gather information through the input layer and then propagate it through some hidden layer to get the model's output. This process is referred to as forward-propagation, and it is commonly used to evaluate the performance of a neural network model by computing the predicted output and comparing it to the given target using a selected error function (see Section 1.6).

In order to find the set of parameters that allows us to approximate the desired function closely, we need a way to change them and gradually obtain outputs closer to the actual values. The gradient-based approach for training a neural network consists of computing the gradient of the error function with respect to the neural network's parameters and changing them following the inverse direction of such gradient to minimize the error. The optimization process of a neural network is therefore composed of two iterative steps, repeated until the desired convergence to the problem solution is reached:

1. calculating the gradient of the error function with respect to all the parameters of the neural networks;
2. using such gradient to update the parameter of the neural network.

For what concerns the first step, it is simple to derive the analytical expression of the gradient of the error function. In contrast, it is tough to evaluate such an expression and obtain the gradient to update the parameters. Rumelhart et al., 1986 proposed the backpropagation algorithm as an efficient way to compute the considerable amount of partial derivative required, which is known to increase with the width and the depth of the neural network. In this section, we explain its use to compute derivatives of the error function with respect to the neural network parameters, even though the algorithm is general enough for differentiating any continuous function. In the field of machine learning, it is also used for computing derivatives with respect to the inputs of the neural networks in order to perform sensitivity analysis of the model (Dixon et al., 2020; Hugu and Savine, 2020).

Usually, the neural network weights are randomly initialized as a starting point. Several initialization techniques are available and represent one of the design choices for the neural network to implement. A common way to initialize the random weights of a neural network is provided by Glorot

and Bengio, 2010, which states that the signal needs to flow correctly in both directions through the network architecture. To achieve that, the input variance of each layer needs to be almost equal to the output variance. Also, the gradients need to have equal variance before passing through a layer in the reverse direction. This condition is not guaranteed unless each layer has the same number of neurons, which is not always the case, although one can have a reasonable compromise. The Glorot initialization, from the primary author name, considers the random weights initialized equivalently as:

- Normal distribution $\mathcal{N}(0, \sigma^2)$, where $\sigma^2 = \frac{1}{\text{fan}_{\text{avg}}}$;
- Uniform distribution $\mathcal{U}(-r, r)$ with $r = \sqrt{\frac{3}{\text{fan}_{\text{avg}}}}$

where $\text{fan}_{\text{avg}} = \frac{\text{fan}_{\text{in}} + \text{fan}_{\text{out}}}{2}$, being fan_{in} and fan_{out} respectively the number of inputs and the number of outputs in each layer of the network. Other parameters initialization methods works in similar way, but modifies the scale of the variance in the Gaussian case or the width of the interval in the uniform case. For more details and references, we refer to Géron, 2019.

The backpropagation algorithm allows the training of a neural network through the chain rule of calculus. After each forward pass through a network, backpropagation performs a backward pass that allows modifying the neural network parameters. In the example, we delve into the optimization process of a feedforward neural network that accepts N values as inputs, outputs K values, and has L layers, whose $L - 2$ are hidden layers each composed by M neurons.

The training examples accepted by the model are in the form of a vector $\mathbf{x} \in \mathbb{R}^N$. The values associated with the neurons for each hidden layer l from 1 up to $L - 1$ are computed as

$$\mathbf{z}^{(l+1)} = \mathbf{W}^{(l)}\mathbf{h}^{(l)} + \mathbf{b}^{(l)} \tag{1.9}$$

$$\mathbf{h}^{(l+1)} = a(\mathbf{z}^{(l+1)}), \tag{1.10}$$

where $\mathbf{W}^{(l)}$ and $\mathbf{b}^{(l)}$ are respectively the matrix of weights and the associated vector of biases arranged by layers. Specifically, $\mathbf{W}^{(l)}$ connects the layer l to the layer $l + 1$. We remark that when $l = 1$ we have $\mathbf{h}^{(1)} = \mathbf{x}^{(1)}$ since the first layer is the input layer. The function $a(\cdot)$ in Eq 1.10 represents an activation function which is typically a nonlinear function of its inputs that we explain in more details in the next section.

The output in the final layer, which is computed as $\hat{\mathbf{y}} = \mathbf{W}^{L-1}\mathbf{a}^{L-1}$, is compared to the true output \mathbf{y} through an error function $J(\hat{\mathbf{y}}, \mathbf{y})$. The backpropagation algorithm uses the values of the error function to minimize them by adjusting the network's weights and biases and getting closer to the true output.

The partial derivative of the error function for a specific weight $w_{jk}^{(l)}$, that connects the j -th unit of the layer $l - 1$ to the k -th unit of the layer l , is computed as

$$\frac{\partial J}{\partial w_{jk}^{(l)}} = \frac{\partial J}{\partial z_j^{(l)}} \frac{\partial z_j^{(l)}}{\partial w_{jk}^{(l)}} \quad (1.11)$$

where $z_j^{(l)} = \sum_{k=1}^M w_{jk}^{(l)} a_k^{l-1} + b_j^{(l)}$, which derived with respect to $w_{jk}^{(l)}$ results in $\frac{\partial z_j^{(l)}}{\partial w_{jk}^{(l)}} = a_k^{l-1}$. It follows that we can write the Eq. 1.11 as

$$\frac{\partial C}{\partial w_{jk}^{(l)}} = \frac{\partial C}{\partial z_j^{(l)}} a_k^{(l-1)}. \quad (1.12)$$

The same calculation applies for the partial derivative of the error function with respect to a single bias, so that

$$\frac{\partial C}{\partial b_j^{(l)}} = \frac{\partial C}{\partial z_j^{(l)}}, \quad (1.13)$$

because $\frac{\partial z_j^{(l)}}{\partial b_j^{(l)}} = 1$.

In this way, we can retrieve any partial derivative with respect to any weight or bias included in the neural network. We notice that the calculation proceeds backward, as the algorithm's name suggests. For a detailed explanation of the backpropagation algorithm and its variants, we refer to (Goodfellow et al., 2016; LeCun et al., 1989).

Once the gradient is computed, we can use every partial derivative to update the value of the neural network weights according to a selected optimization algorithm, which is commonly a gradient descent algorithms. We refer to Ruder, 2016 for a comprehensive overview of these methods and the different optimization algorithms available for training neural networks. The general idea of a gradient descent method is to update the weights towards the inverse direction of the gradient as follows,

$$\mathbf{W}^{(l)} = \mathbf{W}^{(l)} - \eta \cdot \nabla_{\mathbf{W}^{(l)}} J(\mathbf{W}^{(l)}), \quad (1.14)$$

where the scalar η is a parameter called learning rate that scales the size of the step to the pointed direction and $\nabla_{\mathbf{W}^{(l)}} J(\mathbf{W}^{(l)})$ is the gradient of the error function with respect to the set of weights $\mathbf{W}^{(l)}$ defined as

$$\nabla_{\mathbf{W}^{(l)}} J(\mathbf{W}^{(l)}) = \begin{pmatrix} \frac{\partial J(\mathbf{w}_{11}^{(l)})}{\partial \mathbf{w}_{11}^{(l)}} & \cdots & \frac{\partial J(\mathbf{w}_{1M}^{(l)})}{\partial \mathbf{w}_{1M}^{(l)}} \\ \vdots & \vdots & \vdots \\ \frac{\partial J(\mathbf{w}_{N,1}^{(l)})}{\partial \mathbf{w}_{N,1}^{(l)}} & \cdots & \frac{\partial J(\mathbf{w}_{N,M}^{(l)})}{\partial \mathbf{w}_{N,M}^{(l)}} \end{pmatrix}. \quad (1.15)$$

The updated equation for each set of biases in the network is analogous to that of the weights.

1.10 Gradient issues during training

Even though it is true that no activation function is considered to be universally better than another, some of them may cause gradient issues during the training process of a neural network. We have highlighted that neural networks are primarily trained using a gradient-based optimizer. Therefore it becomes crucial to avoid problems and numerical instabilities when computing the gradient of the error function with respect to a large number of parameters of the family of models. This section summarizes some of the numerical issues one may encounter when a particular activation function is used in place of the others and, eventually, how to avoid those issues.

As already discussed, an activation function is applied componentwise to each unit in a single layer of the neural network architecture to get the desired output. Its role is to scale the result of an affine function by adding a nonlinear effect. Here we mainly compare the two groups of activation functions from Section 1.6: the logistic sigmoid and the hyperbolic tangent functions in contrast to the rectified linear activation functions.

The left panel of figure Figure 1.7 shows the range of values that the logistic sigmoid and the hyperbolic tangent function assume, while the right panel of the same figure shows their first derivatives on the same domain.

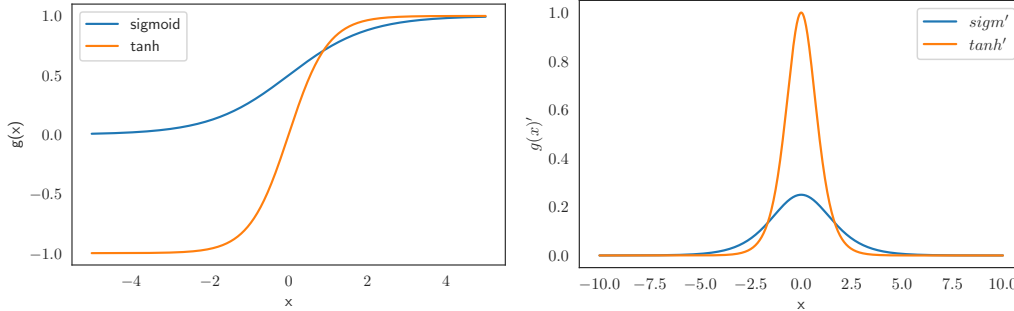


Figure 1.7: Sigmoid and hyperbolic tangent function on the left panel and their first derivatives on the right panel.

When computing the gradient of the error function, a significant issue arises when these two functions are saturated, i.e., they become flat in correspondence with largely negative or positive values. In order to show the problem, we resort again to the chain rule of calculus employed by the back-propagation algorithm. We assume we want to train a neural network with L layers by minimizing an error function $J(\theta)$. In the following examples of this section, we also account for the biases to propose a complete overview. Indicating all the weights of the output layer as w^L , then the general formula for the partial derivative of $J(\theta)$ with respect to them is:

$$\frac{\partial J}{\partial w^L} = \frac{\partial J}{\partial a^L} \frac{\partial a^L}{\partial z^L} \frac{\partial z^L}{\partial w^L} \quad (1.16)$$

where z^L is the output of an affine function and a^L is the output of an activation function. Except for the partial derivative $\frac{\partial a^L}{\partial z^L}$, the other components of the chain rule are not affected by the choice of the activation function. If the selected nonlinear function is the sigmoid function $\sigma(x) = \frac{1}{1+\exp^{-x}}$, we can write such partial derivative as:

$$\frac{\partial a^L}{\partial z^L} = \sigma'(wa^{L-1} + b). \quad (1.17)$$

It is easy to show with a simple calculation that the first derivative of the sigmoid function turns out to be

$$\sigma'(x) = \frac{e^{-x}}{(e^{-x} + 1)^2} = \sigma(x)(1 - \sigma(x)), \quad (1.18)$$

which means that a large positive or negative input x into the sigmoid function returns a value that is almost zero, as it is also shown the right panel of Figure 1.7. Hence, by computing the flow of the gradient, a value close to zero is multiplied with the other partial derivatives in the chain rule

$$\frac{\partial \mathcal{L}}{\partial w^L} = \frac{\partial \mathcal{L}}{\partial a^L} \underbrace{\frac{\partial a^L}{\partial z^L}}_{\text{almost zero}} \frac{\partial z^L}{\partial w^L}. \quad (1.19)$$

If there are many partial derivatives of the error function with respect to the weights that are zero or close to zero, their updates through a gradient descent optimizer add no information because the values of the weights do not change, and the network cannot reach an optimal configuration by minimizing its error function. This behavior of the gradient flow in a neural network is referred to as the *vanishing gradient* problem.

A similar problem applies to the hyperbolic tangent function, whose first derivative exhibits the same behavior as the derivative of the sigmoid function, although on a slightly different range of the domain. The problem of the vanishing gradients exacerbates when the neural network has many hidden layers, because a layer on top of the architecture learns faster than a layer at the bottom, since the former has comparably bigger gradients. As a matter of fact, the error function depends on the variation of the weights connected to every other layer stacked on it. Therefore, the derivative with respect to the weights in a layer at the bottom of the network exhibits a chain of dependence from a number of other weights in the subsequent layer of the neural network architecture. For instance, in a four hidden layer architecture, earlier layers in a feedforward neural network reuse computation from their stacked layers as follows

$$\frac{\partial \mathcal{L}}{\partial w^{(1)}} = \underbrace{\frac{\partial \mathcal{L}}{\partial a^{(4)}} \frac{\partial a^{(4)}}{\partial z^{(4)}}}_{\text{From } w^{(4)}} \underbrace{\frac{\partial z^{(4)}}{\partial a^{(3)}} \frac{\partial a^{(3)}}{\partial z^{(3)}}}_{\text{From } w^{(3)}} \underbrace{\frac{\partial z^{(3)}}{\partial a^{(2)}} \frac{\partial a^{(2)}}{\partial z^{(2)}}}_{\text{From } w^{(2)}} \frac{\partial z^{(2)}}{\partial a^{(1)}} \frac{\partial a^{(1)}}{\partial z^{(1)}} \frac{\partial z^{(1)}}{\partial w^{(1)}} \quad (1.20)$$

while the last layer only depends on one set of changes when computing the chain rule

$$\frac{\partial \mathcal{L}}{\partial w^4} = \frac{\partial \mathcal{L}}{\partial a^4} \frac{\partial a^4}{\partial z^4} \frac{\partial z^4}{\partial w^4} \quad (1.21)$$

An opposite problem appears when the value of the gradients increases rapidly in value, which is commonly referred to as the *exploding gradients*

problem. As for the previous issue, it appears when a relevant part of the weights in the network exhibit the same behavior. A common rule of thumb to diagnose the presence of one of these situations is to look at the magnitude of the single weight so that the chances of running into a vanishing gradient problem are higher when $0 < w < 1$ while encountering exploding gradients could happen when $w > 1$.

To understand how the exploding gradients problem works, let denote weights and biases respectively as w_i , biases b_i and the error function is J . Analogously to the previous case, the rate of change in a weight or bias is measured in relation to the error function, by computing the partial derivative. For instance, focusing on the first bias b_1 of the model, such ratio is

$$\frac{\partial J}{\partial b_1} = \frac{\partial J}{\partial a} \frac{\partial a}{\partial z} \frac{\partial z}{\partial b_1} \quad (1.22)$$

We know that the first bias feeds into the first activation a_1 , recalling that $a_1 = \sigma(z_1) = \sigma(w_1 a_0 + b_1)$. If b_1 changes, and we call its variation Δb_1 , the activation a_1 also varies accordingly to $\partial a_1 / \partial b_1 = \partial \sigma(z_1) / \partial b_1$. The change of a_1 is called Δa_1 and which is approximately

$$\Delta a_1 \approx \frac{\partial \sigma(w_1 a_0 + b_1)}{\partial b_1} \Delta b_1 = \sigma'(z_1) \Delta b_1 \quad (1.23)$$

However, variations in a_1 cause variations in $z_2 = w_2 a_1 + b_2$, which is the input of the next layer, and we can go through the same process again to get the change in z_2

$$\Delta z_2 \approx \frac{\partial z_2}{\partial a_1} \Delta a_1 = w_2 \Delta a_1 \quad (1.24)$$

Iterating this process, we get ΔJ as a function of all the changes in the error function relative to every parameter in the network

$$\Delta J \approx \sigma'(z_1) w_2 \sigma'(z_2) \dots \sigma'(z_4) \frac{\partial J}{\partial a_4} \Delta b_1 \quad (1.25)$$

From this, we simply plug into the $\partial J / \partial b_1$ equation and get the final chain rule of partial derivatives with respect to the parameter of the whole network

$$\frac{\partial J}{\partial b_1} = \sigma'(z_1) w_2 \sigma'(z_2) w_3 \sigma'(z_3) w_4 \sigma'(z_4) \frac{\partial J}{\partial a_4} \quad (1.26)$$

If we assume that many weights are greater than one, the computation of the gradients involves the multiplication of large quantities, which leads to a considerable value for the gradients. Using a gradient descent optimizer would update the parameter too far away from the optimal solution.

Exploding gradient issues, like the one presented here, are partially avoided by limiting gradients to some fixed values or imposing a constraint on the norm of the gradient itself (Géron, 2019). However, these heuristics cannot fully control the vanishing gradients problem. For this purpose, different activation functions such as the ReLU have been adopted for training a neural network. The equation for the ReLU activation and its derivative with respect to the input is as follows:

$$\text{ReLU}(x) = \max(0, x) \tag{1.27}$$

$$\text{ReLU}'(x) = \begin{cases} 1 & \text{if } x > 0 \\ 0 & \text{if } x \leq 0 \end{cases} \tag{1.28}$$

On the one hand, this type of nonlinearity fixes the vanishing gradient problem by returning zero gradients instead of minimal values corresponding to negative inputs. On the other hand, it creates another issue, referred to as the *dead ReLU* problem, since when many partial derivatives are lower bounded at zero, the parameter updates still lack information, and the model is not able to learn a proper configuration. More specifically, the neuron of the network which receives a negative input is set to zero, and it is somehow “dead”, because it is not able to provide any additional information when partial derivatives with respect to its connected weights are computed. The dead ReLU problem also introduces sparsity in the weight matrices, increasing the efficiency of time and space complexity, which has been observed that leads to better optimization performances.

One of the variants of the ReLU is called Leaky-ReLU which are defined as

$$\text{LReLU}(x) = \begin{cases} x & \text{if } x > 0 \\ \alpha x & \text{if } x \leq 0 \end{cases} \tag{1.29}$$

$$\text{LReLU}'(x) = \begin{cases} 1 & \text{if } x > 0 \\ \alpha & \text{if } x \leq 0 \end{cases} \tag{1.30}$$

where α is a scalar parameter commonly set between 0.1 to 0.3. The Leaky-ReLU solves the dead ReLU problem, because the values of the gradients are no longer be stuck at zero and at the same time still avoids the vanishing gradient problem.

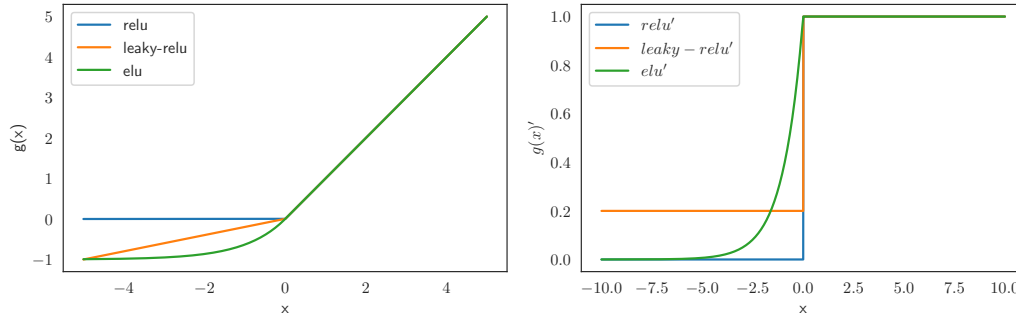


Figure 1.8: ReLU variants activations on the left panel and their first derivatives on the right panel

Another common and effective variant is the Exponential linear unit (ELU) with such form:

$$\text{ELU}(x) = \begin{cases} x & \text{if } x > 0 \\ \alpha(e^x - 1) & \text{if } x < 0 \end{cases} \quad (1.31)$$

$$\text{ELU}'(x) = \begin{cases} 1 & \text{if } x > 0 \\ \text{ELU}(x) + \alpha & \text{if } x \leq 0 \end{cases} \quad (1.32)$$

As for the previous ReLU variant, negative inputs assume values below zero, even though it does not become a linear function when differentiated. Figure 1.8 provides a visualization of ReLU and the variants mentioned above, together with their first derivative function. These kinds of activation functions have been successfully employed in deep learning applications, and they took the place of sigmoid functions as the standard activation for every neural network implementation.

1.11 Summary

This chapter serves to introduce concepts and algorithms that are necessary to deal with any machine learning application. Firstly, we present the three machine learning paradigms to provide a broad perspective of the domains of application of this field of research and what set of problems it can help to solve. We specify the generalization purpose of the machine learning algorithms, in contrast to more classical approaches, and we outlined some background theory for the field. Then we introduce neural networks as one

of the leading mathematical tools for machine learning tasks. The chapter structure facilitates a clear comprehension of the origin and the motivation of those models, the procedure to design them, and the algorithm needed for optimization. We decided to focus only on general machine learning concepts and present just a family of machine learning models, neural networks, because we need the reader to be equipped with instrumental notions to introduce reinforcement learning in the next chapter. Any detailed insights can be retrieved in Friedman et al., 2001 for what concerns classical machine learning and in Goodfellow et al., 2016 regarding deep learning and neural networks.

Chapter 2

Reinforcement learning: theory and background

This work aims to show the role that the reinforcement learning framework can play in the domain of finance. This chapter outlines the background theory and the necessary concepts to understand the origin and functioning of such a machine learning paradigm. At the end of this chapter, the reader should know the connection between the neural networks model family and the reinforcement learning paradigm for applications in several domains, whose finance is one. We complete a broad overview of the methods used in the second part of the work.

”Reinforcement learning, like many topics whose names end with “ing”, such as machine learning and mountaineering, is simultaneously a problem, a class of solution methods that work well on the problem, and the field that studies this problem and its solution methods”(Sutton and Barto, 2018). In these words, the authors of one of the most comprehensive books on reinforcement learning summarize an essential point to grasp in order to understand the subject since the word *reinforcement learning* is equally used to identify both the problem and its solution. In what follows, we will pay particular attention to clearly defining the former and the latter.

As anticipated in the introduction, reinforcement learning tackles the problem of sequential decision making, providing a framework to learn how to behave in a given environment. This setting contrasts with the other machine learning paradigm that we described in the previous chapter. In a supervised learning framework, a fixed dataset is provided, and the selected algorithm search for the desired mapping to represent an available

relationship in the data. In an unsupervised framework, the same dataset is investigated to retrieve its hidden structure and characteristics. Under this perspective, both these frameworks do not resemble the human learning paradigm because they miss the general notion of learning by trial and error. On the contrary, reinforcement learning attempts to mathematically formalize the idea of learning by interaction, which is typical for the human nature. Simple examples are a child who receives feedback and responses from the environment and learns what is safe and dangerous. An adult learns to drive a car by following a set of rules and collecting experiences down the road. Indeed, in many cases learning by interaction is at the core of any human acquired knowledge, from driving a vehicle to using a computer.

The goal of reinforcement learning is broader, even though it shares common traits with the other two machine learning paradigms. For instance, there is still the willingness to find a desired mapping between the situation experienced by the agent and the result of its interaction with the external environment. Reinforcement learning provides a different paradigm in which an algorithm can learn from the gathered experience and is somewhat able to solve the problem of the availability of the data, as long as it is possible to interact with the environment and collect them. The reinforcement learning problem is hence commonly represented by a goal-directed agent that moves in an uncertain environment, and it is particularly suitable to tackle a variety of problems in which the concept of sequential decision-making is crucial.

This paradigm represents one of the core research threads in machine learning and artificial intelligence. Its modern expression is entrenched with other mathematical disciplines such as statistics and optimization. The use of neural networks and reinforcement learning has led to deep reinforcement learning, which contributed to solving the curse of dimensionality problem of the classical control theory.

2.1 Basic notions

As a learning problem, reinforcement learning refers to controlling the behavior of an agent, which is often stochastic, to maximize a scalar value that represents a measure of its long-term objective. The usual interactions in this kind of stochastic system are represented in Figure 2.1: an agent receives information about the current state of the environment, S_t , and it takes action A_t based on that informative set. The environment responds to the action

with a scalar reward signal R_{t+1} and passes through a new state S_{t+1} before the loop starts again. This general scheme holds for each reinforcement learning algorithm, which differs in the details of each scheme component, such as, for instance, the way the performances are measured, the way the state of the environment is expressed, or the type of action performed by the agent.

As anticipated, we always consider the environment stochastic so that there is uncertainty in the state transitions. The first peculiarity of the reinforcement learning problem is the absence of a supervisor, i.e., a ground truth on which the learning problem is based. In contrast, the scalar reward signal plays the most relevant role by providing feedback to the learning agent. This feedback can eventually be delayed over some periods so that an action performed now can be fruitful later in the future. This aspect introduces the second characteristic of reinforcement learning, which is the role of time in the learning process. All the tasks that we aim to solve using reinforcement learning techniques are sequential decision-making tasks, where the training data are not IID but are connected in time. It is essential to understand that each agent's action can affect all the following data that he gathers while interacting with the environment.

Based on the representation of the loop of interactions, we identify the core components of the reinforcement learning problems and give them a clear definition. The first component is a scalar feedback signal referred to as reward, which is helpful to indicate how well the agent is performing at a specific time step in the learning process. Such a measure of performance plays a similar role to the accuracy metric commonly employed in supervised learning in the sense that they measure progress. However, more than that, it can direct and shape the learning process itself. In reinforcement learning, the *reward hypothesis* states that *all of what we mean by goals and purposes can be well thought of as the maximization of the expected value of the cumulative sum of a received scalar signal* (Sutton and Barto, 2018). Even though a recent work of Silver et al., 2021 supports this thesis, it should not be taken as absolute truth, and there is a substantial philosophical debate behind it, claiming that a single scalar is not sufficient to represent multi-modal human behaviors. Such a hypothesis may not attain the goal of general intelligence. However, it is sufficient to express a consistent variety of specific problems in different domains and therefore represents one of the core strengths of the reinforcement learning approach. The formalization and the type of reward function can change dramatically depending on the task that the agents need

2. Reinforcement learning: theory and background

to solve. For certain board games, like Go, the reward is only given at the end of the game, measuring if a winning performance has been achieved. Otherwise, the agent would be encouraged to play well by cumulating its reward, but not necessarily to win the game. In some videogames, which are standard testbeds for reinforcement learning algorithms, the reward depends on the increment or decrement of the score, and the game itself naturally measures it. In other contexts, it is less trivial to formalize the reward concepts, like when the reinforcement learning agent needs to learn how to manage a power station: we reward the agent for producing the necessary power, but we penalize for the waste in excess production or for surpassing the safety thresholds. In a financial application, the reward formalization comes very naturally with the concept of cumulated money over time, which helps shape the problem and already gives the perception of how well the reinforcement learning framework suits financial problems of that sort.

Since the goal is to maximize the cumulative future reward, the so-called exploration-exploitation tradeoff balances the greediness in the behavior apt to maximize the total reward while choosing between the two opportunities. We know that the reward of action may be delayed and affect all future states of the system so that the agent would sacrifice some part of the current reward to achieve a better total reward in the long run. The role of exploration is critical in reinforcement learning, which consists in choosing actions that are not necessarily the most rewarding at the current time, but that could bring a huge benefit in the future. Being the scalar signal is the only feedback that the agents have to learn, its chances are to try different actions and understand where they bring before acting with a whole greedy behavior that exploits the opportunities in the short-term.

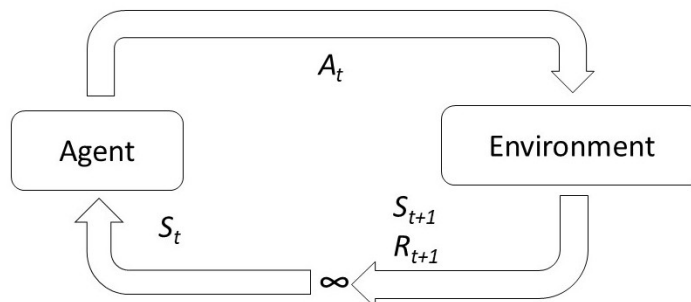


Figure 2.1: The agent–environment loop of interactions.

Following the stylized scheme in Figure 2.1, at each discrete time step t , the agent observes the variables that characterize the current state of the environment S_t and executes an action A_t based on such information. As a response, the environment provides the agent with a scalar reward R_t and the new a state representation S_{t+1} . To underline the agent-environment relationship from the opposite viewpoint, we can also say that at the very same time step the environment is in a current state S_t , it receives an action A_t and outputs a scalar reward R_t and a subsequent state representation S_{t+1} . The relationship between the agent and the environment produces a set of sequential information $S_t, A_t, R_{t+1}, S_{t+1}, A_{t+1}, R_{t+1}, S_{t+2} \dots$ which are not IID. The transition from one state of the environment to the other is given by the model of the environment, which provides what the next state will be, given the current state.

Although the agent's action can be intuitively understood as anything that can bring value to him in the long-term, such as a smart move in a board game or a profitable trade in a financial problem, the state of the environment is somewhat blurred and requires some specifications. In reinforcement learning, we generally refer to the state S_t as the representation of the state of the environment at time t . Therefore, it is the information used by the reinforcement learning agent to decide which action to take to maximize its cumulative reward. The state representation can include any information that is believed to be relevant for the agent to succeed in the defined environment.

The connection between the state of the environment and the agent's action is formalized by the policy of the agent, which reflects its behavior. It is a map from the state representation to the action. A policy π can be either deterministic, $a = \pi(s)$, or stochastic, $\pi(a | s) = \mathbb{P}[A_t = a | S_t = s]$.

Given the existence of delayed feedback, it is important to have a way to measure the goodness of being in a certain state of the environment and hence behaving according to a specific strategy, i.e., a policy. In reinforcement learning, the value function plays such a role, representing an estimate of the future rewards that depends on the agent's policy, so that $v_\pi(s) = \mathbb{E}[R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots | S_t = s]$ represents the expected cumulative discounted sum of rewards based on behaving according to the policy π , given the current state of the information. The next section will provide more details about its use for solving a general reinforcement learning problem.

As we will see in the following section of this chapter, the elements described here are always present in each reinforcement learning algorithm.

Going beyond the high-level definition provided in this section helps in categorizing those algorithms into different families. The introduction of reinforcement learning taxonomy clarifies the choice of a given reinforcement learning algorithm to solve a specific problem. In the next section, we will present the mathematical framework that encloses all these elements and helps make the reinforcement learning problem tractable and even solvable in closed-form under certain conditions.

2.2 Markov Decision Processes

A Markov Decision Process (MDP) is the mathematical formalization of the reinforcement learning problem where the agents aim at learning from interactions. The interaction exists between the agent, which is the decision-maker, and the environment, which includes everything that is not under the direct control of the former. Our focus will be on discrete MDP so that the interaction takes place at discrete time steps $t = 0, 1, 2, 3, \dots$, but the results can be generalized also to continuous time (Bertsekas and Tsitsiklis, 1996). A characteristic of the reinforcement learning problem formalized as a discrete MDP is the possibility to consider the time steps as subsequent possible moments in the decision-making process without necessarily assuming that they are uniformly distributed over time.

At each time step t in a discrete MDP, the agent observes a representation of the current state of the environment, $S_t \in \mathcal{S}$, and selects an action, $A_t \in \mathcal{A}$, on the basis of such information. At the subsequent time step $t + 1$, the agent receives a scalar reward, $R_{t+1} \in \mathcal{R} \subset \mathbb{R}$, together with a new state representation, S_{t+1} . If the MDP is finite, this means that the set of possible states, action and rewards (\mathcal{S} , \mathcal{A} , and \mathcal{R}) is composed by finite elements. Hence, the dynamics of the random variables representing the state and the rewards, S_{t+1} and R_{t+1} , is represented by discrete probability distributions that depends on the previous value of the state and the action that causes the transition from S_t to S_{t+1}

$$p(s', r | s, a) \doteq \Pr \{S_{t+1} = s', R_{t+1} = r | S_t = s, A_t = a\}. \quad (2.1)$$

This definition holds from every finite elements of the three sets (\mathcal{S} , \mathcal{A} , and \mathcal{R}). It worth noting that the probability distribution is specified for every

possible choice of the pair (s, a) so that

$$\sum_{s' \in \mathcal{S}} \sum_{r \in \mathcal{R}} p(s', r | s, a) = 1, \text{ for all } s \in \mathcal{S}, a \in \mathcal{A} \quad (2.2)$$

The function $p : \mathcal{S} \times \mathcal{R} \times \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$ determines the dynamics of the MDP, so that the probability of the occurrence for S_t and R_t depends only on the information available at the preceding time step, hence S_{t-1} and A_{t-1} . In a MDP, the state encloses all the relevant information about the past interactions, hence it respects the Markov property.

From the very general Eq. 2.1, one can compute the state-transition probabilities $p : \mathcal{S} \times \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$ simply by summing over all possible rewards,

$$p(s' | s, a) \doteq \Pr \{S_{t+1} = s' | S_t = s, A_t = a\} = \sum_{r \in \mathcal{R}} p(s', r | s, a), \quad (2.3)$$

or the expected rewards $r : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ for a specific pair (s, a)

$$r(s, a) \doteq \mathbb{E} [R_{t+1} | S_t = s, A_t = a] = \sum_{r \in \mathcal{R}} r \sum_{s' \in \mathcal{S}} p(s', r | s, a). \quad (2.4)$$

The scalar feedback provided to the agent by the environment allows to describe a measure of the total performance of the agent in the learning problem as an infinite discounted sum of rewards:

$$G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots = \sum_{k=0}^{\infty} \gamma^k R_{t+k-1}, \quad (2.5)$$

so that the value of receiving the reward R after $k + 1$ steps is $\gamma^k R$. The role of the discount factor γ is important for evaluating instantaneous reward with respect to future rewards. A value close to zero resembles a short-term attitude and pays more attention to immediate rewards, while a value close to one enhances more long-term rewards. The role of the discount is also to avoid infinite return values in infinite Markov processes and helping to represent the human behavior preference for immediate reward. In our financial and economic context, it is also directly involved in representing the time value of money.

In a (stochastic) MDP, a policy is considered to be the distribution of the actions given the states:

$$\pi(a | s) \doteq \Pr \{A_t = a | S_t = s\}, \quad (2.6)$$

so that, if the agent follows a policy π and is experiencing s as the current state, the action a will be chosen with probability $\pi(a | s)$. The ordering of policies is guaranteed by specifying a total return function as in equation 2.5, which represents the objective function of the theoretical reinforcement learning problem. The agent searches for the policy that maximizes the expectation of total return

$$\max_{\pi} \mathbb{E}[G_t]. \quad (2.7)$$

Value functions help in searching for the optimal policy since they represent how good it is to be in a given state in terms of the expected future cumulative sum of rewards and help estimate the future amount G_t . They are usually defined as related to a policy that can be modified during the learning process. The state-value function $v(s)$ of an MDP represents the expected future return of being in a given state and following a policy π onwards. An alternative to the state-value function is the action-value function, representing the expected return function, assuming to start in the state s , take action a and follow some fixed policy π onwards. Both concepts are mathematically expressed as follows:

$$v_{\pi}(s) := \mathbb{E}_{\pi}[G_t | S_t = s] \quad (2.8)$$

$$q_{\pi}(s, a) := \mathbb{E}_{\pi}[G_t | S_t = s, A_t = a]. \quad (2.9)$$

For any policy π and any state s , the following consistency condition holds between the value of the current state and the value of its possible successor states:

$$\begin{aligned} v_{\pi}(s) &= \mathbb{E}_{\pi}[G_t | S_t = s] = \\ &= \mathbb{E}_{\pi}[R_{t+1} + \gamma G_{t+1} | S_t = s] = \\ &= \sum_{a \in \mathcal{A}} \pi(a | s) \left(r(s, a) + \gamma \sum_{s' \in \mathcal{S}} p(s' | s, a) \mathbb{E}_{\pi}[G_{t+1} | S_{t+1} = s'] \right) = \\ &= \sum_{a \in \mathcal{A}} \pi(a | s) \left(r(s, a) + \gamma \sum_{s' \in \mathcal{S}} p(s' | s, a) v_{\pi}(s') \right), \end{aligned} \quad (2.10)$$

which is called Bellman equation from the work of Bellman, 1966. The same recursion can also be written for the action-value function and the two value functions are tied by the following relationship

$$v_{\pi}(s) = \sum_{a \in \mathcal{A}} \pi(a | s) q_{\pi}(s, a) \quad (2.11)$$

2. Reinforcement learning: theory and background

We know that the state-value function allows for ranking the policies, so that $\pi > \pi'$ if $v_\pi(s) > v_{\pi'}(s) \forall s \in \mathcal{S}$. Although an optimal policy is defined to be at least as good as any other policy and does not need to be unique if two policies are optimal, they certainly share the same optimal value functions:

$$v_*(s) = \max_{\pi} v_{\pi}(s) \tag{2.12}$$

$$q_*(s, a) = \max_{\pi} q_{\pi}(s, a) \tag{2.13}$$

Noting that $v_*(s) = \max_a q_*(s, a)$, we can conclude that the optimal action-value function is a more general concept than the optimal state-value function. If we knew the $q_*(s, a)$, the optimal policy would be guaranteed by choice of a that maximizes q_* for the current state s . This policy is called greedy, and the reinforcement learning problem is usually reduced to find q_* in an iterative manner.

However, the greedy policy is a particular case of a deterministic policy. In a more general case, one can express a stochastic policy that specifies certain probabilities $\pi(a | s)$ of taking action a given the current state s . When there are several actions at which the maximum action-value function is achieved, a stochastic policy expresses the agent's behavior by assigning probabilities to each maximizing action. Any weighting scheme is suitable as long as the suboptimal action, i.e., those that do not maximize the value function, are given zero probabilities. Hereafter, we can alternatively use deterministic or stochastic policies. The former is just a general case in which the weighting scheme is chosen so that one maximizing action is given full probability and the other is given zero. Therefore, all the results presented in this section are extendable to the stochastic case.

The Bellman optimality equations are both satisfied by the optimal value functions as follows:

$$v_*(s) = \max_a \left(r(s, a) + \gamma \sum_{s' \in \mathcal{S}} p(s' | s, a) v_*(s') \right) \tag{2.14}$$

$$q_*(s, a) = r(s, a) + \gamma \sum_{s' \in \mathcal{S}} p(s' | s, a) \max_{a'} q_*(s', a'), \tag{2.15}$$

where the recursion is not averaged over all the action in the space \mathcal{A} as in Eq. 2.10, but instead, it is expressed in term of the action that maximizes

it. The optimality equation for $q_*(s, a)$ is easily obtained from Eq. 2.11. The majority of reinforcement learning algorithms associate a quantity $q_*(s, a)$ to the state-action pair (s, a) . The main problem is that we are not able to calculate $q_*(s, a)$ since, at any current time, the optimal action-value function is unknown. Therefore, we need to resort to an approximation of the target value, using our current best guess of the function q_* . In general, any algorithm belonging to the reinforcement learning family consists of the following iterative steps:

1. generate samples (s, a, r, s') by running a policy
2. estimate the return of the sampled trajectory
3. improve the policy.

Depending on the way those steps are approached, different classes of reinforcement learning algorithms are defined.

The solution to the reinforcement learning problem formalized by an MDP is obtained by solving the Bellman optimality equation, which is nonlinear by nature and seldom has a closed-form solution. In the next section, we will present different methods to solve the general problem described above, along with their advantages and disadvantages.

2.3 Tabular Reinforcement Learning

This section enters into the details of tabular reinforcement learning algorithms, which represent the first building block to understanding more sophisticated families of algorithms. For this purpose, it is helpful to describe the dynamic programming approach and its differences with respect to reinforcement learning to fully appreciate the flexibility of the former class of methods for sequential decision-making problems.

2.3.1 Dynamic Programming

Classical reinforcement learning, also referred to as tabular, is a discipline related to dynamic programming, which encloses a set of methods for stochastic control problems. Examining the former name, we understand that there is the attempt to solve a sequential problem, hence *dynamic*, by optimizing

a *program*, which in our case is the way of behaving of the agent. Understanding dynamic programming helps in grasping the basic algorithm of reinforcement learning, which is at the core of modern and more complex ones. The primary reference for this set of algorithms is (Bellman, 1966), although the same material is also presented in (Sutton and Barto, 2018).

Each dynamic programming problem has a clear and defined structure. There is an optimal solution to be found by decomposing the main problem into subproblems that recur many times. An MDP provides a standard framework to represent this decomposition in sub-problems.

The main characteristic of the dynamic programming approach is to assume the full knowledge of the MDP, i.e., the probability distribution in Eq. 2.1 at each discrete time, which gives the possibility to compute optimal policies. Reinforcement learning methods often share some characteristics with dynamic programming techniques, even though for the former, the dynamics of the state transition are generally unknown and need to be discovered by interacting with the environment. In general, since dynamic programming exploits the knowledge of the MDP's dynamics, it is more computationally expensive than a reinforcement learning approach because it does not resort to approximation.

Dynamic programming uses value functions to structure the search for good policies since it is relatively easy to obtain optimal policies once obtained the value function by solving the Bellman equation as in Eq. 2.14 or 2.15. The problem is split in two parts, of which the first is called *policy evaluation* and regards the computation of the state-value function v_π for a given policy π , to assess the goodness of the current state. When the model of the environment is known, as it is the case of many dynamic programming applications, the Eq. 2.14 consists in solving a system of $|N|$ linear equations in $|N|$ variables, which are those composing the state space representation. This system of equation is solved by iteration, so that an initial value v_0 is chosen, and then the subsequent approximation are obtained through the following update rule for all $s \in \mathcal{S}$

$$v_\pi^{k+1}(s) = \sum_{a \in \mathcal{A}} \pi(a | s) \left(r(s, a) + \gamma \sum_{s' \in \mathcal{S}} p(s' | s, a) v_\pi^k(s') \right), \quad (2.16)$$

that is the Bellmann equation iteratively updated every k step. Repeated updates of the value function are referred to as *iterative policy evaluation*,

whose computational costs grows with the number of possible states in which the environment can make its transition.

Evaluating a given policy, i.e., computing its associated value function, also allows improving it by searching for a better one. One can wonder if for a given state s at the time t , it would be better to choose a different action with respect to what the policy suggests, $a \neq \pi(s)$, and then follow the policy π itself thereafter. The action value function in Eq. 2.9 evaluates this way of behaving as

$$q_{\pi}(s, a) = r(s, a) + \gamma \sum_{s' \in \mathcal{S}} p(s' | s, a) v_{\pi}(s'), \quad (2.17)$$

which provides a measure of how much is convenient to change the current policy by taking a different action. If $q_{\pi}(s, a) > v_{\pi}(s)$, then it would be better to take the action a in the current state s and follows the given policy thereafter. On the other hand, the opposite case discourages changing the policy for the current state.

The general case for the improvement of a dynamic programming policy is formalized by the *policy improvement theorem* which states that, given a pair of deterministic policies π and π' such that for all $s \in \mathcal{S}$ holds,

$$q_{\pi}(s, \pi'(s)) \geq v_{\pi}(s), \quad (2.18)$$

then the policy π' must be as good as, or better than, π . Being better in reinforcement learning terminology means that the policy π' must obtain at least an equal amount of expected return from all states $s \in \mathcal{S}$, i.e.

$$v_{\pi'}(s) \geq v_{\pi}(s) \quad \forall s \in \mathcal{S} \quad (2.19)$$

If $q_{\pi}(s, a) > v_{\pi}(s)$, then the modified policy is indeed better than π , since the two initial policies differ in the way of behaving only at the state s , while are entirely equal for all the other possible states in \mathcal{S} . The proof of the policy improvement theorem goes beyond the scope of this section and can be easily retrieved in Sutton and Barto, 2018. Notwithstanding the use of deterministic policy to introduce the policy improvement theorem, all results can also be extended to stochastic policies $\pi(a | s)$ that specify probabilities of taking actions.

The alternating process of evaluating a policy and improving it by obtaining a better one produces a series of monotonically improving policies and

value functions. Thanks to the policy improvement theorem, strict improvement is guaranteed at any step. The finiteness of the MDP under which dynamic programming works allows convergence to optimality in a finite number of iterations. The serious drawback of such a procedure is the need for an entire sweep, i.e., evaluating the policy for all the possible states in the MDP. However, it has been shown that the policy evaluation step of policy iteration can be truncated without losing the convergence guarantees of the policy iteration algorithm. The value iteration algorithm stops the policy evaluation exactly after one complete update of each state, which takes the following form:

$$v_{\pi}^{k+1}(s) = \max_a \left(r(s, a) + \gamma \sum_{s' \in \mathcal{S}} p(s' | s, a) v_{\pi}^k(s') \right), \quad (2.20)$$

for all $s \in \mathcal{S}$. The algorithm stops once the value function changes by only a small amount in a sweep. Value iteration allows evaluating and updating the same policy in one sweep over the possible state space, allowing to do fewer computations and achieving faster convergence. The previous algorithm converges to an optimal policy if used to solve a discounted finite MDP.

The main drawback of dynamic programming that has been discussed in the literature is the so-called *curse of dimensionality* that makes this family of methods impractical for high-dimensional problems. The time required for a dynamic programming method to solve a policy is polynomial in the dimensionality of the state and the action spaces. Therefore, when the number of possible states grows, dynamic programming has hard time to solve an MDP. Even though high-performance calculus machines partially solve the timing issues, the curse of dimensionality represents one of the main reasons for the surge of reinforcement learning techniques.

2.3.2 Model-free Reinforcement Learning

Contrary to the need for dynamic programming to fully know the dynamics of the environment, model-free reinforcement learning encloses a set of techniques and algorithms to solve MDPs in the case of partial or absence of information about such dynamics. The approach to solving the MDP is still articulated in two steps, i.e., prediction, estimating the value function, and control, to optimize the value function and eventually improve the policy.

The difference with respect to dynamic programming refers to how these two steps are accomplished.

The first set of model-free reinforcement learning algorithms we want to present is Monte-Carlo reinforcement learning, which learns from an episode of experiences. An episode is a set of consecutive time steps in which the agent and the environment interact up to a final time step. After that, the system is restarted, and a new episode begins. Such a definition implies that these techniques can be applied only to episodic MDPs, representing a sequential problem that needs to end, as it happens in almost any game.

After collecting several episodes of experience using a policy π , the goal of a Monte-Carlo algorithm is again to learn the state-value function v_π . The experience obtained from an episode appears like a sequence of signals that goes back and forth between the agent and the environment $S_1, A_1, R_2, \dots, S_k \sim \pi$. Recalling the definition of total return as the total amount of discounted reward, which in this case is a finite sum due to the episodic MDP, we obtain $G_t = R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{T-1} R_T$. From Eq. 2.8, we have defined the state-value function as the expected return under the current policy π . Monte-Carlo reinforcement learning aims to use the empirical mean return over the episode instead of the expected one.

When one wants to evaluate the state s , the algorithm keeps a counter $N(s)$ of the number of times that such state has been visited during an episode. The value of the state is then estimated as the empirical mean return $V(s) = S(s)/N(s)$, where $S(s)$ is the aggregated return obtained after every visit of the state s . Then we know that as long as the agent visits such state for a sufficient amount of time, the solution of the MDP is going to converge to the searched value function, i.e., $V(s) \rightarrow v_\pi(s)$ as $N(s) \rightarrow \infty$. This approach is called every-visit because it averages the return after each state visit in the set of all episodes. As an alternative, the first-visit approach averages the return only after the first appearance of such a state in each episode. Usually the value function is incrementally updated after each episode $S_1, A_1, R_2, \dots, S_T$, so that for each state S_t with return G_t we have

$$V(S_t) \leftarrow V(S_t) + \frac{1}{N(S_t)} (G_t - V(S_t)), \quad (2.21)$$

where $N(S_t)$ is still the counter of the visit for the state S_t . The update rule can be modified to adapt to non-stationary environment, where it could be convenient to gradually forget past episode by giving more weights to the

recent one, so that

$$V(S_t) \leftarrow V(S_t) + \alpha(G_t - V(S_t)). \quad (2.22)$$

Temporal-difference learning is a different model-free approach, which still does not require knowledge of the transition function associated with the MDP that one needs to solve. However, it can learn from incomplete episodes. Temporal-difference reinforcement learning adopts bootstrapping techniques that update the value function, which is a guess, towards another guess estimated according to what has been seen so far during the agent-environment interaction. The aim is to learn online, hence faster, without waiting until the end of the episode and collecting a full set of sequences up to a termination state. This allows for solving MDP, which is not expressed episodic but eventually has an infinite time horizon. Recalling Eq. 2.22, temporal-difference reinforcement learning updates the value function $V(S_t)$ toward an estimation of the return $R_{t+1} + \gamma V(S_{t+1})$ called TD target, so that

$$V(S_t) \leftarrow V(S_t) + \alpha(R_{t+1} + \gamma V(S_{t+1}) - V(S_t)). \quad (2.23)$$

where $\delta_t = R_{t+1} + \gamma V(S_{t+1}) - V(S_t)$ is the TD error.

When it comes to comparing the two model-free approaches presented in this section, there is no prevalent approach in terms of performance. All that matters is the proper selection of the approach for the MDP that we are trying to solve. In order to recap, the two approaches are not competing but complementary since they provide algorithmic methods for solving both finite and infinite horizon MDPs. The temporal-difference approach can learn after every step, without going too far up to the end of the episode, using incomplete sequences of experience. In this sense, it works for both episodic and continuing environments, while Monte-Carlo reinforcement learning can only solve those that are episodic.

Looking at the way the value function is updated, there is an intrinsic tradeoff between bias and variance since the total return $G_t = R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{T-1} R_T$ is an unbiased estimate of $v_\pi(S_t)$, while the TD target $R_{t+1} + \gamma V(S_{t+1})$ is a biased estimate of $v_\pi(S_t)$, but it has lower variance than the former. This holds because the total return depends on the full set of state-action-reward transitions, while the TD target depends only on one of them. Monte-Carlo has high variance but low bias and usually has good convergence properties. The main drawback is collecting an entire episode of sequences before updating the value function. On the other hand, the

2. Reinforcement learning: theory and background

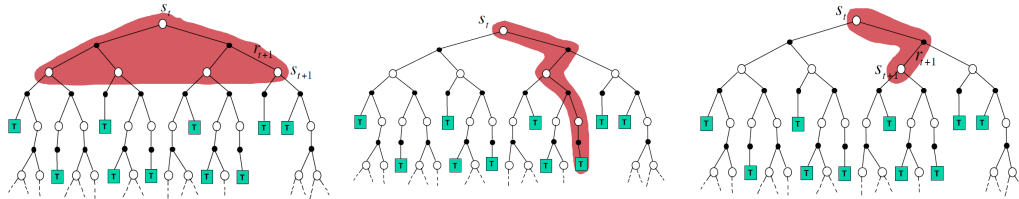


Figure 2.2: From the left to the right, backup diagrams for dynamic programming, Monte-Carlo and temporal-difference approach. These diagrams represent the update operations, also referred to as backups, that are carried out by these classes of algorithms. Source: David Silver lecture on reinforcement learning available [here](#)

temporal difference is more efficient but has fewer convergence guarantees. In this regard, temporal difference exploits the Markov property. It works well in Markov environments where all the information is enclosed in the system's current state. In contrast, Monte-Carlo does not use that property and requires the full stream of experience to update the value function.

Figure 2.2 helps to clarify these distinctions also compared to the dynamic programming approach. The dynamic programming update operation, commonly known as backup (see the left part of the Figure), underlines that to use the update rule $V(S_t) \leftarrow \mathbb{E}_\pi [R_{t+1} + \gamma V(S_{t+1})]$, one needs to consider all the possible successive state starting from S_t and compute an expectation, which is possible due to the given knowledge of all the transition probabilities in the environment. A Monte-Carlo update, $V(S_t) \leftarrow V(S_t) + \alpha (G_t - V(S_t))$, needs the collection of the entire path of the episode up to a terminal state denoted with a T in the figure. A temporal-difference approach requires a truncated version of the same path, which does not go beyond the first successive step, as highlighted on the right end side of the Figure, and it is commonly referred to as TD(0).

There are available options between the two model-free approaches that move the tradeoff between bias and variance and allow using temporal difference up to several steps in the future. Algorithms like eligibility traces or n-step temporal difference are treated in more detail in Sutton and Barto, 2018, and we omit their details in this work since they are not primarily used.

Indeed, the temporal difference algorithm with a single step, also called TD(0), has been shown to work well. It is simple to implement, representing the basis of most of the complex deep reinforcement learning algorithms that

we will describe in the next section.

In order to conclude our exposition of the tabular reinforcement learning method, we want to present Q-learning, which is one of the most known algorithms, originating from a seminal work of (Watkins and Dayan, 1992). After introducing the prediction approach in a model-free context, we will use Q-learning as a representative example of model-free control.

Q-learning

Q-learning is an off-policy learning control method, which means learning a policy π using experience gathered from the interaction with the environment through a different policy. This contrasts with the so-called on-policy learning, where the policy used to sample experience is the same one that wants to optimize. The equivalent algorithm of Q-learning in an on-policy setting is called SARSA, whose name precisely refers to the sequence that needs to be collected to perform an update. In this section, we just outline the off-policy approach, while the on-policy can be retrieved in Sutton and Barto, 2018; Szepesvári, 2010.

Off-policy learning consists in evaluating a policy $\pi(a|s)$, that we want to learn, by computing either the value functions $v_\pi(s)$ or $q_\pi(s, a)$, while in the meantime, following a different policy $\mu(a|s)$. The off-policy concept is important because it allows learning situations where we can observe different behaviors than strictly the one that the agent is taking, and it is also more efficient in the use of collected data. Off-policy learning offers a way to circumvent the known issue of the exploration and exploitation tradeoff in reinforcement learning by allowing to learn optimal policies while following policies that tend to explore the action space.

Q-learning consider off-policy learning of action-values $Q(s, a)$ by choosing an action A_t according to a behavior policy $\mu(A_t | S_t)$ and updating towards $Q(s, a)$ computed with a greedy target policy. This means picking the action that maximizes the action-value function given the current state so that the iterative update of the learning algorithm is

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha \left(R_{t+1} + \gamma \max_{a'} Q(S_{t+1}, a') - Q(S_t, A_t) \right). \quad (2.24)$$

In this setting the agent can behave according to a stochastic policy $\mu(A_t | S_t)$ to learn a deterministic policy $\pi(S_t) = \arg \max_{a'} Q(S_{t+1}, a')$. The balance between exploration and exploitation is usually guaranteed by a behavior

policy that is ϵ -greedy with respect to $Q(s, a)$, representing a simple but effective idea for ensuring constant exploration in this setting. It allows trying all the possible actions in the space \mathcal{A} with nonzero probability by choosing the greedy action with probability $1 - \epsilon$ or picking an action at random with probability ϵ . In a more complex algorithm, the parameter ϵ that controls exploration can be tuned and decreased as the learning process progresses, becoming more greedy towards the end of the update process.

According to the greedy target policy, the Q-learning target takes the form of:

$$\begin{aligned} & R_{t+1} + \gamma Q(S_{t+1}, A') \\ &= R_{t+1} + \gamma Q\left(S_{t+1}, \underset{a'}{\operatorname{arg\,max}} Q(S_{t+1}, a')\right) \\ &= R_{t+1} + \max_{a'} \gamma Q(S_{t+1}, a') \end{aligned} \tag{2.25}$$

and therefore the iterative update is

$$Q(S, A) \leftarrow Q(S, A) + \alpha \left(R + \gamma \max_{a'} Q(S', a') - Q(S, A) \right). \tag{2.26}$$

Looking at the Eq. 2.26, Q-learning is a TD(0) algorithm so that it is sufficient to collect one step of experience to perform an update, which makes the use of such an algorithm very appealing for its efficiency and ease of use. A detailed outline of the algorithmic routine is provided below. Q-learning is also backed up by some theoretical guarantee around the convergence to the optimal solution, as is shown in the seminal paper of Watkins and Dayan, 1992.

Algorithm 1 Q-learning algorithm pseudocode

```
 $Q(s, a), \forall s \in \mathcal{S}, a \in \mathcal{A}$  initialized arbitrarily and  
 $Q(\text{terminalstate}, \cdot) = 0$   
for each episode do  
  Initialize  $S$   
  for each step of episode do  
    Choose  $A$  given  $S$  using policy  $Q(s, a)$  e.g. epsilon-greedy  
    Take action  $A$ , observe  $R, S'$   
     $Q(S, A) \leftarrow Q(S, A) + \alpha (R + \gamma \max_{a'} Q(S', a') - Q(S, A))$   
     $S \leftarrow S'$   
  until  $S$  is terminal  
  end for  
end for
```

Nonetheless, the theoretical convergence guarantee, Q-learning is a tabular reinforcement learning method, which means that it only works when the state space \mathcal{S} and the action space \mathcal{A} are finite-dimensional. Every tabular reinforcement learning algorithm expresses its value function as a table that has an entry for each possible state-action pair (s, a) . Hence, these methods are suitable when both the state and the action space can be discretized in a finite number of elements. However, even a fine-grained discretization may cause a loss of information in problems whose state variables or actions are naturally expressed as real values. A similar issue verifies when the state variables and the action can be expressed as discrete variables, but the number of the possible state-action pair is so large that it is computationally expensive to maintain a huge table with all such entries. To this end, tabular methods still share the curse of dimensionality drawback of the dynamic programming approach. Therefore in the next section, we introduce the concept of function approximation for reinforcement learning, which is crucial to overcoming the issue mentioned above.

2.4 Approximate Reinforcement Learning

Tabular model-free approaches, such as Q-learning, are scaled up by using function approximation to solve high-dimensional problems, which otherwise would require a huge table to be represented. Some examples of intractable

2. Reinforcement learning: theory and background

problems with tabular methods are the game of Go, which has 10^{170} possible states or even maneuvering a truck, which is a problem in continuous state space.

The straightforward approach is to approximate the previously defined value functions using a convenient parametrization to overcome the curse of dimensionality

$$\hat{V}(s, \mathbf{w}) \approx v_\pi(s) \tag{2.27}$$

$$\hat{q}(s, a, \mathbf{w}) \approx q_\pi(s, a), \tag{2.28}$$

instead of relying on a huge table of values, which also cause memory storage problems when it comes to putting it into practice. The principle of function approximation in reinforcement learning is based on generalization over previous unseen states, without the need to have an exact value for every state-action pair. Algorithmic methods such as Monte-Carlo or TD are still useful when the value functions are properly approximated, as we will see in the next subsections.

Among the several possible function approximators, we are going to use neural network models described in Section 1.4, which have been proven to be powerful and effective in solving reinforcement learning problems. In the literature, one can retrieve different approaches, such as linear approximators (Baird, 1995; Melo et al., 2008; Zhou et al., 2020), Fourier bases (Konidaris et al., 2011), regression trees (Uther and Veloso, 1998; Zhou et al., 2020) and nearest neighbors (de Lope, Maravall, et al., 2011; Shah and Xie, 2018). Neural networks, differently from simpler linear approximators, are differentiable function approximators that are usually in competition. However, neural networks have the superior capability in solving high-dimensional sequential decision-making problems, so they contribute to the rise of a relatively recent field of research called deep reinforcement learning, from deep neural networks to solve complex MDPs.

Before introducing the deep reinforcement learning models that we will use in the second part of our work, it is worth giving a broad perspective on how an approximate reinforcement learning algorithm is updated, either using a linear or a nonlinear type of approximation. Since we are referring to functions differentiable in the parameter vector \mathbf{w} , gradient descent methods as described in Section 1.9 are of practical use.

Let denote $J(\mathbf{w})$ a differentiable function of parameter vector \mathbf{w} , then its gradient with respect to such parameters is $\nabla_{\mathbf{w}}J(\mathbf{w})$. To find the local

minimum of the objective function, the parameters are updated towards the negative direction of the gradient, according to the following update function

$$\Delta \mathbf{w} = -\frac{1}{2}\alpha \nabla_{\mathbf{w}} J(\mathbf{w}) \quad (2.29)$$

being α a step-size parameter.

When the objective function $J(\mathbf{w})$ is the mean-square error between the approximate value function $\hat{v}(s, \mathbf{w})$ and the true value function $v_{\pi}(s)$, i.e.

$$J(\mathbf{w}) = \mathbb{E}_{\pi} [(v_{\pi}(S) - \hat{v}(S, \mathbf{w}))^2], \quad (2.30)$$

gradient descent allows then to find a local minimum by computing the following gradient

$$\begin{aligned} \Delta \mathbf{w} &= -\frac{1}{2}\alpha \nabla_{\mathbf{w}} J(\mathbf{w}) \\ &= \alpha \mathbb{E}_{\pi} [(v_{\pi}(S) - \hat{v}(S, \mathbf{w})) \nabla_{\mathbf{w}} \hat{v}(S, \mathbf{w})]. \end{aligned} \quad (2.31)$$

The gradient descent update is known to be inefficient because it requires to average over the whole set of collected states in order to update the parameter of the value function according to Eq. 2.31. This type of update is referred to as a full gradient descent update. A sophisticated reinforcement learning algorithm uses variants of the stochastic gradient descent that does not require to average over all states and at the same time has been proven to be efficient. On the extreme opposite in terms of approach, there is the stochastic gradient descent update which samples a single state transition from the collected experience and performs the following update

$$\Delta \mathbf{w} = \alpha (v_{\pi}(S) - \hat{v}(S, \mathbf{w})) \nabla_{\mathbf{w}} \hat{v}(S, \mathbf{w}). \quad (2.32)$$

Between the two approaches the batch gradient descent samples a set \mathcal{D} of transition of a given size and perform a gradient update by averaging those samples as

$$\Delta \mathbf{w} = \alpha \mathbb{E}_{\mathcal{D}} [(v_{\pi}(S) - \hat{v}(S, \mathbf{w})) \nabla_{\mathbf{w}} \hat{v}(S, \mathbf{w})]. \quad (2.33)$$

Almost all the commonly used deep reinforcement learning algorithms use the batch version of the gradient descent, which is more computationally efficient than a full update and leads to less oscillation and noise with respect to the stochastic update because it still averages the gradients over some samples rather than a single value. Despite the efficiency, a noisy step of a stochastic

update can be beneficial to avoid being locked in the local optima of the objective function. Therefore, the size of the batch for the gradient descent update is a crucial hyperparameter to consider in reinforcement learning, as much as it is for supervised and unsupervised learning applications.

When using function approximation, it is still easy to write the temporal-difference update rule for the approximated state-value function. Let the TD target be $R_{t+1} + \gamma \hat{v}(S_{t+1}, \mathbf{w})$, which is a biased sample of the true value $v_\pi(S_t)$, then the update rule takes the following form:

$$\Delta \mathbf{w} = \alpha (R + \gamma \hat{v}(S', \mathbf{w}) - \hat{v}(S, \mathbf{w})) \nabla_{\mathbf{w}} \hat{v}(S, \mathbf{w}). \quad (2.34)$$

A similar approach can be taken to approximate the action-value function $\hat{q}(S, A, \mathbf{w}) \approx q_\pi(S, A)$, so that a squared error measure between the approximate action-value function $\hat{q}(S, A, \mathbf{w})$ and true action-value function $q_\pi(S, A)$ is minimized

$$J(\mathbf{w}) = \mathbb{E}_\pi [(q_\pi(S, A) - \hat{q}(S, A, \mathbf{w}))^2], \quad (2.35)$$

where the true value function is usually replaced by a biased sample as above. The minimization is carried out via a stochastic gradient descent approach, so that

$$\begin{aligned} -\frac{1}{2} \nabla_{\mathbf{w}} J(\mathbf{w}) &= (q_\pi(S, A) - \hat{q}(S, A, \mathbf{w})) \nabla_{\mathbf{w}} \hat{q}(S, A, \mathbf{w}) \\ \Delta \mathbf{w} &= \alpha (q_\pi(S, A) - \hat{q}(S, A, \mathbf{w})) \nabla_{\mathbf{w}} \hat{q}(S, A, \mathbf{w}). \end{aligned} \quad (2.36)$$

The main drawback of the function approximation with the temporal difference approach is that the gradient updates do not follow the gradient of an existing function since we update a guess towards another guess using a bootstrapping approach. The update of the approximated $\hat{q}(S, A, \mathbf{w})$ points towards a biased sample of the true $q_\pi(S, A)$ which is obtained in terms of the same parametric approximator, i.e., a neural network, as it happens for the state-value function in Eq. 2.34. This update rule implies that the target moves at each step since the parameters \mathbf{w} express both the current value and the target at the same time. The consequence is that there is no convergence guarantee that the biased sample of the true value leads the update of the value function towards the correct value. This explains the divergence issues that function approximation with reinforcement learning is known to have, as summarized in the deadly triad of reinforcement learning (Sutton and Barto, 2018).

2. Reinforcement learning: theory and background

The set of algorithms shown so far in this section approximates the value function and generates a policy accordingly, for example, adopting an ϵ -greedy behavior for selecting actions.

A more direct approach is to parametrize the policy $\pi_\theta(s, a) = \mathbb{P}[a \mid s, \theta]$, which has beneficial effects in the case of MDPs with continuous actions spaces. The policy $\pi_\theta(s, a)$ can be represented in different ways, as long as it is differentiable with respect to its parameters θ . This implies that $\nabla_\theta \pi_\theta(a, s)$ exists for all the states and actions respectively in \mathcal{S} and \mathcal{A} . Another important aspect to consider when choosing the parametric form of the policy is that it is desirable for exploration purposes that the policy never becomes fully deterministic. Given these two properties, the parametric form of the policy is generally chosen according to the MDP that the policy itself is trying to solve. When the MDP has a discrete action space \mathcal{A} that is not too large, a common form is to parametrize the policy with a softmax distribution.

$$\pi_\theta(a, s) = \frac{e^{\phi(s, a, \theta)}}{\sum_{a' \in \mathcal{A}} e^{\phi(s, a', \theta)}}, \quad (2.37)$$

where $\phi(s, a, \theta) \in \mathbb{R}$ are numerical preferences for each state-action pair (s, a) . In this way, the actions with the highest preferences in each state are associated with higher probability. The numerical preferences can be obtained through a neural network, that computes the softmax function in the last layer and directly output the probabilities associated to the actions. On the contrary, when the space \mathcal{A} of the MDP is continuous, the policy is commonly parametrized as a Gaussian, so that $a \sim \mathcal{N}(\phi(s, a, \theta), \sigma^2)$, where $\phi(s, a, \theta)$ is the single output of a neural network that learns the mean of the Gaussian policy. The variance can be fixed or learned together with the mean.

From the introduction provided in this section, we can characterize two different families of algorithms that exploit the function approximation methodology within the reinforcement learning framework. We will refer to these algorithms as deep reinforcement learning since they use neural networks as a function approximation tool. These algorithms allow overcoming the curse of dimensionality problems shared by dynamic programming and tabular reinforcement learning algorithms in Section 2.3 and are largely used because they allow for solving and generalizing over problems in continuous state and action spaces. We introduce one representative algorithm for both sets of methods in deep reinforcement learning, according to the common taxonomy used in the literature. However, several algorithms compose each group,

and each of them can be either a slightly variant of another algorithm or can be very specific for the solution of a particular class of problems. In general, deep reinforcement learning methods are divided into two distinct groups

1. *value-based methods*, which solves an MDP by approximating a value function through the use of neural networks. As presented above, the function to be approximated can be either a state-value function or an action-value function. Once the value function is approximated, the optimal policy for the given problem is inferred from it.
2. *policy-based methods*, which solves the MDP by directly approximating the policy function without going through a value function.

Even though the policy-based approach could appear more natural in learning a good way to behave in a given environment, value functions are still beneficial. Several successful methods in reinforcement learning are hybrid so that they approximate both a value function and a policy function. This hybrid class is referred to as an *actor-critic*, as we are going to see in the following sections. To summarize, recalling the no free lunch theorem in machine learning, there is no optimal choice for the algorithm in reinforcement learning. The choice depends on the type of problem and the setting required by the MDP. Therefore, we will present both sides of the landscape to give a broad perspective of the possibilities provided by the deep reinforcement learning framework.

2.4.1 Value-based methods

The use of a neural network to approximate a function that depends on the current state of the environment allows for a continuous state representation so that $s \in \mathbb{R}^N$, where N represents the number of state-space variables. Arguably the most known deep reinforcement learning algorithm in the value-based class has been proposed by Mnih et al., 2015, which obtained ground-breaking results in learning how to play arcade video games at a superhuman level. It also set the basis for the subsequent development of new algorithms in deep reinforcement learning. It is widely recognized as the first successful application in this relatively new field of research. Mnih et al., 2015 proposed an extension of the Q-learning approach explored in Section 2.3.2 by approximating the action-value function with a neural network. The algorithm is commonly referred to as Deep-Q-Network (DQN),

but it is sometimes called Deep Q-learning. It represents a natural extension of the former tabular approach since a neural network can receive real values as input, giving rise to new opportunities for reinforcement learning applicability.

Value-based methods estimate the action-value function by using the Bellman equation as an iterative update,

$$Q_{i+1}(s, a) = \mathbb{E} \left[r + \gamma \max_{a'} Q_i(s', a') \mid s, a \right] \quad (2.38)$$

where i is the number of iteration.

Value iteration algorithms should, in principle, converge to the optimal function (Watkins and Dayan, 1992), $Q_i \rightarrow Q^*$ as $i \rightarrow \infty$, but this simple approach lacks any kind of generalization, i.e., it must be run with a table of all the possible state-actions pairs. Using a neural network to parametrize the value function helps to avoid the use of a table, which would be infeasible for high dimensional state space or in the case of real-valued states. This approach aims to generalize over previous unseen states so that the learning process consists in finding the right set of parameters for the neural network in a way that it will be able to recognize similar states without storing an entire database of possible occurrences. In DQN, the approximated action-value function $Q(s, a; \theta)$ is a multi-layer neural network that for a given input state s returns a vector of action values. Hence, the state space \mathcal{S} becomes continuous with respect to the Q-learning, while the action space \mathcal{A} remains discrete. Indeed, the benefit introduced by DQN is the possibility to do discrete control in an environment that is properly represented by continuous state variables.

However, a nonlinear function approximator to estimate the action-value function, $Q^*(s, a) \approx Q(s, a; \theta)$ is known to produce unstable behavior during the training process and lacks background theory to guarantee convergence. Deep reinforcement learning methods are prone to diverge because of the deadly triad (Sutton and Barto, 2018) whenever function approximators are combined with an off-policy algorithm and learning by current estimates as in classical Q-learning. Like Q-learning, DQN is an off-policy method, and it can use any behavioral policy while learning the greedy behavior, so it is subjected to the deadly triad.

A Q-network can be trained by sampling uniformly sequences of experience $(s_t, a_t, r_{t+1}, s_{t+1}) \sim D$ from a buffer $\mathcal{D} = e_1, \dots, e_N$ of fixed dimensionality N . The dataset \mathcal{D} represents the replay memory of the agent, which

can generate sample streams of experience collected through the behavioral policy and perform minibatch updates of the network. After \mathcal{D} reaches its maximum dimension N , new streams are stored in place of the oldest, or alternatively one can keep the dataset growing in size. Sampling experiences from a previously filled buffer represents a training heuristic referred to as experience replay (Lin, 1992) to help with the stability of the learning process. Indeed, suppose the gradient is updated using batches of sequential data. There is the risk of diverging from the optimum because the samples are repeatedly collected using the current policy, which could become suboptimal during training and collect misleading sequences of experience for the learning algorithm. Instead, the experience replay allows each sequence to be potentially involved in many updates and be sampled several times. Random sampling between a stream of experiences that also occurred at a very distant point tends to break the correlation between consecutive samples in an online algorithm and reduce the variance of the update.

Recalling the structure for the gradient update with action-value function approximation in 2.35, DQN uses an estimate $y_i = r + \gamma \max_{a'} Q(s', a'; \theta_{i-1})$ as the target for the true value function. At each time step i , the target y_i is time-varying because it depends on the parameter θ_{i-1} at the previous time step. To avoid the problem of learning by current estimates because the true value target is expressed in terms of the current value function approximation, DQN fixes the neural network's weights to compute the target for the action-value function updates. The target weights can be fixed for a certain number of iterations, a single step in the simplest case. This heuristic helps gain stability during training, although theoretical proof does not back it. This mechanism is put into practice by concurrently maintaining two neural networks while training the algorithm, one for the current value function estimate and one for the target estimate, whose parameters are fixed.

It follows that the algorithm minimizes a loss function $L_i(\theta_i)$ that changes at each iteration i

$$L_i(\theta_i) = \mathbb{E}_{(s,a,r,s') \sim \mathcal{D}} [(y_i - Q(s, a; \theta_i))^2] \quad (2.39)$$

The parameters θ_{i-1} of the target value function are held fixed to cast the problem as in a supervised setting, where the targets are fixed from the beginning of the learning process.

The loss function differentiated with respect to the network weights re-

sults in the following gradient,

$$\begin{aligned} & \nabla_{\theta_i} L_i(\theta_i) \\ &= \mathbb{E}_{(s,a,r,s') \sim \mathcal{D}} \left[\left(r + \gamma \max_{a'} Q(s', a'; \theta_{i-1}) - Q(s, a; \theta_i) \right) \nabla_{\theta_i} Q(s, a; \theta_i) \right] \end{aligned} \quad (2.40)$$

As outlined in the previous section, it is more efficient to compute the gradient of the Eq. 2.40 through a batch method that samples a certain amount of sequences e_t , rather than averaging over all of those included in \mathcal{D} at each time.

DQN shares a significant problem with standard Q-learning, which is an overestimation of the action-value function that has been investigated by Van Hasselt et al., 2016. The estimation of the action-value function performed by DQN produces a much larger value than the real one. In the bootstrapping estimates there are two sources of noise which are correlated. Indeed, we recall that the target for the DQN update is computed as

$$y_i^{\text{DQN}} = r + \gamma \max_{a'} Q(s', a'; \theta_{i-1}). \quad (2.41)$$

Then we know that

$$\max_{a'} Q(s', a'; \theta_{i-1}) = Q\left(s', \arg \max_{a'} Q(s', a'; \theta_{i-1}); \theta_{i-1}\right), \quad (2.42)$$

which shows that the same noise affect the maximization over the action space and the value function estimates. Decoupling these two sources of noise helps obtain an estimate which is less biased, although with more variance, hence avoiding the overestimation problem. The proposed solution is called Double Q-learning, formulated by Hasselt, 2010 for the tabular case and then extended by Van Hasselt et al., 2016 in the deep reinforcement learning framework. The main novelty of Double Q-learning relies on the calculation of the target y_i for each iteration in the Eq. 2.39. The maximization in the Eq. 2.41 uses the same value function to select and evaluate an action, whereas Double DQN decouples the selection of the action from the evaluation:

$$y_i^{\text{Double DQN}} = r + \gamma Q\left(s', \arg \max_{a'} Q(s', a'; \theta_i); \theta_{i-1}\right) \quad (2.43)$$

The original Double Q-learning algorithm employs two tabular value functions to compute the update. At each step in the learning process, one

Q-table is used to determine the greedy action and the other to determine its value. To keep the algorithm simple and similar to the original DQN, Double DQN still employs two neural networks: one computes the target, and the other computes the current action-value function. In the deep reinforcement learning algorithm, the computation of the target is split between the current neural network that selects greedily the action and the target neural network that evaluates such action. The selection of the action in equation 2.43 is still due to the current weights θ_i , whereas the fixed set of weights of the target network is only used to evaluate that selected policy. The pseudocode for the original version of the DQN algorithm is provided below.

Algorithm 2 Deep Q-learning with experience replay pseudocode

```

Initialize replay memory  $\mathcal{D}$  to capacity  $N$ 
Initialize action-value function  $Q(s, a; \theta_0)$  with random weights  $\theta$ 
Initialize target action-value function  $Q(s, a; \theta')$  with random weights  $\theta' = \theta$ 
for episode = 1,  $M$  do
  Initialize sequence  $s$ 
  for each step of episode  $t = 1, T$  do
    With probability  $\epsilon$  select a random action  $a$ 
    otherwise select  $a = \arg \max_a Q(s, a; \theta)$ 
    Execute action  $a$  in the environment, observe reward  $r$  and next state  $s'$ 
    Store transition  $e = (s, a, r, s')$  in  $\mathcal{D}$ 
    Sample random batch of  $P$  transitions  $\{e_1 \dots e_P\}$  from  $\mathcal{D}$ 
    Set  $y = r$  if episode terminates at step  $t + 1$ 
    Otherwise set  $y = r + \gamma \max_{a'} Q(s', a'; \theta')$ 
    Perform a gradient descent step  $(y - Q(s, a; \theta))^2$  on with respect to network parameters  $\theta$ 
    Every  $c$  steps reset  $\theta' = \theta$ 
  end for
end for

```

Other known algorithms in the value-based family of deep reinforcement learning are Hindsight Experience Replay (Andrychowicz et al., 2017), which proposes a different use of the replay memory with respect to DQN, distribu-

tional variants of value-based such as Quantile Regression DQN (Bellemare et al., 2017) or Dabney et al., 2018, where they attempt to learn the entire distribution of the value function, instead of just an estimate. For what concern hybrid actor-critic algorithm, deep deterministic policy gradient (Lillicrap et al., 2015) and their improvements (Fujimoto et al., 2018) represents extensions of DQN that also use a policy function approximation. More details about these techniques are given in the next section.

2.4.2 Policy-based methods

The policy-based methods directly parametrizes a policy without inferring it from a value function and represent a possible different approach for the same set of reinforcement learning problems. Once chosen the parametric form of the policy, the goal of each policy-based algorithm is to find the optimal set of parameters θ for the policy $\pi_\theta(s, a)$. The first important point to stress is the choice of the metrics used to measure the quality of a policy. Depending on the type of MDP, one can use the initial value in the first state s_1

$$J_1(\theta) = V^{\pi_\theta}(s_1) = \mathbb{E}_{\pi_\theta}[v_1] \quad (2.44)$$

when the environment is episodic, or use the average value instead

$$J_{avg}(\theta) = \sum_s d^{\pi_\theta}(s) V^{\pi_\theta}(s) \quad (2.45)$$

when the environment is continuing, with $d^{\pi_\theta}(s)$ the stationary distribution over the state associated to the MDP given π_θ .

Solving an MDP through a policy-based method is essentially an optimization problem where one has to find the set of parameters θ that maximizes $J(\theta)$. Among the possible different approaches, such as hill climbing (Wardell and Peterson, 2006), simplex algorithm (Lee et al., 2017), Nelder-Mead (Sarakhsi et al., 2016), genetic algorithm (Daniel and Rajendran, 2005) and conjugate gradient (Baxter and Bartlett, 2001; Cohen et al., 2019), we are going to focus on gradient descent methods, which have proven to be effective in optimizing the neural network family of models (Engstrom et al., 2020; Goodfellow et al., 2016).

After choosing an objective function $J(\theta)$ for the problem that we want to solve through a policy-based method, the algorithm search for a local

maximum of $J(\theta)$ by ascending the gradient of the policy with respect to its parameters θ , so that the parameter update is

$$\Delta\theta = \alpha \nabla_{\theta} J(\theta), \quad (2.46)$$

where α is a step-size parameter, also referred to as learning rate.

Although some methods tends to approximate the gradient of the function $J(\theta)$, as happens for the finite difference methods (Peters and Bagnell, 2010), such gradient can be analytically computed. The policy gradient theorem provides a way to compute the gradient of the selected objective function using the likelihood ratio approach. It is widely applicable to different objective, either the start state objective in Eq. 2.44 or the average value objective in 2.45. The theorem, which derives from two independent works of Sutton et al., 1999 and Marbach and Tsitsiklis, 2001, provides an analytical expression for the gradient of $J(\theta)$ so that for any differentiable policy $\pi_{\theta}(s, a)$ and for any of the policy objective functions described in this section, the policy gradient is

$$\begin{aligned} \nabla_{\theta} J(\theta) &= \mathbb{E}_{\pi_{\theta}} \left[\frac{\nabla_{\theta} \pi_{\theta}(s, a)}{\pi_{\theta}(s, a)} Q^{\pi_{\theta}}(s, a) \right] \\ &= \mathbb{E}_{\pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta}(s, a) Q^{\pi_{\theta}}(s, a)], \end{aligned} \quad (2.47)$$

The derivation of Eq. 2.47 follows from the proof of the policy gradient theorem. We present the episodic formulation of the theorem, although there is also the version for continuing tasks. We recall that in an episodic problem, the on-policy distribution depends on the initial state of the single episode. We denote $h(s)$ as the probability for the episode to begin in a state s and $\eta(s)$ the average number of steps spent by the agent in that specific state. Therefore we consider an agent as spending time in a state s , if the start of the episode is exactly that or if there is a transition to s from a preceding state \bar{s} . The time spent in a state is hence written as

$$\eta(s) = h(s) + \sum_{\bar{s}} \eta(\bar{s}) \sum_a \pi(\bar{s}, a) p(s | \bar{s}, a), \quad \text{for all } s \in \mathcal{S}. \quad (2.48)$$

It follows that the on-policy distribution is the amount of time spent in each state normalized to sum up to one:

$$\mu(s) = \frac{\eta(s)}{\sum_{s'} \eta(s')}, \quad \text{for all } s \in \mathcal{S} \quad (2.49)$$

2. Reinforcement learning: theory and background

Then we remark that the gradient of the state-value function can be written in terms of the action-value function $\forall s \in \mathcal{S}$ as

$$\begin{aligned}
\nabla_{\theta} v_{\pi_{\theta}}(s) &= \nabla_{\theta} \left[\sum_a \pi_{\theta}(s, a) q_{\pi_{\theta}}(s, a) \right], \\
&= \sum_a [\nabla_{\theta} \pi_{\theta}(s, a) q_{\pi_{\theta}}(s, a) + \pi_{\theta}(s, a) \nabla_{\theta} q_{\pi_{\theta}}(s, a)] \\
&= \sum_a \left[\nabla_{\theta} \pi_{\theta}(s, a) q_{\pi_{\theta}}(s, a) + \pi_{\theta}(s, a) \nabla_{\theta} \sum_{s', r} p(s', r | s, a) (r + v_{\pi_{\theta}}(s')) \right] \\
&= \sum_a \left[\nabla_{\theta} \pi_{\theta}(s, a) q_{\pi_{\theta}}(s, a) + \pi_{\theta}(s, a) \sum_{s'} p(s' | s, a) \nabla_{\theta} v_{\pi_{\theta}}(s') \right] \\
&= \sum_a \left[\nabla_{\theta} \pi_{\theta}(s, a) q_{\pi_{\theta}}(s, a) + \pi_{\theta}(s, a) \sum_{s'} p(s' | s, a) \right. \\
&\quad \left. \sum_{a'} \left[\nabla_{\theta} \pi_{\theta}(a' | s') q_{\pi_{\theta}}(s', a') + \pi_{\theta}(a' | s') \sum_{s''} p(s'' | s', a') \nabla_{\theta} v_{\pi_{\theta}}(s'') \right] \right] \\
&= \sum_{x \in \mathcal{S}} \sum_{k=0}^{\infty} \Pr(s \rightarrow x, k, \pi_{\theta}) \sum_a \nabla_{\theta} \pi_{\theta}(a | x) q_{\pi_{\theta}}(x, a). \tag{2.50}
\end{aligned}$$

After repeated unrolling, as in the second last step above, where $\Pr(s \rightarrow x, k, \pi_{\theta})$ is the probability of transitioning from state s to state x in k steps under policy π_{θ} , it follows that

$$\begin{aligned}
\nabla_{\theta} J(\theta) &= \nabla_{\theta} v_{\pi_{\theta}}(s_0) = \\
&= \sum_s \left(\sum_{k=0}^{\infty} \Pr(s_0 \rightarrow s, k, \pi_{\theta}) \right) \sum_a \nabla_{\theta} \pi_{\theta}(s, a) q_{\pi_{\theta}}(s, a) \\
&= \sum_s \eta(s) \sum_a \nabla_{\theta} \pi_{\theta}(s, a) q_{\pi_{\theta}}(s, a) \\
&= \sum_{s'} \eta(s') \sum_s \frac{\eta(s)}{\sum_{s'} \eta(s')} \sum_a \nabla_{\theta} \pi_{\theta}(s, a) q_{\pi_{\theta}}(s, a) \\
&= \sum_{s'} \eta(s') \sum_s \mu(s) \sum_a \nabla_{\theta} \pi_{\theta}(s, a) q_{\pi_{\theta}}(s, a) \\
&\propto \sum_s \mu(s) \sum_a \nabla_{\theta} \pi_{\theta}(s, a) q_{\pi_{\theta}}(s, a) \quad \blacksquare \tag{2.51}
\end{aligned}$$

2. Reinforcement learning: theory and background

In the episodic case, the constant of proportionality in the above derivation is the average length of an episode, while in the continuing case the relationship becomes an equality.

The analytical expression of the policy gradient theorem can be evaluated if we are able to sample experiences whose expectations approximate the analytical expression. In other words, we need a way to approximate the true value $Q^{\pi_\theta}(s, a)$ with an estimate so that the resulting policy gradient is proportional to Eq. 2.47. The REINFORCE algorithm, whose pseudocode is provided below, is one of the simplest policy-based algorithm that exploits the policy gradient theorem by replacing the true value function with the total return of the episode

$$\Delta\theta_t = \alpha \nabla_\theta \log \pi_\theta(s_t, a_t) G, \quad (2.52)$$

where $G = \sum_{k=t+1}^T \gamma^{k-t-1} R_k$. REINFORCE represents one among different policy gradient methods proposed in the literature and its need of the total return as estimate for the true value function makes it only suitable for episodic MDPs. A comprehensive overview of different policy gradient methods is provided in Peters and Bagnell, 2010.

Algorithm 3 REINFORCE pseudocode

Require: A differentiable policy parametrization $\pi_\theta(s, a)$

Set a scalar step size $\alpha > 0$

Initialize policy parameter $\theta \in \mathbb{R}^d$, e.g. to 0

for each episode = 1, M **do**

 Generate an episode $s_0, a_0, r_1, \dots, s_{T-1}, a_{T-1}, R_T$ following $\pi_\theta(\cdot, \cdot)$

for each step of the episode $t = 0, 1, \dots, T - 1$ **do**

$G \leftarrow \sum_{k=t+1}^T \gamma^{k-t-1} R_k$

$\theta \leftarrow \theta + \alpha \nabla_\theta \log \pi_\theta(s_t, a_t) G$

end for

end for

The major drawback of each policy gradient algorithm is the presence of high variance in the estimation of the action-value function. It is common for this family of algorithms to reduce such variance using a *critic*, i.e., an additional parametrization for the action-value function so that

$$Q_w(s, a) \approx Q^{\pi_\theta}(s, a). \quad (2.53)$$

This approach is called actor-critic, which maintains two sets of parameters to approximate both the action-value function through the set w and the policy function through the set θ . Such a hybrid approach is placed between the value-based and policy-based methods since it exploits both techniques to solve the MDP. However, it is generally presented as an extension of the policy-based family, i.e., as a way to drive the learned policy toward a solution closer to the optimum. An actor-critic algorithm follows an approximate policy gradient of the form

$$\nabla_{\theta} J(\theta) \approx \mathbb{E}_{\pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta}(s, a) Q_w(s, a)] \quad (2.54)$$

$$\Delta \theta = \alpha \nabla_{\theta} \log \pi_{\theta}(s, a) Q_w(s, a) \quad (2.55)$$

where the critic is helpful to solve the problem of policy evaluation that we described in the first part of this chapter.

The approximation of the policy gradient could potentially introduce a bias in the policy optimization process, which prevents the possibility of reaching an optimal solution. A theoretical result derived in Sutton et al., 1999 helps ensure a proper approximation for the value function in the actor-critic setting so that no bias is introduced. The theorem proves that a parametrized action-value function that satisfies certain conditions can be used in an actor-critic policy gradient setting because the resulting policy gradient is unbiased with respect to the true policy gradient in Eq. 2.47. We refer to the original paper for the detailed statement and proof. Therefore, we can approximate the value function without the risk of introducing a bias in the policy gradient.

In an actor-critic setting it is also common to subtract a baseline function $B(s)$ from the value function estimate in order to reduce the variance of the estimate without adding a bias. Indeed one can show that $\mathbb{E}_{\pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta}(s, a) Q_w(s, a)] = \mathbb{E}_{\pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta}(s, a) (Q_w(s, a) - B(s))]$ because

$$\begin{aligned} \mathbb{E}_{\pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta}(s, a) B(s)] &= \sum_{s \in \mathcal{S}} d^{\pi_{\theta}}(s) \sum_a \nabla_{\theta} \pi_{\theta}(s, a) B(s) \\ &= \sum_{s \in \mathcal{S}} d^{\pi_{\theta}}(s) B(s) \nabla_{\theta} \sum_{a \in \mathcal{A}} \pi_{\theta}(s, a) \\ &= 0, \end{aligned} \quad (2.56)$$

which therefore does not change the expectation of the policy gradient as stated in Eq. 2.54.

A commonly used baseline is the state-value function $B(s) = V^{\pi_\theta}(s)$, which transform the policy gradient to

$$A^{\pi_\theta}(s, a) = Q^{\pi_\theta}(s, a) - V^{\pi_\theta}(s) \tag{2.57}$$

$$\nabla_\theta J(\theta) = \mathbb{E}_{\pi_\theta} [\nabla_\theta \log \pi_\theta(s, a) A^{\pi_\theta}(s, a)], \tag{2.58}$$

where $A^{\pi_\theta}(s, a)$ is called advantage function. In this way, the critic would consist in estimating both $V^{\pi_\theta}(s)$ and $Q^{\pi_\theta}(s, a)$ with two different estimator at the same time, resulting into

$$V_v(s) \approx V^{\pi_\theta}(s) \tag{2.59}$$

$$Q_w(s, a) \approx Q^{\pi_\theta}(s, a) \tag{2.60}$$

$$A(s, a) = Q_w(s, a) - V_v(s). \tag{2.61}$$

However, some approaches in the literature allow estimating the advantage function for the policy gradient by using a critic with a single set of parameters, as with the generalized advantage estimation (GAE) (Schulman et al., 2015b).

Proximal Policy Optimization

This subsection presents a deep reinforcement learning policy gradient algorithm that follows the framework presented above. Before introducing the algorithm we will use in the practical application of this work, it is worth remembering some pioneering policy gradient methods that adopted neural networks as function approximators to understand the novelty and the problem that such an algorithm proposes to solve.

If we assume the policy function to be deterministic, the output of the policy $\mu(s)$ is a single scalar value that depends on the state. The deterministic policy gradient (DPG) algorithm (Silver et al., 2014) parametrize such a policy as $\mu_\theta(s)$ and optimizes the following objective:

$$J(\theta) = \int_{\mathcal{S}} \rho^\mu(s) Q(s, \mu_\theta(s)) ds \tag{2.62}$$

where $\rho^\mu(s)$ is the discounted state distribution. To compute the gradient of the objective above, we firstly compute the gradient of the function Q with respect to the action a and then we take the gradient of the deterministic

policy function μ with respect to θ as follows:

$$\nabla_{\theta} J(\theta) = \int_{\mathcal{S}} \rho^{\mu}(s) \nabla_a Q^{\mu}(s, a) \nabla_{\theta} \mu_{\theta}(s) ds \quad (2.63)$$

$$= \mathbb{E}_{s \sim \rho^{\mu}} [\nabla_a Q^{\mu}(s, a) \nabla_{\theta} \mu_{\theta}(s)], \quad (2.64)$$

where $a = \mu_{\theta}(s)$. The deterministic policy gradient is considered a particular case of the stochastic policy gradient explored in the previous section when the probability distribution represented by the policy contains only one nonzero value over one action. Silver et al., 2014 show that one can parametrize a stochastic policy $\pi_{\theta}(a | s)$ by learning a deterministic policy μ_{θ} and an additional parameter σ which represents the variance, i.e. the confidence around the deterministic choice of the policy μ_{θ} . Hence, the stochastic policy is equivalent to the deterministic one when the variance σ approaches zero. The advantage of focusing on a deterministic policy rather than a stochastic one is that the latter is known to require more samples of experience because it needs integration over the entire state and action spaces.

However, a deterministic policy suffers from a lack of exploration during training because it always acts greedily according to the current $\mu_{\theta}(s)$ that maximizes the objective in Eq. 2.62 when the agent interacts with the environment. Two solutions are viable to overcome the exploration problem. Either the deterministic policy is perturbed with a noise, which results in transforming the original policy into a stochastic one, or alternatively, a noisy policy is directly used as a behavior to collect experience in an off-policy setting. Hence, in the former approach, the trajectories are generated by a policy $\beta(a | s)$, and the state distribution follows the state density ρ^{β}

$$J_{\beta}(\theta) = \int_{\mathcal{S}} \rho^{\beta} Q^{\mu}(s, \mu_{\theta}(s)) ds \quad (2.65)$$

$$\nabla_{\theta} J_{\beta}(\theta) = \mathbb{E}_{s \sim \rho^{\beta}} [\nabla_a Q^{\mu}(s, a) \nabla_{\theta} \mu_{\theta}(s)]. \quad (2.66)$$

Since we are learning a deterministic policy, we require $Q^{\mu}(s, \mu_{\theta}(s))$ instead of the average $\sum_a \pi(a | s) Q^{\pi}(s, a)$ as an estimate for the reward of a given state s .

Although the exploration issue can be weakened by the off-policy setting, policy gradient methods in general exhibit training instability when the policy and eventually the value function, if in an actor-critic setting, are parametrized by neural networks. It becomes important to limit the updates

of the parameters that could move the policy too far from the current one in just one step of the updating process. The Trust Region Policy Optimization (TRPO) (Schulman et al., 2015a) is one of the first developed methods that attempts to solve the instability problem for deep reinforcement learning by using a Kullback-Leibler divergence constraint on the size of the policy update.

Assuming we are still in the off-policy setting and the behavioral policy β which gathers experiences is different from the policy π that needs to be learned, the objective function is expressed as a measure of the total advantage over states and actions

$$\begin{aligned}
 J(\theta) &= \sum_{s \in \mathcal{S}} \rho^{\pi_{\theta_{\text{old}}}} \sum_{a \in \mathcal{A}} \left(\pi_{\theta}(a | s) \hat{A}_{\theta_{\text{old}}}(s, a) \right) \\
 &= \sum_{s \in \mathcal{S}} \rho^{\pi_{\theta_{\text{old}}}} \sum_{a \in \mathcal{A}} \left(\beta(a | s) \frac{\pi_{\theta}(a | s)}{\beta(a | s)} \hat{A}_{\theta_{\text{old}}}(s, a) \right) \\
 &= \mathbb{E}_{s \sim \rho^{\pi_{\theta_{\text{old}}}}, a \sim \beta} \left[\frac{\pi_{\theta}(a | s)}{\beta(a | s)} \hat{A}_{\theta_{\text{old}}}(s, a) \right] \tag{2.67}
 \end{aligned}$$

where $\rho^{\pi_{\theta_{\text{old}}}}$ is the state distribution according to the policy $\pi_{\theta_{\text{old}}}$ before the update and $\beta(a | s)$ is the behavior policy for collecting trajectories. An advantage estimator $\hat{A}(s, a)$ replaces the true advantage because the true rewards are unknown at the time of the update.

Training the same algorithm in an on-policy setting means that the behavioral and the target policies collapse into the same. The behavior policy is hence denoted as $\pi_{\theta_{\text{old}}}(a | s)$, since it is just the current policy before the newest parameter update so that the objective function of the problem becomes

$$J(\theta) = \mathbb{E}_{s \sim \rho^{\pi_{\theta_{\text{old}}}}, a \sim \pi_{\theta_{\text{old}}}} \left[\frac{\pi_{\theta}(a | s)}{\pi_{\theta_{\text{old}}}(a | s)} \hat{A}_{\theta_{\text{old}}}(s, a) \right]. \tag{2.68}$$

TRPO employs a trust region constraint, hence the name of the algorithm, which enforces the distance between the old policy distribution and the new policy distribution to be small by computing the Kullback-Leibler divergence,

$$\mathbb{E}_{s \sim \rho^{\pi_{\theta_{\text{old}}}}} [D_{\text{KL}}(\pi_{\theta_{\text{old}}}(\cdot | s) \| \pi_{\theta}(\cdot | s)) \leq \delta], \tag{2.69}$$

where the parameter δ represents the tolerance for the distance between the two distributions. The constraint helps in keeping the old and the new policy close enough to avoid large update of the parameters. TRPO still guarantees

a monotonic improvement of the policy (Schulman et al., 2017), but it is complicated to implement the constraint and compute it in practice. For this reason, Schulman et al., 2017 developed an algorithm called Proximal Policy Optimization (PPO), which still enforces the policy update to be small without the need of a constraint, while instead clipping the objective function of the problem.

To understand how this algorithm works, let denote the probability ratio between old and new policies as:

$$r(\theta) = \frac{\pi_{\theta}(a | s)}{\pi_{\theta_{\text{old}}}(a | s)}, \quad (2.70)$$

which allows writing the objective function of the on-policy TRPO algorithm as

$$J^{\text{TRPO}}(\theta) = \mathbb{E} \left[r(\theta) \hat{A}_{\theta_{\text{old}}}(s, a) \right]. \quad (2.71)$$

However, we already know that the maximization of this objective function could result in serious instability issues because we are not constraining the distance between the policies, parametrized respectively by θ_{old} and θ . PPO circumvents the problem by imposing the ratio $r(\theta)$ to be within a small interval around 1, precisely $[1 - \epsilon, 1 + \epsilon]$, where ϵ is a scalar parameter,

$$J^{\text{CLIP}}(\theta) = \mathbb{E} \left[\min \left(r(\theta) \hat{A}_{\theta_{\text{old}}}(s, a), \text{clip}(r(\theta), 1 - \epsilon, 1 + \epsilon) \hat{A}_{\theta_{\text{old}}}(s, a) \right) \right]. \quad (2.72)$$

The operator $\text{clip}(r(\theta), 1 - \epsilon, 1 + \epsilon)$ has the role to maintain the ratio between $1 + \epsilon$ and $1 - \epsilon$. Therefore, the objective function of the PPO algorithm takes the minimum one between the original value and the clipped version. Since the PPO algorithm is usually applied within an actor-critic framework, although it can also be used in a policy-based setting, and the parameters of the policy and the value network are shared, we can write an augmented objective function as

$$J^{\text{CLIP}'}(\theta) = \mathbb{E} \left[J^{\text{CLIP}}(\theta) - c_1 (V_{\theta}(s) - V_{\text{target}})^2 + c_2 H(s, \pi_{\theta}(\cdot)) \right] \quad (2.73)$$

which accounts for the error term in the value estimation and the entropy term that ensures sufficient exploration during training. The parameter c_1 and c_2 are constant and regulate the contribution of each term to the whole objective.

Other than the policy gradient algorithms presented in this chapter, there are other known variants of the similar approximation task, such as Asynchronous Advantage Actor-Critic (A3C) (Mnih et al., 2016), where a single critic network learn the value function. At the same time, multiple actors are trained in parallel and are synced with global parameters for a more efficient training runtime. The same work also shows the synchronous version of the same algorithm, where the global parameters of the actor are updated once each of the multiple actors has finished collecting experiences. Phasic policy gradient (PPG) (Cobbe et al., 2021) is an on-policy actor-critic policy gradient algorithm, which resembles a PPO algorithm with separate training phases for policy and value functions. The actor-critic with experience replay (ACER) (Wang et al., 2016) is an off-policy algorithm that uses experience replay, significantly increasing the sample efficiency and decreasing the data correlation, and represents the off-policy counterpart of the A3C algorithm. Soft Actor-Critic (SAC) (Haarnoja et al., 2018) can incorporate the entropy measure of the policy into the reward to encourage exploration, making the policy learned highly random. It is an off-policy actor-critic model following a maximum entropy approach. However, this set of algorithms is not exhaustive since many slight variants come out to tackle the specific disadvantages of each setting.

2.5 Summary

This chapter introduced reinforcement learning through its common root within the dynamic programming framework. First, some basics notion about the mathematical formalization of the control problem is given so that we introduce the agent-environment interaction within the MDP framework. On top of that, we presented the tabular algorithms, which still suffer the curse of dimensionality of dynamic programming in the case of high dimensional control problems. For that reason, detailed exposure to deep reinforcement learning algorithms is provided. We introduced two families of algorithms, value-based and policy-based, that provide a broad view to solve problems of a different nature. The entire chapter is instrumental to the empirical part of this work, where we employ the described algorithms to solve problems in finance and economics.

Part II

Reinforcement learning: simulation analysis with applications to finance

This part of the work contains three essays regarding different applications of reinforcement learning to the field of finance. Chapter 3 explores the use of reinforcement learning algorithms, either from the value-based and the policy-based family, to solve a financial trading problem in a controlled environment. Through an extensive series of simulations, we show that reinforcement learning algorithms are more flexible in solving the trading problem with respect to a dynamic programming approach, which we use as a benchmark solution, especially when the underlying dynamics of the simulated assets are misspecified. Chapter 4 extends this idea to a multi-asset trading framework, where returns are still simulated according to known dynamics, but there is no closed-form solution. We show that a model-free reinforcement learning approach has a hard time finding the solution to the trading problem in a finite amount of training time, and we resort to a different approach to leverage prior knowledge of the financial problem. Finally, Chapter 5 explores the application of reinforcement learning to an agent-based model of banks that exchanges liquidity to fulfill their obligations. We conclude by highlighting the importance of the reinforcement learning paradigm in finance for improving existing solutions to problems in the financial domain by leveraging the need for fewer modeling assumptions.

Chapter 3

Deep Reinforcement Trading with Predictable Returns

3.1 Introduction

The important milestone represented by modern portfolio theory of Markowitz, 1952 has set the basis for the beginning of financial portfolio optimization as an active field of research. Its original formulation suffers several drawbacks (Kolm et al., 2014) and has been extended from a single-period to a multi-period framework to capture intertemporal effects and to allow dynamical portfolio rebalancing (Engle and Ferstenberg, 2007; Grinold, 2006; Kolm and Maclin, 2012; Kolm and Ritter, 2014; Tutuncu, 2011). However, the addition of the time dimension makes even more complicated the estimation of an optimal strategy, which requires forecasting financial quantities such as risks and returns for several periods in the future. Single-period models are often still adopted because their dynamic counterpart is not practical and the forecasting step may lead to systematic errors due to the uncertainty about the chosen model or the inherent presence of a low signal-to-noise ratio in the financial data. Even when a multi-period model is effective in capturing the market impact or alpha decay, classical optimal control techniques lay over a set of restricting assumptions which cannot properly represent the real financial world.

In this work, we use reinforcement learning (RL) (Sutton and Barto, 2018; Szepesvári, 2010) as a convenient framework to model sequential decision problems of a financial nature without the need of directly modeling the

underlying asset dynamics. RL finds its roots in the optimal control theory along with the dynamic programming literature (Bertsekas, 2005) and has gained a huge revival after the last decade improvement of deep learning (DL) as a field of research. This gave rise to the so-called deep reinforcement learning (DRL) that has already obtained relevant results in application domains such as gaming (Mnih et al., 2015; Silver et al., 2017b) and robotics (Levine et al., 2016). For a comprehensive overview of DRL methods and its fields of application, see Arulkumaran et al., 2017.

The RL approach is not new to the financial domain, and there are examples of practical applications for trading and portfolio management (Jiang et al., 2017; Zhang et al., 2020). However, recent DRL algorithms are very often home-made recipes without theoretical control. For this reason, the study of their performances in real financial trading problems is always an intricate combination of different effects, some of them related to the goodness of the dataset and the signals used to predict returns, some others related to the specific algorithm and its trainability issues.

To the best of our knowledge, there is a lack of research works that investigate DRL performances in financial trading problems besides the issues coming from market efficiency: the search for a good signal to predict returns or the possible lack of any signals in the dataset. For this reason, we consider a controlled environment in which a signal is known to exist and study the capability of DRL agents to discover profitable opportunities in the market.

Similarly to Chaouki et al., 2020; Kolm and Ritter, 2019a, we simulate financial asset returns which contain predictable factors, and we let the agent trade in an environment whose associated optimization problem admits an exact solution (Gârleanu and Pedersen, 2013). The optimal benchmark strategy allows us to evaluate strengths and flaws of a DRL approach, both in terms of accuracy and efficiency.

As a main novelty of our work, we exploit a data-driven setting of DRL in which the agents not only compete against classical strategies but can also leverage on their experience to optimize the state-action space and increase the learning speed. We test different DRL approaches on a variety of financial data with different properties to investigate their flexibility when the simulated dynamics is misspecified with respect to the assumptions of the benchmark model. We show that DRL algorithms can reach the performance of the benchmark strategy, when it is optimal, and also outperform it in the case of model misspecifications like the presence of extreme events and volatility clustering. We also show that classical strategies can help DRL

agents by giving them information about the typical scale of a good strategy to start and adjust.

3.2 Financial Market Environment

The agent operates in a financial market where at each time $t \in \mathbb{Z}$ it can trade N securities whose excess returns $y_{t+1} = p_{t+1} - (1 + r_f)p_t$ are given by

$$y_{t+1} = Bf_t + u_{t+1}, \quad (3.1)$$

where f_t is a $K \times 1$ vector of return-predicting factors, B is a matrix of factor loadings and u_{t+1} is a noise term with $\mathbb{E}[u_{t+1}] = 0$ and $\text{Var}[u_{t+1}] = \Sigma$.

The factors can be either value factors, which describes the profitability of the asset relatively to some fundamental measure, or momentum factors, which rely on past price movements to predict the future (Moskowitz et al., 2012). We assume they evolve according to a discretization of a mean-reverting process (Uhlenbeck and Ornstein, 1930)

$$\Delta f_{t+1} = -\phi f_t + \epsilon_{t+1}, \quad (3.2)$$

where ϕ is a $K \times K$ matrix of mean-reversion coefficients and ϵ_{t+1} represents a stochastic shock component with $\mathbb{E}[\epsilon_{t+1}] = 0$ and $\text{Var}[\epsilon_{t+1}] = \Omega$.

Trading in this environment produces transaction costs which we assume to be a quadratic function of the traded amount $\Delta h_t = h_t - h_{t-1}$, i.e.

$$C(\Delta h_t) = \frac{1}{2} \Delta h_t^T \Lambda \Delta h_t, \quad (3.3)$$

where Λ is a symmetric positive definite matrix ensuring transaction costs convexity as generally required by empirical literature (Gârleanu et al., 2008; Lillo et al., 2003) and is consistent with the assumption of a linear price impact. In the following, we also assume that $\Lambda = \lambda \Sigma$, i.e. that trading costs are actually the compensation for the dealer's risk that takes the other part of the transaction. In this context, λ can be interpreted as the dealer's risk aversion and controls the degree of liquidity of the asset.

The agent's goal is to find a dynamic portfolio strategy (h_0, h_1, \dots) by maximizing the present value of all future returns, penalized for risk and net of transaction costs, i.e.

$$\max_{(h_0, h_1, \dots)} \mathbb{E}_0 \left[\sum_t \rho^{t+1} (h_t^T y_{t+1} - \frac{\gamma}{2} h_t^T \Sigma h_t) - + \frac{\rho^t}{2} \Delta h_t^T \Lambda \Delta h_t \right], \quad (3.4)$$

3. Deep Reinforcement Trading with Predictable Returns

where $\rho \in (0, 1)$ is a discount rate and γ is the risk aversion coefficient.

When the noise terms u_t and ϵ_t are assumed to be distributed as a Gaussian, the model coincides with Gârleanu and Pedersen, 2013 that has a closed-form solution as

$$h_t = \left(1 - \frac{a}{\lambda}\right) h_{t-1} + \frac{a}{\lambda} h_t^{aim}. \quad (3.5)$$

The optimal strategy is then a convex combination of holding the previous portfolio and trading towards the objective portfolio h_t^{aim} with trading rate a/λ . Such trading rate $\frac{a}{\lambda} < 1$, where

$$a = \frac{-(\gamma(1 - r_f) + \lambda r_f) + \sqrt{(\gamma(1 - r_f) + \lambda r_f)^2 + 4\gamma\lambda(1 - r_f)^2}}{2(1 - r_f)} \quad (3.6)$$

is a decreasing function of the transaction costs by the effect of λ and increasing in the risk aversion γ . The objective portfolio h_t^{aim} is defined by

$$h_t^{aim} = (\gamma\Sigma)^{-1} B \left(I + \frac{a}{\gamma} \Phi \right)^{-1} f_t, \quad (3.7)$$

being a generalization of the well-known Markowitz portfolio (Markowitz, 1952)

$$h_t^M = (\gamma\Sigma)^{-1} B f_t, \quad (3.8)$$

which is optimal only in the static case and in absence of transaction costs. Instead, the aim portfolio in Eq. (3.5) represents a dynamic strategy and can be shown to be a weighted average of all future Markowitz portfolios.

If the matrix Φ is diagonal, the aim portfolio become

$$h_t^{aim} = (\gamma\Sigma)^{-1} B \left(\frac{f_t^1}{1 + \Phi^1 \frac{a}{\gamma}}, \dots, \frac{f_t^K}{1 + \Phi^K \frac{a}{\gamma}} \right)^T, \quad (3.9)$$

where the K factors are scaled down by their speed of mean-reversion Φ . A factor i with slower speed of mean-reversion is scaled less than a faster factor j and the relative weight of f^i with respect to the weight of f^j , $\frac{1 + \Phi^j \frac{a}{\gamma}}{1 + \Phi^i \frac{a}{\gamma}}$ increases with the transaction cost λ . In fact, the cost friction leads the investor to slow down the rate of portfolio rebalancing and faster factors require closing out the position in a shorter time frame.

In the following we will use the optimal strategy (3.5) of the Gaussian model as a benchmark for the DRL performance in solving the problem

(3.4) but we also consider other possible model specifications for which an explicit optimal solution is not available. In particular, we introduce fat-tailed distributed shocks and heteroskedastic volatility as interesting model misspecifications that reflect general properties of empirical asset returns (Cont, 2001).

A riskier environment with many extreme events is constructed by assuming the asset noise to depart from a Gaussian distribution. In particular, we consider u_t and ϵ_t distributed as a Student's T distribution with ν degrees of freedom. On the other hand, heteroscedasticity is introduced according to a generalized autoregressive conditional heteroskedastic (GARCH) process (Bollerslev, 1987) for the variance of asset returns to model volatility clustering. In the case of a single asset, it means that $u_t = \sigma_t z_t$ where

$$\sigma_t^2 = \omega + \sum_{j=1}^p \alpha_j |u_{j-i}|^2 + \sum_{k=1}^q \beta_k \sigma_{t-k}^2 \quad (3.10)$$

and z_t is a noise term that can be either a standard Gaussian or a Student's T with ν degrees of freedom.

3.3 Deep Reinforcement Learning Methods

The aim of RL is to solve a decision-making problem in which the timing of costs and benefits is relevant. Financial portfolio optimization comprises a set of problems where current actions can influence the future, even at a very distant point in time. RL approaches the resolution of this problem by trial and error, learning by obtaining feedback after each sequential decision.

A RL problem can be formulated in the context of a Markov Decision Process (MDP), which is defined by a set of possible states $S_t \in \mathcal{S}$, a set of possible actions $A_t \in \mathcal{A}$ and a transition probability $\mathcal{P}_{ss'}^a = P[S_{t+1} = s' | S_t = s, A_t = a]$. Therefore, it is the (stochastic) control problem of finding

$$\max_{\{\pi\}} \mathbb{E} \left[\sum_{t=0}^{\infty} \rho^t R_{t+1}(S_t, A_t, S_{t+1}) \right] \quad (3.11)$$

where π defines the agent's strategy that associates a probability $\pi(a | s)$ to the action $A_t = a$ given the state $S_t = s$. A RL agent aims at maximizing the expected sum of (discounted) rewards by finding the best action given

the current state. We consider the model-free context in which the agent has no knowledge of the internal dynamics of the environment, i.e. the transition probability is not known, and the only source of information is the sequence of states, actions and rewards.

Value based methods are defined by introducing the action-value function

$$Q_\pi(s, a) \equiv \mathbb{E} \left[\sum_{k=0}^{\infty} \rho^k R_{t+1+k} \mid S_t = s, A_t = a, \pi \right], \quad (3.12)$$

which reflects the long-term reward associated with the action a taken in the state s if the strategy π is followed hereafter. The estimation of (3.12) allows deriving a deterministic optimal policy as the highest valued action in each state. Depending on how the agent estimates the action-value function (3.12), different classes of value-based algorithms can be introduced. Conversely, direct policy search approaches are alternative methods that try to explore directly the policy space (or some subset of it), being the problem a particular case of stochastic optimization.

3.3.1 Tabular Reinforcement Learning

Tabular RL methods are practical when the possible states and actions are few enough to be represented in a table, which has an entry for every (s, a) pair. In this case, the agent can explore many possible state-action pairs within a reasonable amount of computational time and obtain a good approximation of the value function.

Q-learning (Watkins and Dayan, 1992) is a tabular method in which at each time step the agent tries an action A_t , receives a reward R_{t+1} and updates the current estimate of the action-value function $Q(S_t, A_t)$ as

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha(T_t^Q - Q(S_t, A_t)), \quad (3.13)$$

where α is a learning rate and the target

$$T_t^Q = R_{t+1} + \rho \max_a Q(S_{t+1}, a) \quad (3.14)$$

is a decomposition of the value function in terms of the current reward and the current estimate of the future value discounted by ρ . At the end of the learning process, the optimal policy is the greedy strategy $A_t = \arg \max_a Q(S_t, a)$ but Q-learning is trained off-policy because the agent chooses the action A_t

following a ϵ -greedy policy that ensures adequate exploration of the state-action space.

When states and actions are continuous, as it is for a realistic financial environment, Q-learning barely obtains good estimates of the value function in a feasible computational time. Moreover, the discretization of the state space itself may cause loss of relevant information depending on its granularity. In this context, the DRL framework becomes particularly necessary.

3.3.2 Approximate Reinforcement Learning

DRL algorithms tackle previously intractable problems by approximating Eq. (3.12) through a neural network that allows a continuous state space representation.

Deep Q-Network (DQN) (Mnih et al., 2015) is an extension of Q-learning and allows learning a parametrized value function $Q^*(s, a) \approx Q(s, a; \theta)$. $Q(s, a; \theta)$ is a multi-layer neural network that for a given input state s returns a vector of action values. The standard update of Eq. (3.13) therefore becomes,

$$\theta_{t+1} = \theta_t + \alpha(T_t^{\text{DQN}} - Q(S_t, A_t; \theta_t)) \nabla_{\theta_t} Q(S_t, A_t; \theta_t) \quad (3.15)$$

which resembles a standard gradient descent toward the target

$$T_t^{\text{DQN}} = R_{t+1} + \rho \max_a Q(S_{t+1}, a; \theta_t). \quad (3.16)$$

Even if tabular methods converge to the optimal function (Watkins and Dayan, 1992), they fail to generalize over previously unseen states. Instead, DRL has good generalization capabilities, but produces unstable behaviors during the training whenever function approximation is combined with an off-policy algorithm and learning by estimates (Sutton and Barto, 2018). The issue of training instability is partially solved by adding two ingredients: an experience replay buffer and a fixed target. An experience buffer is a finite set $\mathcal{D} = \{e_1, \dots, e_N\}$ of fixed cardinality N , where at each time t the agent's stream of experience $e_t = (S_t, A_t, R_{t+1}, S_{t+1})$ is stored, replacing one of the old ones. The replay buffer is then used to perform a batch update of the network parameters. A fixed target is exactly as the online target, except that its parameters θ^- are updated ($\theta_t^- = \theta_t$) and then kept fixed for τ iterations. Combining the two ingredients, the gradient step of Eq. (3.15)

becomes

$$\mathbb{E}_e[(r + \rho \max_{a'} Q(s', a'; \theta_t^-) - Q(s, a; \theta_t)) \nabla_{\theta_t} Q(s, a; \theta_t)] \quad (3.17)$$

where $e = (s, a, r, s')$ is uniformly sampled from \mathcal{D} .

In what follows we adopt a variant of the algorithm called double DQN (DDQN) (Van Hasselt et al., 2016) which prevents some overestimation issues of the value function. For convenience, we still refer to the chosen value based algorithm as DQN, even if the implementation follows the variant of DDQN. In the appendix, we recap the technical details of the value based algorithms used in the numerical experiments.

The optimization problem in Eq. (3.11) can be equivalently solved using a policy gradient algorithm like the Proximal Policy Optimization (PPO) (Schulman et al., 2017). A policy gradient algorithm directly parametrizes the optimal strategy within a given policy class $\pi_\theta = \pi(A_t | S_t; \theta)$, for example a multilayer neural network with the set of parameters θ . The optimization problem is approximately solved by computing the gradient of the performance measure $J(\theta) = \sum_{t=0}^{\infty} \rho^t R_{t+1}(S_t, A_t, S_{t+1}; \pi_\theta)$ and then carrying out gradient ascent updates according to

$$\theta_{t+1} = \theta_t + \alpha \nabla_{\theta} J(\theta_t), \quad (3.18)$$

where α is still a scalar learning rate. The policy gradient theorem (Marbach and Tsitsiklis, 2001; Sutton et al., 1999) provides an analytical expression for the gradient of $J(\theta)$ as

$$\begin{aligned} \nabla_{\theta} J(\theta) &= \mathbb{E}_{\pi_{\theta}} \left[\frac{\nabla_{\theta} \pi(A_t | S_t; \theta)}{\pi(A_t | S_t; \theta)} Q_{\pi_{\theta}}(S_t, A_t) \right] \\ &= \mathbb{E}_{\pi_{\theta}} [\nabla_{\theta} \log \pi(A_t | S_t; \theta) Q_{\pi_{\theta}}(S_t, A_t)], \end{aligned} \quad (3.19)$$

where the expectation, with respect to (S_t, A_t) , is taken along a trajectory (episode) that occurs adopting the policy π_θ . It can be proven that it is possible to modify the action value function $Q_\pi(s, a)$ in (3.19) by subtracting a baseline $V_\pi(s)$ that reduces the variance of the empirical average along the episode, while keeping the mean unchanged. A popular baseline choice is the state-value function

$$V_\pi(s) \equiv \mathbb{E} \left[\sum_{k=0}^{\infty} \rho^k R_{t+1+k} \mid S_t = s, \pi \right], \quad (3.20)$$

which reflects the long-term reward starting from the state s if the strategy π is adopted onwards. The gradient thus can be rewritten as

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\pi_{\theta}} [\nabla_{\theta} \log \pi (A_t | S_t; \theta_t) \mathbb{A}_{\pi_{\theta}}(S_t, A_t)] \quad (3.21)$$

where

$$\mathbb{A}_{\pi}(s, a) \equiv Q_{\pi}(s, a) - V_{\pi}(s), \quad (3.22)$$

is called advantage function and quantifies the gain obtained by choosing a specific action in a given state with respect to its average value for the policy π .

Different policy gradient algorithms depend on how the advantage function is estimated. In PPO, the advantage estimator $\mathbb{A}(s, a; \psi)$ is parametrized by another neural network with parameters ψ . This approach is known as actor-critic: the actor is represented by the policy estimator $\pi(a | s; \theta)$ that outputs the mean and the standard deviation of a Gaussian distribution which the agent uses to sample actions, the critic is the advantage function estimator $\mathbb{A}(s, a; \psi)$ whose output is a single scalar value. The two neural networks interact during the learning process: the critic drives the updates of the actor, which successively collects new sample sequences that will be used to update the critic and again evaluated by it for new updates. The PPO algorithm can therefore be described by the extended objective function

$$J^{\text{PPO}}(\theta, \psi) = J(\theta) - c_1 L^{\text{AF}}(\psi) + c_2 H(\pi(a | s; \theta)). \quad (3.23)$$

The second term is a loss between the advantage function estimator $\mathbb{A}(s, a; \psi)$ and a target \mathbb{A}^{targ} , represented by the cumulative sum of discounted reward, needed to train the critic neural network. The last term represents an entropy bonus to guarantee an adequate level of exploration. Details about the specific choice of the losses, the target and the neural network parametrization, together with additional information about the general algorithm implementation, are given in the appendix.

3.4 Numerical Experiments

In this section, we conduct synthetic experiments in the controlled financial environment outlined in Section Section 3.2. We present two different groups of experiments where the agents observe financial time series that come from different data generating processes. The first group is related to the case

where the return dynamics is driven by Gaussian mean reverting factors, as in Eq. (3.2) and the optimal strategy is known to be Eq. (3.5). The second group includes a set of cases where the model that generates the dynamics allows for a bigger amount of extreme events and heteroskedastic volatility. In this case, Eq. (3.5) is still used as a representative classical strategy of dynamic portfolio optimization.

In all experiments, the agents trade a single asset, but the framework is general enough to allow for multi asset trading. We test Q-learning and DQN in parallel in the same environment, while PPO training is slightly different. The first two algorithms are trained in-sample for a number of updates equal to the length T_{in} of the simulated series, while the learned policy is evaluated out-of-sample at several intermediate training moments on different series of length T_{out} . The same logic operates for PPO, which however works in an episodic way: the algorithm is trained in-sample and evaluated out-of-sample respectively for E_{in} and E_{out} number of episodes of length equal to 2000 time steps. Each agent operates in a model-free context so that no prior information about the data generating process is provided.

In order to bring the RL formalism to the portfolio optimization problem of Eq. (3.4) we choose the actions as the amount of shares traded $A_t = \Delta h_t$, while the state is defined as the pair return-holding $S_t = (y_t, h_{t-1})$. We include the asset return to the state representation instead of the predicting factors because it is our interest to assess DRL as a pure data-driven approach. The choice of financial factors is known to be a non-trivial task, and it can be highly discretionary. For every experiment, we also adapt the boundaries of the action space \mathcal{A} according to the magnitude of the action performed by the benchmark. More specific details about this heuristic are provided in the appendix.

After taking an action and causing a change in portfolio position, the agent observes the next price movement and the reward signal that is

$$R_{t+1}(y_{t+1}, h_{t-1}, \Delta h_t) = h_t^T y_{t+1} - \frac{\gamma}{2} h_t^T \Sigma h_t - \frac{1}{2} \Delta h_t^T \Lambda \Delta h_t. \quad (3.24)$$

Note that we decided to allow the benchmark agent to be perfectly informed, so that it knows exactly the predicting factors of the price dynamics. On the contrary, RL agents can just gather information from the observed return, which is affected by an additional source of noise. This choice allows the RL agent to be completely agnostic with respect to the price dynamics. This represents a clear disadvantage for the RL agent, but it is a step towards a

more flexible approach when the dynamics is not known and the performance is strongly dependent on the factors selection. Alternatively, one can assume that the RL agents are completely informed by replacing y_t with f_t in the definition of the state variables. For the purpose of comparison, we highlight that the value-based algorithms in this study perform discrete control, while the benchmark solution can adopt a continuous strategy according to Eq. (3.5). Although PPO can express both discrete and continuous policies, we test the continuous version to allow for a more expressive policy and compare the differences of the two settings.

The details about the parameters used to simulate the financial data and the hyperparameters setting for training the neural networks are provided in the appendix. The experiments are run in parallel on a 64-cores Linux server which has an Intel Xeon CPU E5-2683 v4 @ 2.10GHz. The training runtime for a single value-based method experiment of length $T_{in} = 300000$ ranges from two to four hours when the neural network architecture is not deeper than two hidden layers. Approximately the same runtime is needed for the PPO algorithm when $E_{in} = 300$. The source code written in Python is available on GitHub¹. The following subsections discuss the results of the two groups of experiments.

3.4.1 Tracking the Benchmark

From Figure 3.1 we observe the evolution of the out-of-sample performance of several Q and DQN agents in the case of a return dynamics driven by one or two Gaussian factors. After about half of the training runtime, DQN reaches on average a close-to-optimal cumulative net PnL , which is expressed as the gross return of the portfolio deducted from the transaction costs, i.e.

$$PnL_{t+1}^{\text{net}}(y_{t+1}, h_t, \Delta h_t) = h_t y_{t+1} - \frac{1}{2} \Delta h_t \Lambda \Delta h_t. \quad (3.25)$$

The trained DQN agents are then able to retrieve the mean reverting signals in the data and to control the amount of transaction costs without knowing the data generating process of the underlying dynamics. On the other hand, Q-learning agents hardly reach half of the cumulative net PnL of the optimal benchmark in the same training time.

¹<https://github.com/Alessiobrini/Deep-Reinforcement-Trading-with-Predictable>Returns>

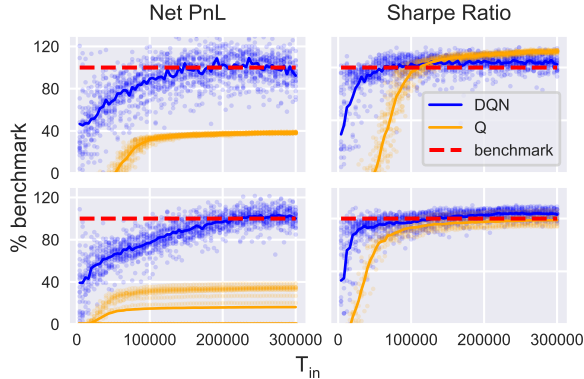


Figure 3.1: Results for DQN and Q-learning in the case of Gaussian dynamics driven by one (first row) and two (second row) mean-reverting factors. Cumulative net PnL (first column) and SR (second column) are displayed as the size of the training time series increases up to $T_{in} = 300000$ on the x-axis. Every dot represents the average over 10 out-of-sample tests of length $T_{out} = 5000$ for a specific agent out of the 20 tested in total. The horizontal dashed line represents the optimal benchmark, while the solid lines represent the average performance of all the agents in relative percentage to the benchmark.

The performances of the tabular algorithm are strictly dependent on the granularity of the state discretization. Q-learning can reach the benchmark performance only when $T_{in} \rightarrow \infty$ and \mathcal{S} is dense enough to closely represent the continuous trading environment. However, even if we keep the Q-table relatively small, usually below 100000, it can still be very sparse for this range of T_{in} . This is particularly evident in the case of two Gaussian factors, where many agents have a negligible cumulative net PnL simply because they do not perform any buy or sell actions. Increasing the size of the Q-table for experiments of fixed length T_{in} leads to even worse results.

DQN avoids the inefficient tabular parametrization of the action-value function by using fewer parameters with respect to the amount of entries in the Q-table. The use of a neural network as an action-value function approximator is crucial in this financial environment because the agent learns faster when the state space is entirely observable and the parameters can be updated by batches of experience.

In order to compare the risk of different strategies, we use the annualized

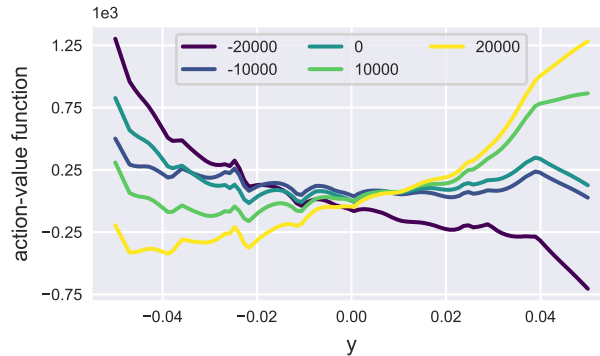


Figure 3.2: Learned action-value function for a DQN agent when the asset return varies and the holding is fixed at 0. Different colors represent different actions.

Sharpe ratio (SR) (Sharpe, 1994) which is computed as

$$SR = \frac{\mathbb{E}[PnL^{\text{net}}]}{\sqrt{\text{Var}[PnL^{\text{net}}]}} * \sqrt{252}, \quad (3.26)$$

defining the expected return of the portfolio per unit of risk on a yearly basis. It is a common metric to evaluate the trade-off between risk and return of financial strategies, especially in a mean-variance optimization framework.

The second column of Figure 3.1 showcases the evolution of the SR of the agents and highlights that DQN obtains on average the same level of benchmark profit adjusted for risk since the beginning of the training. The performances of Q-learning in this case are strongly biased, since often the tabular agents choose not to trade and avoid increasing the risk of its portfolio position. The DQN agents first learn how to obtain low-risk portfolios, then they start making higher profits, as it is shown from the faster convergence of the SR with respect to the cumulative net PnL .

Figure 3.2 provides insights about the learned behavior of the DRL agents, showing the learned action-value function of the best performing DQN agent at the end of the period of training represented in Figure 3.1. The agent is trained over a return dynamics driven by one predicting factor, but the findings are valid also in the case of multiple factors. The estimated action-value function $Q((y, h), a; \theta)$ is displayed for all the actions in the discrete space \mathcal{A} and implicitly represents the behavior of the agent when different

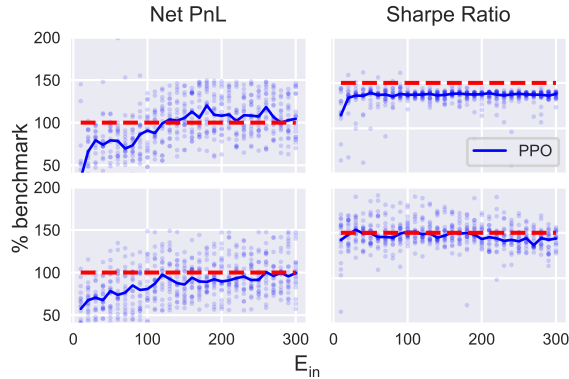


Figure 3.3: Results for PPO in the case of Gaussian dynamics driven by one (first row) and two (second row) mean-reverting factors. Cumulative net PnL (first column) and SR (second column) are displayed as the training episodes increases up to $E_{in} = 300$ on the x-axis. Every dot represents the average over 10 out-of-sample tests of length $T_{out} = 2000$ for a specific agent out of the 20 tested in total. The horizontal dashed line represents the optimal benchmark, while the solid line represents the average performance of all the agents in relative percentage to the benchmark.

levels of returns are experienced. If the agent acts greedily and chooses the highest Q-value for every level of y , positive actions appear prevalent when returns are positive, while the opposite holds for negative actions.

Figure 3.3 shows analogous experiment results for PPO trained with financial returns driven by mean-reverting Gaussian dynamics. PPO retrieves the signal in the data and converges to the benchmark, but exhibits higher variance in the Net PnL measure compared to DQN. This is motivated by the different type of policy that the algorithm is describing. Working in a continuous action space allows the possibility to trade any fraction of the synthetic asset, but also complicates the exact convergence to the benchmark because the sampling space is large. In practice, there is no theoretical guarantee to find the proper way to sample actions from \mathcal{A} in a finite time.

Figure 3.4 represents the average greedy policy learned by the agents for DQN and PPO presented in Figure 3.1 and Figure 3.3. Both algorithms can discover the inherent arbitrage introduced in the market, since the average policy follows the sign of the returns by buying low and selling high. The learned policies appear to be monotonic as the one of the benchmark.

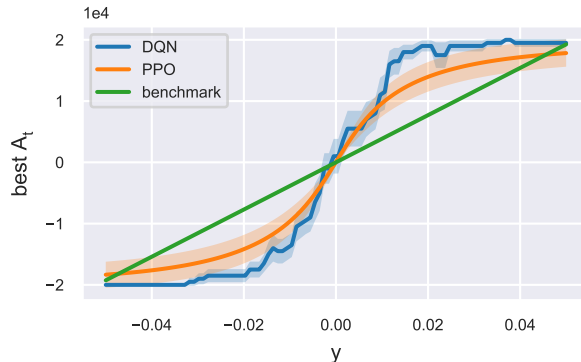


Figure 3.4: Greedy policy function for a DQN and a PPO agent together with the policy of the benchmark when the asset return varies and the holding is fixed at 0. Different colors represent different algorithms. The curves present a confidence interval because the average maximum action is obtained from the results of 20 trained agents for each algorithm.

3.4.2 Outperforming the Benchmark

In order to show the flexibility of the DRL approach, we study its performances with respect to the benchmark when the return dynamics depart from the original model specification. In particular, we introduce two types of model misspecifications: the presence of extreme events and noise heteroscedasticity. In both cases the strategy in Eq. (3.5) is no more optimal, but since it performs well, and it is often used in practice, it can be considered as a benchmark representing a broader class of factor trading strategies. It is therefore natural to investigate whether DQN and PPO are able to outperform the benchmark other than just reaching it.

We consider two different reference strategies that are respectively referred to as fully informed when the benchmark is provided with the simulated factors and partially informed when instead it needs to extract them from the observed returns. These different settings should not affect the DRL performance, except for the boundaries of the action space \mathcal{A} that we decided to adapt to the benchmark for a better comparison (see the appendix).

The fully informed benchmark agent can directly use Eq. (3.5) to trade, just by estimating the speeds of mean reversion and factor loadings from the observed return predicting factors. Instead, in the partially informed case, the benchmark agent does not know exactly which are the best predicting

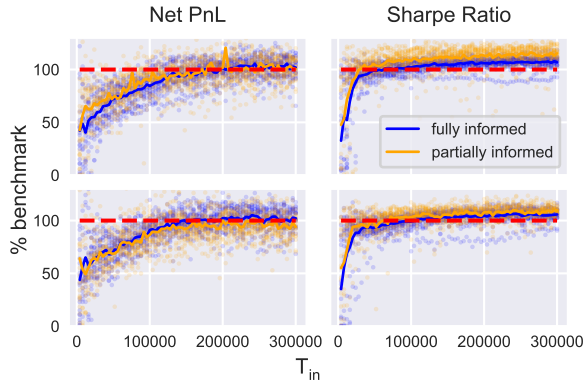


Figure 3.5: Results for Student’s T dynamics in the case of one mean-reverting factor. DQN is tested over Student’s T distributed returns with $\nu = 6$ (first row) and $\nu = 8$ (second row) respectively. The figure should be read with the same logic of Figure 3.1 since the number of agents and the length of in-sample and out-of-sample experiments are equal. The solid lines represent different PPO performances with respect to different benchmark strategies, which in this context are not optimal anymore.

factors and needs to guess or extract them from what is observed in the state space. One of the typical choices in financial literature is to use lagged past returns (Asness et al., 2013) as factors to predict future returns. We resort to a simple heuristic to select the best possible lagged variables by fitting the Eq. (3.1) to a set of candidate lags. Then, we select the best one by minimizing the average squared residuals.

From Figure 3.5 we note that in presence of extreme events (T-student noise) the DQN agents are able to control the trading costs and obtain equal or superior cumulative net PnL with respect to the two benchmark agents, especially for lower degrees of freedom where the misspecification has a greater impact and extreme events are more frequent. The SR fairly outperforms the benchmark towards the end of the training process in both the cases. The RL agents learn how to control the higher amount of risk introduced in the environment, while model based strategies like the benchmark should have considered this in advance. Figure 3.6 shows the same misspecified case for PPO agent, which is able to consistently manage the transaction costs and obtain higher net PnL than the benchmark, still exhibiting greater variance than DQN. We note that PPO performs better with respect to the

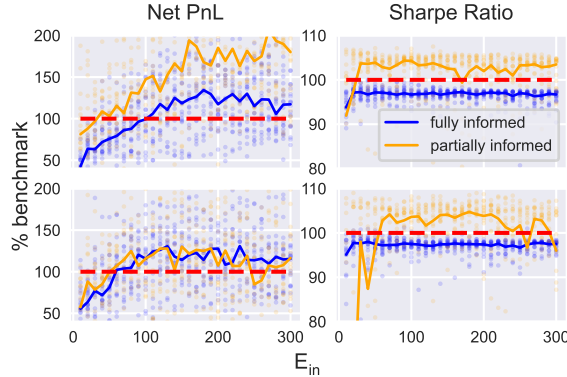


Figure 3.6: Results for Student’s T dynamics in the case of one mean-reverting factor. PPO is tested over Student’s T distributed returns with $\nu = 6$ (first row) and $\nu = 8$ (second row) respectively. The figure should be read with the same logic of Figure 3.3 since the number of agents and the length of in-sample and out-of-sample experiments are equal. The solid lines represent different PPO performances with respect to different benchmark strategies, which in this context are not optimal anymore.

benchmark when the latter is provided with partial information and needs to discover the persistence of the signal on its own.

The second proposed misspecification introduces heteroscedasticity in the asset returns by considering a GARCH process with $p = 1$ and $q = 1$ for the asset variance. For simplicity, we assume the predictable component of the returns to be an autoregressive model with lag of order 1.

Figure 3.7 showcases that DQN obtains on average a higher cumulative net PnL with respect to the benchmark. We compare the difference, instead of the ratio, between the two cumulative net PnL s because in some cases the net PnL obtained by the benchmark agent is negative. Differently from the previous set of experiments, the increment of performance in the presence of heteroscedasticity regards the control of the amount of transaction cost. Looking at the SR , DQN mostly tracks the benchmark and outperforms it only in the case of Gaussian noises. The increased amount of extreme events in the Student’s T case causes a worsening in the DQN performance relative to the benchmark. It has to be noted that we use the same set of hyper-parameters for all these experiments. This is a signal that the performance in the fat-tailed case could be improved by tuning a more effective config-

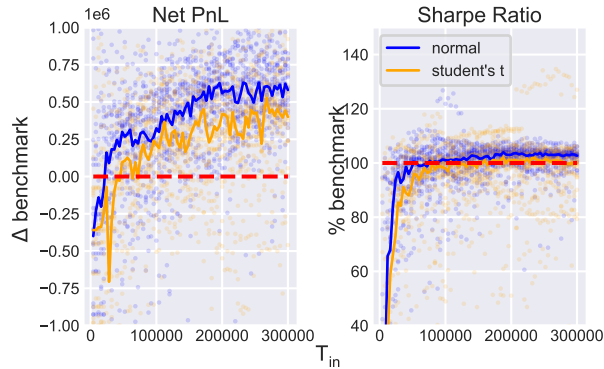


Figure 3.7: Results for AR-GARCH return dynamics when the noise is distributed as a standard Normal or a Student’s T distribution ($\nu = 8$). The performance of DQN is compared with the dashed line benchmark. On the y-axis of the left plot, there is the cumulative net PnL difference between DQN and the benchmark, so that the dashed line is set at 0. Then, the figure should be read with the same logic of Figure 3.1 since the number of agents and the length of in-sample and out-of-sample tests are equal.

uration. Figure 3.8 further confirms that PPO deals more effectively with model misspecifications with respect to DQN. When trained over GARCH dynamics, PPO achieves better than the benchmark performance, controlling associated risks and costs.

Figure 3.9 shows the realized out-of-sample holdings for some DQN agents. When the underlying dynamics can be predicted by mean reverting factors, as for the Gaussian and Student’s T cases, the inversion of the factor sign causes the inversion of the sign of the portfolio itself. The oscillation between short and long positions confirms that the DRL algorithm has learned to follow the signal present in the data. In particular, when compared with a partially informed benchmark agent, as in the bottom left plot of Figure 3.9, the DQN algorithm obtains a higher cumulative net PnL by anticipating the reverting movement of the returns. Figure 3.10 presents out-of-sample holdings for PPO in the same cases already shown for DQN. When the returns are driven by a Gaussian or a Student’s T dynamics, PPO tracks the portfolio benchmark well and in the partially informed case seems to anticipate the mean-reversion of the signal as DQN does. In the case of GARCH dynamics, both the algorithmic approaches show a portfolio holding which

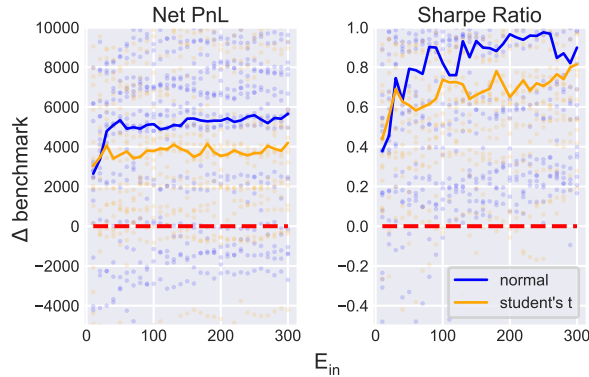


Figure 3.8: Results for AR-GARCH return dynamics when the noise is distributed as a standard Normal or a Student’s T distribution ($\nu = 8$). The performance of PPO is compared with the dashed line benchmark. On the y -axis of both sub figures there is the difference between PPO and the benchmark respectively for cumulative net PnL and the Sharpe ratio. Then, the figure should be read with the same logic of Figure 3.3 since the number of agents and the length of in-sample and out-of-sample tests are equal.

greatly differs from the benchmark. In both figures, we have two different y -axes for the bottom right panel. In order to visualize portfolio holdings that are on different scales, we let the left y -axis be associated with the algorithm and the right y -axis be associated with the benchmark. The latter does not adapt well to heteroskedastic peaks in the time-series of simulated returns, and it happens that the benchmark portfolio is more sensitive to extreme events. On the other hand, RL is able to limit the trading even in the presence of heteroscedasticity and obtain a lower portfolio size that produces less transaction costs when it needs to be rebalanced.

It is important to check the robustness of the RL performance with respect to the choice of the dynamics parameters, or conversely its sensitivity with respect to some of them. Figure 3.11 and Figure 3.12 show the level of SR for DQN and PPO agent as a function of dynamics half life, factor loading and fat-tailed return distribution parameters. The result is an average over 10 different agents for each parameter configuration, this also allows creating confidence intervals around the average performance. In the top row of both figures, the effect of variation in half-life of mean-reversion and in the factor loading b of a one-factor Gaussian dynamics is similar for DQN and PPO.

3. Deep Reinforcement Trading with Predictable Returns

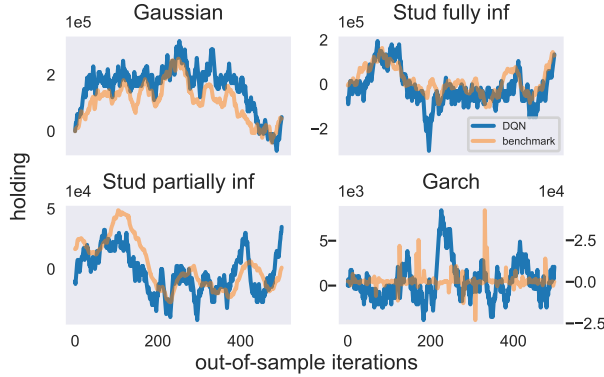


Figure 3.9: Portfolio Holdings for a snapshot of length 500 for some out-of-sample tests performed. The DQN agents selected are the best performing for each group in terms of cumulative net PnL : the Gaussian case with one factor, the Student’s T case with 6 degrees of freedom (fully and partially informed) and the GARCH(1,1) with normal noise.

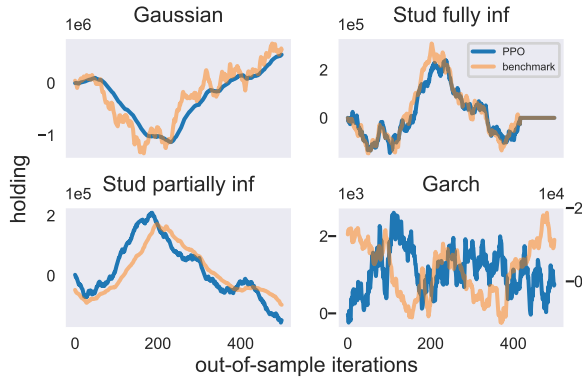


Figure 3.10: Portfolio Holdings for a snapshot of length 500 for some out-of-sample tests performed. The PPO agents selected are the best performing for each group in terms of cumulative net PnL : the Gaussian case with one factor, the Student’s T case with 6 degrees of freedom (fully and partially informed) and the GARCH(1,1) with normal noise.

They both tend to obtain the same SR of the benchmark strategy, which is in both cases an increasing function. This is due by one side from the fact that a higher half life produces a more persistent return sign that agents can

more easily exploit to make profit. From the other side, when the factor loading is small, then the return dynamics contains no meaningful signal and the Eq. (3.1) is driven purely by noise. The left panel in the bottom row of Figure 3.11 and Figure 3.12 proposes the same sensitivity analysis when we simulate one factor Student's T dynamics with increasing degrees of freedom ν . All the strategies get worse as the percentage of extreme events increases, but both PPO and DQN deal with riskier events effectively and consistently achieve a greater SR than the benchmark strategy. The right panel in the bottom row of both figures presents the variation of the performance when the kurtosis of the GARCH(1,1) process distribution increases. The fourth standardized moment for the stochastic volatility process in Eq. (3.10) is computed as

$$\frac{E(\sigma_t^4)}{[E(\sigma_t^2)]^2} = \frac{3[1 - (\alpha_1 + \beta_1)^2]}{1 - (\alpha_1 + \beta_1)^2 - 2\alpha_1^2} \quad (3.27)$$

knowing that when $1 - 2\alpha_1^2 - (\alpha_1 + \beta_1)^2 > 0$, the tails in the distribution of the GARCH(1,1) are heavier than a Gaussian. A GARCH(1,1) model with heavy tails represents a consistent misspecification of the original conditions and only PPO consistently outperforms the benchmark while the tails of the simulated return distribution become thicker.

3. Deep Reinforcement Trading with Predictable Returns

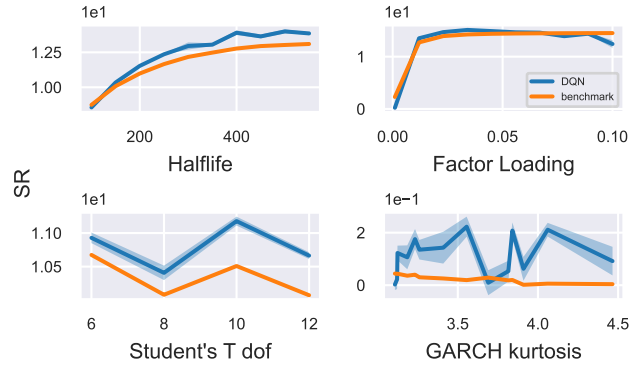


Figure 3.11: DQN performances measured by percentage or differential SR with respect to the benchmark when some relevant parameter of the simulated dynamics varies. The top row shows the one Gaussian mean-reverting factor dynamics when the level of half-life of mean-reversion and factor loadings increases. The bottom row shows the one mean-reverting Student's T factor dynamics and the GARCH dynamics when respectively the degrees of freedom (dof) and the kurtosis of the returns distribution increase.

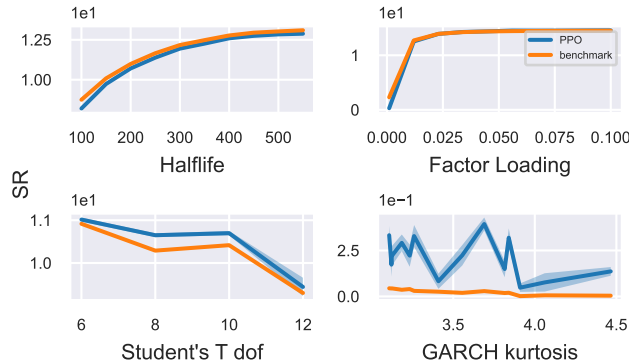


Figure 3.12: PPO performances measured by percentage or differential SR with respect to the benchmark when some relevant parameter of the simulated dynamics varies. The top row shows the one Gaussian mean-reverting factor dynamics when the level of half-life of mean-reversion and factor loadings increases. The bottom row shows the one mean-reverting Student's T factor dynamics and the GARCH dynamics when respectively the degrees of freedom (dof) and the kurtosis of the returns distribution increase.

3.5 Conclusions

In this work, we have used different RL algorithms to solve a trading problem in a financial environment where trading is costly. When the optimization problem is known to have an exact solution, DQN and PPO are able to track this benchmark, but they can also adapt to variations of the original environment setting and find a way to control portfolio risks and costs in a data-driven manner. While value-based DRL results to be accurate in following the trading signals and controlling the market frictions, policy-based DRL appears to be more robust to extreme events and heteroskedastic volatility.

Although DQN is able to learn the direction of the trades, the discretization of the action space still represents a major concern because the traded size is a multiple of a fixed quantity chosen in advance. Instead of moving towards actor-critic frameworks, which usually shares the instability issues with DQN, PPO helps in solving this issue by expressing the policy in a continuous action space.

RL algorithms are demanding in terms of training data that can be quite scarce, especially at low frequencies. We believe that the use of a financial model with a known optimal solution can offer a workaround to this problem by allowing to pretrain DRL agents on synthetic data and then to fine tune them on real time series. Moreover, classical strategies can facilitate the training of DRL agents by providing information about good (even if suboptimal) strategies. This can help in terms of a rationalization of the state-action space, which results in the lightening of the training process itself.

3.6 Appendix : Algorithms and Hyperparameters

In this appendix, we provide some details regarding the implementation of Q-learning and DQN algorithms used in the numerical experiments. Then we outline the choice of the parameters for simulating the financial data and of the relevant hyperparameters to set up the training of the algorithms.

3.6.1 Q-learning

Q-learning requires the discretization of the state space \mathcal{S} and the action space \mathcal{A} , which affects the dimensionality of the Q-table that produces estimates $Q(S_t, A_t)$ of the optimal action-value function. For every possible state and action variable, we need to choose a proper discrete range that we believe is adequately large to capture the relevant information and solve the problem.

Since real traders usually operate by trading quantities of assets that are multiples of a fixed size called lot, the dimensions of the Q-table are bounded by setting the traded quantity Δh_t to be at most K round lots and the portfolio holding h_t to a maximum of M round lots. The discrete set of returns is represented by an upper and lower bounded set of values that are linearly spaced by the size of a basis point, denoted as bp . The bounded sets and their dimensionality are respectively:

$$\mathcal{A} = \{-K, -K + 1, \dots, K\}, \quad |\mathcal{A}| = 2K + 1 \quad (3.28)$$

$$\mathcal{H} = \{-M, -M + 1, \dots, M\}, \quad |\mathcal{H}| = 2M + 1 \quad (3.29)$$

$$\mathcal{R} = bp \cdot \{-T, -T + 1, \dots, T\}, \quad |\mathcal{R}| = 2T + 1 \quad (3.30)$$

In our financial environment, the basis point and the lot are respectively the size of minimum return movement and the minimum tradable quantity of the asset at each discrete time. The sizes of the three sets are defined respectively by the hyperparameters K , M and T , which are a crucial choice to define the magnitude of the synthetic financial problem.

Denoting the size of the table as $d = |\mathcal{R}| \times |\mathcal{H}| \times |\mathcal{A}|$, our simulated experiments show that Q-learning is not even able to reach a positive profit when d approaches the length of the simulated series T_{in} . The more the dimensionality of the table increases, the worse are the cumulative net PnLs and rewards obtained, when the T_{in} is fixed.

If T_{in} is not sufficiently long to allow the agent to visit the entire state space and update the Q-table in each corresponding entry, the algorithm approximates the action-value function with a sparse Q-table. Thus, it is not able to represent the effect of slight changes in the state space variables. Such a bottleneck becomes even worse if we increase the number of actions that the agent can perform.

In principle, we could let Q-learning experience longer simulated series to

partially avoid the exploration issue, but this would not be of any practical use for two reasons: (i) it requires an increasingly long training runtime to match the benchmark performance and still this result would be obtained under a discretized state-space of a more complex financial environment; (ii) training for a high number of iteration the experiment would not even resemble a real financial application since there is no way to retrieve such a massive amount of financial data, especially at a daily frequency.

To ensure proper exploration of the state space, the agent acts according to a ϵ -greedy policy, such that at each time a greedy action $a = \arg \max_a Q(S_t, a)$ is selected with probability $1 - \epsilon$, while occasionally with probability ϵ a random action is sampled from the set \mathcal{A} . As a common approach in the RL literature, the value of ϵ decays linearly during training until it reaches a small value that is kept fixed until the end.

3.6.2 DQN

DQN requires a discretization of the action space \mathcal{A} , which is approached as for the tabular case. This discretization could represent an issue when one wants to represent the choice of the agent at a more granular level. The more one increases the size of \mathcal{A} , the more the computational cost of the algorithm increases and its efficiency in solving the financial problem decreases. However, we believe the discrete control can still be adequate for a set of financial problems since usually market orders are executed in multiples of a fixed quantity.

Since the agent learns offline by choosing past batches of experience from a buffer with fixed size, we set this dimension as a percentage of the total updates in-sample T_{in} . We have found that letting the buffer size to increase can improve the performance, therefore we do not discard any sequence. The exploration-exploitation trade-off is balanced as in Q-learning, using a ϵ -greedy policy where the ϵ decreases linearly to a low value towards the end of the training. Despite the original DQN implementation (Mnih et al., 2015) suggests updating the target network parameters at every fixed discrete step, we choose to continuously update the target parameter so that they slowly track the learned networks as follows:

$$\theta_t^- \leftarrow \tau \theta_t^- + (1 - \tau) \theta_t \tag{3.31}$$

where τ is the chosen step size for the update, θ_t^- are the target network

parameters and θ are the parameters for the current action-value function estimator.

A problem of the overestimation of the action-value function is known to arise in the classical DQN algorithm (Hasselt, 2010; Van Hasselt et al., 2016). Thus, we adopt the double DQN (DDQN) variant suggested in Van Hasselt et al., 2016. Recalling that the target of a DQN update is computed as

$$T_t^{\text{DQN}} = R_{t+1} + \rho \max_a Q(S_{t+1}, a; \theta), \quad (3.32)$$

where we can write

$$\max_{a'} Q(s', a'; \theta) = Q\left(s', \arg \max_{a'} Q(s', a'; \theta); \theta\right),$$

it happens that the same noise affects both the maximization over the action space and the value function estimates. Removing the correlation between the sources of noise coming into these two operations is beneficial to avoid an overestimation of the value function.

Double Q-learning decouples the selection of the action from the evaluation as,

$$T_t^{\text{DDQN}} = R_{t+1} + \rho Q(S_{t+1}, \arg \max_{a'} Q(S_{t+1}, a'; \theta_1); \theta_2), \quad (3.33)$$

i.e., DDQN uses two neural networks: one computes the target and the other computes the current action-value function. The computation of the target is split between the current neural network that greedily selects the action and the target neural network that evaluates such action. Therefore, the selection of the action in Eq. (3.33) is due to the current weights $\theta_1 = \theta_t$, while the target network is used to evaluate the value function for that action $\theta_2 = \theta_t^-$.

We have found that DDQN outperforms DQN in all the tests we carried out, meaning that also in this specific financial application it is a profitable procedure.

Regarding the shape of the loss function, a common choice for the DQN family of algorithms is the Huber loss rather than the mean squared error (MSE), which is typical for regression tasks. The Huber loss is less sensitive to the presence of outliers, and it is expressed as

$$\mathcal{L}_\delta(y, \hat{y}) = \begin{cases} \frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i)^2 & \text{for } |y_i - \hat{y}_i| \leq \delta \\ \delta \frac{1}{N} \sum_{i=1}^N |y_i - \hat{y}_i| - \frac{1}{2} \delta^2 & \text{otherwise} \end{cases} \quad (3.34)$$

The Huber loss is quadratic for small values of the squared difference and linear for larger values. This kind of loss function adds a lower penalty to large errors, and it is better than MSE for this kind of problems because learning by estimates as in DQN could produce unexpectedly high errors. Even if in the main we showed the case of an MSE loss, in practice our implementation utilizes a Huber loss. This choice does not change the update rule of the presented algorithms, but the computation of the gradient will differ in the presence of large quadratic errors.

The neural networks used to approximate the action-value function are 2-layer fully connected networks with ELU (Clevert et al., 2015) activation and uniform weight initialization as in He et al., 2015. We have tried with different types of rectified nonlinear activation, but ELU outperforms a more usual choice as ReLU. The sizes of the hidden layers are 256 and 128 for the first and the second respectively, but also smaller hidden layers have proven to be effective (Xu et al., 2015).

The gradient descent optimizer is Adam (Kingma and Ba, 2014), which performs a batch update of size 256. The original implementation proposes default values for β_1 , β_2 , and ϵ^{adam} , which are respectively the exponential decay rates for the first and the second moment estimates of the true gradient and a small constant for numerical stability. Those parameters required some tuning for improving performances, so that we set them as $\beta_1 = 0.5$, $\beta_2 = 0.75$ and $\epsilon^{\text{adam}} = 0.1$ for all experiments. The learning rate α usually starts around 0.005 and then decays exponentially towards the end of the training.

Since in a RL setting the data are not all available at the beginning of the training, we can not normalize our input variables as usual in the preprocessing step of a supervised learning context. Hence, we add a Batch Normalization layer (Ioffe and Szegedy, 2015) before the first hidden layer to normalize the inputs batch by batch and obtain the same effect.

3.7 PPO

PPO allows expression continuous policies through an algorithm which is easier to implement than a trust-region method (Schulman et al., 2015a) and easier to tune with respect to the continuous counterpart of DQN (Lillicrap et al., 2015). In principle, continuous policies are more expressive than discrete policies, but are also harder to learn. Our implementation of PPO follows Andrychowicz et al., 2020 which performs a large empirical study of the effect

of implementation and parameters choices on the PPO performances. Even if our financial problem is different from their test bed, we also follow the direction of their results in order to tune our hyperparameters, since we have limited computational resources to do this search from scratch. Another relevant source for an effective implementation is Engstrom et al., 2020.

As described in the main, we implement PPO in an actor-critic setting without shared architectures. Differently from the standard implementation, the actor outputs a single scalar value, which is the mean of a Gaussian distribution, while the standard deviation is then learned as a global parameter in the optimization process and updated using the same gradient optimizer. Learning a global standard deviation for all the state representation has proven to be as much as effective as learning a state dependent parameter, with the benefit of being slightly less computational expensive (Andrychowicz et al., 2020). Policy gradient methods like PPO allow inserting some prior knowledge on the form of the policy with respect to value-based methods. The real-valued output of the actor usually passes through a hyperbolic tangent function in order to bound the action in the interval $[-1, 1]$. Then is rescaled to directly express a range of possible trading actions, whose extreme values are selected according to a heuristic described in the next subsection. Exploration during training is guaranteed by the learned standard deviation parameter and by the entropy bonus in the objective function. When doing out-of-sample tests, the PPO policy is tested as if it were deterministic by just picking the mean of the Gaussian instead of sampling from it. We take this choice because we do not want the test results to be affected by stochasticity.

The on-policy feature of PPO makes the training process episodic, so that experience is collected by interacting with the environment and then discarded immediately once the policy has been updated. The on-policy learning appears in principle a more obvious setup for learning, even if it comes with some caveats because it makes the training less sample efficient and more computationally expensive since a new sequence of experiences need to be collected after each update step. In this process, the advantage function is computed before performing the optimization steps, when the discounted sum of returns over the episode can be computed. In order to increase the training efficiency, after one sweep through the collected samples, we compute again the advantage estimator and perform another sweep through the same experience. This trick reduces the computational expense of recollecting experiences and increases the sample efficiency of the training

process. Usually we do at most 3 sweeps (epochs) over a set of collected experiences before moving on and collecting a new set.

The optimizer and the normalization of the inputs through a Batch Normalization layer are the same used for DQN, with the only exception that in the PPO case we do not tune the hyperparameters of the Adam optimizer.

Maximizing the objective function that returns the gradient in Eq. (3.21) is known to be unstable, since updates are not bounded and can move the policy too far from the local optimum. Similarly to TRPO (Schulman et al., 2015a), PPO optimizes an alternative objective to mitigate the instability,

$$J^{\text{CLIP}}(\theta, \psi) = \mathbb{E}_{\pi_{\theta}} \left[\min \left(r(\theta) \hat{\mathbb{A}}(s, a; \psi), \text{clip} \left(r(\theta), 1 - \epsilon, 1 + \epsilon \right) \hat{\mathbb{A}}(s, a; \psi) \right) \right] \quad (3.35)$$

where $r(\theta) = \frac{\pi(A_t|S_t;\theta)}{\pi(A_t|S_t;\theta_{\text{old}})}$ is a ratio indicating the relative probability of an action under the current policy with respect to the old one. Instead of introducing a hard constraint as in TRPO, the ratio is bounded according to a tolerance level ϵ to limit the magnitude of the updates. The combined objective function in Eq. (3.23) can be easily optimized by PyTorch’s automatic differentiation engine, which quickly computes the gradients with respect to the two sets of parameters θ and ψ . The implemented advantage estimator depends on the parametrized value function V_{ψ} and is a truncated version of the one introduced by Mnih et al., 2016 for a rollout trajectory (episode) of length T :

$$\hat{\mathbb{A}}_t = \delta_t + (\gamma\tau)\delta_{t+1} + \dots + \dots + (\gamma\tau)^{T-t+1}\delta_{T-1} \quad (3.36)$$

where $\delta_t = r_t + \gamma V_{\psi}(s_{t+1}) - V_{\psi}(s_t)$, γ is a discount rate with the same role as ρ in DQN and τ is the exponential weight discount which controls the bias variance trade-off in the advantage estimation. The generalized advantage estimator (GAE) uses a discounted sum of temporal difference residuals, similarly to the one-step target value of DQN in Eq. (3.16).

3.7.1 Environment choices

In all the simulated experiments we set $|\mathcal{A}| = 5$ for both Q-learning and DQN, so that every agent can perform two buy actions, two sell actions and a zero action. In order to make the results comparable with those of the benchmark solution, we adopt a systematic way to choose the size of the action space \mathcal{A} from which we obtain the possible actions. Basically, we let the dynamic programming solution run for some iterations before starting

the training of the RL algorithms, and we look at the distribution of the continuous action performed. Then we select the lower and upper boundary of \mathcal{A} as respectively the quantiles located at the 0.1% and the 99.9

For what concerns the simulated environment, the cost multiplier λ is chosen equal to 0.001 for all experiments. Then we assume zero discount rate and all the starting positions for the holding are zero. For any experiment involving mean-reverting factors, we compute the speed of mean reversion ϕ as $\phi = \frac{\log(2)}{\log(h)}$ where h is known as the half-life of mean-reversion and represents the time it is expected to take for half of the trading signal to disappear. This allows to simulate the predicting factors and aggregate their effect to compute the asset returns. We tried many set up for the half-lives, factor loadings and volatility of simulated assets and the findings are quite robust. The half-lives of the mean-reverting factors are 350 for the case of a single factor and (170,350) for the case of two factors. Factor loadings are chosen of the order of magnitude of 10^{-3} , specifically as 0.00535 and 0.005775 in the proposed cases. The volatilities of the factor are respectively 0.2 and 0.1, while the volatility of the unpredictable part of the asset return is always set to 0.01. Suitable ranges for these hyperparameters are $[0.5, 0.05]$ for the former and $[0.05, 0.005]$ for the latter. In general, DQN is also able to retrieve the underlying dynamics in the case of two concurrent factors with different speeds, as long as those factors do not include one which is really fast (e.g. half-life of mean reversion lower than 10 days) and also highly noisy with a volatility above 0.2. This is acceptable because the signal-to-noise ratio would be very low, and it may require more sophisticated layers for feature extraction, other than a feedforward network structure. The parameters for the AR-GARCH simulation are $\omega = 0.01$, $\alpha_1 = 0.05$ and $\beta_1 = 0.94$ which are common GARCH parameters to simulate a stable financial market. The autoregressive parameter is set to $\phi_{Lr} = 0.9$ and the degrees of freedom in the Student's T case are $\nu = 10$.

Chapter 4

Reinforcement learning for trading in a market with frictions: a residual approach

4.1 Introduction

Portfolio optimization represents selecting the best allocation of capital in a given universe of financial assets out of all the available asset combinations while maximizing a utility objective that reflects the investor preferences. It has been a well-known and studied problem in the field of finance since the introduction of the modern portfolio theory in the seminal work of Markowitz, 1952 that first presented the concept of diversification in a portfolio of assets by trading-off risk and return. From there, financial literature has proposed several approaches to find an efficient way to optimize a financial portfolio, using advanced mathematical techniques and accounting for realistic aspects of the financial markets. Kolm et al., 2014 provides a comprehensive survey of the development of portfolio optimization techniques where they present and discuss the significant drawbacks of classical Markowitz's framework. One main point is overcoming the single-period optimization (SPO) setting, which performs poorly when the trader incurs transaction costs. A proper framework that accounts for the costly portfolio rebalancing is represented by a multi-period optimization (MPO) framework, as proposed by Merton, 1969 where the investors seek to maximize the total utility over time. Boyd et al., 2017 provides a detailed definition of the framework under the MPO

4. Reinforcement learning for trading in a market with frictions: a residual approach

setting.

MPO portfolio problems have been generally solved by quadratic programming under the mean-variance framework (Çelikyurt and Özekici, 2007; Mei et al., 2016), stochastic programming (Gondzio and Grothey, 2007; Samuelson, 1969) and dynamic programming (Cai et al., 2013; Gârleanu and Pedersen, 2013; Wang et al., 2015). Beyond these common approaches, reinforcement learning (RL) (Sutton and Barto, 2018; Szepesvári, 2010) is a convenient framework for financial MPO problems with a mathematical formalization that resembles the optimal control theory and dynamic programming (Bertsekas, 2005). RL obtained benefits from the last decade’s improvement of deep learning (DL) (Goodfellow et al., 2016) as a field of research that use neural networks as powerful function approximators (Cybenko, 1989). The recent technological advancement favored the development of sophisticated algorithms to solve high-dimensional problems, among which deep reinforcement learning (DRL) that became an active field of research in many domains of applications, such as gaming (Schrittwieser et al., 2020; Silver et al., 2017b) and robotics (Levine et al., 2016).

The model-free setting is generally adopted in the financial literature to solve financial portfolio problems through RL. Sato, 2019 provides a broad overview of the literature at the intersection between finance and RL. Some works are worth mentioning to understand the current state of the financial RL literature. Du et al., 2016 optimizes a portfolio with a risky and a riskless asset using a policy search method, whereas Jin and El-Saawy, 2016 uses the Deep-Q-Network algorithms to select a portfolio between high and low beta companies. One of the most-cited applications of RL in finance is Jiang et al., 2017, which proposes a policy-based framework with an ensemble of neural networks to discover the optimal way of trading a portfolio of cryptocurrencies. Liang et al., 2018 use the same framework to build portfolios in the equity space. Similar works are also Bertoluzzo and Corazza, 2012; Gao et al., 2020; Hu and Lin, 2019; Park et al., 2020; Wang et al., 2019; Xiong et al., 2018; Zhang et al., 2020.

The major drawback of the model-free approach for financial reinforcement learning is the training instability due to the difficulty of devising the underlying dynamics of the problem without knowing its structure. As already tackled on a single asset level by Brini and Tantari, 2021 and Chaouki et al., 2020, synthetic experiments in a controlled environment can help in understanding the effectiveness of reinforcement learning algorithm when it comes to trading a single or a portfolio of assets. Disentangling the problem’s

4. Reinforcement learning for trading in a market with frictions: a residual approach

outcome from the quality of the data and performing controlled experiments using a known closed-form solution represent the first step to understanding the proper use of model-free RL methods. However, model-free methods often fail at learning the optimal solution because of their training instability, partially due to the insufficient information provided to the RL algorithm, which cannot devise the underlying dynamics of the problem. Another known issue is the algorithmic approximation performed by the deep RL algorithms, which exacerbates when the problem’s dimensionality increases by solving multi-asset portfolio optimization problems.

At the other end of the spectrum of RL techniques, model-based RL represents the class of algorithms that learn an optimal policy when provided with a model of the environment. Hence, differently from model-free learning, they exploit prior knowledge of the environment or the problem to help the agent succeed faster. Han et al., 2019; Wei et al., 2019; Yu et al., 2019 represent some of the attempts to learn profitable trading actions through model-based RL. However, the literature lacks approaches that can guarantee the flexibility of a model-free setting while still converging faster as in a model-based environment.

Our work aims to investigate how to adapt a model-free reinforcement learning approach to leverage domain knowledge of the problem without imposing restricting assumptions on the policy that the agent should learn by providing a complete model of the environment. Since the novelty of the model-free RL approach rather than the optimal control literature is avoiding the modeling of the underlying dynamics of the system, we want to maintain our approach as close as possible to this setting without adopting a comprehensive model-based approach. Gruenstein et al., 2021; Johannink et al., 2019; Silver et al., 2018 study reinforcement learning approach for robots and controllers for which it is difficult to obtain a proper model, and they propose a way to leverage approximate models by learning only partially the underlying dynamics of the environments. This approach, called residual RL, helps in reducing the sample complexity associated with learning an accurate robot model and leads to faster convergence. In this work, we similarly tackle the MPO problem by leveraging available solutions to the portfolio trading problem to learn how to deviate from it and reach the optima. We propose a domain-specific approach for our selected policy learning algorithm, and we also provide some other solutions at the end of the analysis as topics for further research.

The paper is structured as follows: the Section 4.2 presents the classical

multiperiod framework for portfolio optimization common to RL and optimal control problems. In Section 4.3 we introduce the financial model used to simulate alpha return. The Section 4.4 presents the numerical experiments conducted to prove the effectiveness of the residual RL approach in contrast to a simple model-free framework. The last section introduces ideas for further research in RL for portfolio optimization and trading and draws some conclusions on the results obtained so far.

4.2 Multiperiod problem

This section explores the structure of the multiperiod portfolio problem and describes the essential elements of the MDP for linking it to RL. The framework is general enough to be adapted to various problem structures by properly choosing the risk and the cost functions.

Let a financial agent trades in a universe of N assets by allocating a fixed amount of initial capital c_0 that represents the cash asset. Trading occurs at several equally-spaced steps $t, t + 1, t + 2, \dots$ in an infinite time horizon, where the time step t denotes the entire time spanning between t and $t + 1$. The holding amount of each asset at time t is expressed as a vector of dollar values $\mathbf{h}_t = (h_t^1, h_t^2, \dots, h_t^{N+1}) \in \mathbb{R}^{N+1}$ whose first N terms are the dollars invested in the risky assets and the last term is the risk-free cash asset. Positive entries of the vector \mathbf{h}_t represent long positions, while negative entries represent short positions in the corresponding asset. The holding vector is evaluated for each asset i using the reference price p_t^i that usually is the average quote between the best bid and the best ask at the time t .

The net asset value (NAV) of the portfolio is $v_t = \mathbf{1}^T \mathbf{h}_t$ and the gross exposure of the portfolio is the sum of all position in absolute value $e_t = \sum_{i=1}^N |h_t^i|$, excluding the asset cash position. Therefore, the leverage of the portfolio is $l_t = \frac{e_t}{v_t}$. The portfolio returns is then the percentage variation of portfolio value through time $R_t^P = \frac{v_{t+1} - v_t}{v_t}$. One can also express the portfolio holding as a vector of weights $\mathbf{w}_t = \frac{\mathbf{h}_t}{v_t}$ with the usual constraint that they sum up to one, $\mathbf{1}^T \mathbf{w}_t = 1$.

The agent rebalances the portfolio at the beginning of each period t according to a vector of trades $\mathbf{a}_t \in \mathbb{R}^{N+1}$. They moves the portfolio holding

4. Reinforcement learning for trading in a market with frictions: a residual approach

according to

$$\mathbf{h}_{t+1} = (1 + \mathbf{r}_t) \circ \mathbf{h}_t + \mathbf{a}_t \quad (4.1)$$

where the vector \mathbf{r}_t represents the percentage returns $r_t^i = \frac{p_t - p_{t-1}}{p_{t-1}}$ for each security i with p_t as the current price and \circ indicates the element-wise product. The update rule in Eq. 4.1 assumes that the trade is carried out an instant before $t+1$, so the additional part of each holding a_t doesn't earn or lose value. The return for the cash asset is the risk free rate of return available on the market.

The total variation of the portfolio value on a subsequent time step is

$$\Delta v_t = \mathbf{1}^T \mathbf{a}_t + \mathbf{1}^T (\mathbf{r}_t \circ \mathbf{h}_t) \quad (4.2)$$

which composes of the traded value and the additional portfolio variation caused by the portfolio holding at the beginning of the period. While trading, the agent incurs in some costs that arises from two sources: market frictions and holding short positions. The first set of cost is parametrized by the following function

$$c_t^{\text{trade}}(\mathbf{a}_t) = \sum_{i=1}^N c_t^{\text{trade}}(a_t^i),$$

which is separable, since trading costs arise from each asset independently. One should note that the trading costs do not depend on the variation of the cash asset. The holding costs are instead parametrized by a function $c_t^{\text{hold}}(\mathbf{h}_{t+1})$ that depends on the portfolio holdings after the trades occurs. This function can take into account holding costs related to the short position contained in the portfolio or the borrowing cost for additional cash when leveraging the position is possible.

The portfolio is rebalanced according to the following self-financing condition,

$$\mathbf{1}^T \mathbf{a}_t + c_t^{\text{trade}}(\mathbf{a}_t) + c_t^{\text{hold}}(\mathbf{h}_{t+1}) = 0, \quad (4.3)$$

which means that the expenses for trading the risky assets are entirely covered by the variation of the cash asset value without any other external cash inflows.

We frame the portfolio optimization as an MPO problem, where the agent plans a set of decisions for the future given all the current information at disposal. An MPO problem requires forecasting relevant quantities to make a future trading decision, typically the asset returns and their covariance matrix. As outlined in Kolm et al., 2014, these estimates are, in practice, hard

4. Reinforcement learning for trading in a market with frictions: a residual approach

to obtain with a high degree of accuracy since a good model for the underlying dynamics of returns is required. Therefore, a model-free RL approach would be appealing because it does not require a considerable modeling effort for the problem. The quality of the forecasts highly influences the solution of an MPO and represents a challenging problem on its own.

We employ the MPO framework to better deal with transaction costs and assess their impact on the multiple rebalancing of the asset holdings. The SPO setting, as in the framework of Markowitz, 1952, consists in planning the trades just for the subsequent period without assessing any impact on the future state of the portfolio. Therefore, in the presence of transaction costs, SPO results are not optimal because the approach could cause large portfolio movements that create a massive amount of transaction costs if the trader needs to unwind the position in the following few rebalancing periods.

The MPO problem consists in planning the sequence of trades \mathbf{a}_t for each t over an infinite horizon. We denote the estimate of future quantities as $\hat{\mathbf{r}}_t$ conditioned on all the information available up to time t . The objective of the problem is the risk adjusted PnL of the portfolio

$$\text{PnL}_t = \hat{\mathbf{r}}_t^T \mathbf{h}_t - \kappa \psi_t(\mathbf{h}_t) - c_t^{\text{trade}}(\mathbf{a}_t) - c_t^{\text{hold}}(\mathbf{h}_t) \quad (4.4)$$

where ψ is a general function of the portfolio holding that accounts for the risk, κ is a risk aversion coefficient and c_t^{trade} and c_t^{hold} are the estimated cost functions.

The MPO problem is then,

$$\max_{\mathbf{h}_0, \mathbf{h}_1, \dots} \mathbb{E}_0 \left[\sum_{t=0}^{\infty} \gamma^t \text{PnL}_t \right] \quad (4.5)$$

where $\gamma \in [0, 1)$ is a discount factor for the infinite sum of values. In this objective only the initial \mathbf{h}_0 is known, while $\mathbf{h}_1, \mathbf{h}_2, \dots$ depend on the future trades $\mathbf{a}_1, \mathbf{a}_2, \dots$ and the returns estimate through the Eq. 4.1. The problem in Eq. 4.5 is convex, as long as the risk and cost functions are all convex. In that case, the problem is solved by convex optimization techniques as described in Boyd and Vandenberghe, 2009 and Boyd et al., 2017. The Eq. 4.5 represents an optimization problem where a trading plan $\mathbf{a}_t, \mathbf{a}_{t+1}, \dots$ is chosen so that it generates a series of portfolio holdings $\mathbf{h}_t, \mathbf{h}_{t+1}, \dots$. Only the first trade in the series is executed in practice because the trading plan is optimized again at the subsequent step. The purpose of optimizing the

whole sequence of trades is to evaluate the effect of a single trade in a longer time horizon. For instance, a trader could have difficulty dismissing a huge trade carried out at time t in a relatively short period. On the contrary, he can adequately balance the movement of its portfolio by obtaining the whole path of trades up to the end of the trading horizon.

The flexibility of the RL framework is beneficial for solving the problem in Equation 4.5, especially when the choices for the risk and the cost functions make the problem not convex. In such cases, optimal control algorithms cannot find any closed-form solutions due to the curse of dimensionality. We refer to the Chapter 2 and Chapter 3 for a detailed description of the MDP framework that mathematically formalizes the RL problem.

4.3 Financial model for the excess return

This section outlines the underlying assumptions of the financial model we use to simulate synthetic excess return with respect to the risk-free market rate and test the RL approaches. We refer to these returns as alpha returns or alphas, as it is common in the financial literature (Schneeweis and Spurgin, 1999). We start from the model of Kolm and Ritter, 2014, and we present a multiperiod problem for a single asset that is easily extendable to a multi-asset framework. Our proof of concept example employs a stylized dynamics consisting of an exponential decay with rate τ from a starting level r_0 , so that the excess return is expressed as $r_t = r_0 e^{-\frac{\log(2)}{\tau}t}$. Beyond this simple case, the alpha return can also be expressed as a sum of exponential decays with a different initial starting point and halflives of decay. In what follows, we assume these alpha returns to be the forecasted return provided by a portfolio manager's model. Therefore, the RL agent needs to determine the portfolio's size and trading rate by using previously generated forecasts. We believe that an exponential decay is a stylized representation of an alpha return predicted by a portfolio manager that inevitably approaches zero over time due to the market that arbitrage out the profitable opportunity.

The reward function for our RL problem is obtained from the risk adjusted PnL in Eq. 4.4

$$R_t = \hat{r}_t h_t - \frac{\kappa}{2} h_t^2 \sigma^2 - c_t(a_t) \quad (4.6)$$

where the risk function is quadratic, the risk aversion parameter is chosen $\kappa = 1e - 5$ as in Kolm and Ritter, 2014 and the cost function follows the

4. Reinforcement learning for trading in a market with frictions: a residual approach

model of Kyle and Obizhaeva, 2016

$$c_t(a_t) = \kappa_1 |a_t| + \kappa_2 a_t^2 / (0.01pV) \quad (4.7)$$

being p the average daily price of the synthetic asset and V its average daily traded volume. For the purpose of transaction costs, the price of the asset is fixed as the daily trading volume, which is also rescaled of by the number of simulated periods in a day.

As a comparison for the performance of our RL approach, we use the solution of Gârleanu and Pedersen, 2013, as in Brini and Tantari, 2021, since it is based on a renowned model in the financial literature and it is also optimal when the excess returns are modeled as a sum of mean reverting dynamics. To be consistent in computing the solution under our setting, we need to choose the parameter λ (cost multiplier) for Eq. 3.3 according to our choice for the cost function. Indeed, we assume that $\Lambda = \frac{\kappa_2}{0.01pV}$. From Gârleanu and Pedersen, 2013, we know that $\Lambda = \lambda\sigma^2$, so we set $\lambda = \frac{\kappa_2}{0.01pV\sigma^2}$. In this way, we can compute the trading rate a/λ according to Eq. 3.6.

Calibrating the quadratic cost function of our problem to the total cost function in Kolm and Ritter, 2014 is necessary to let the financial agent bear the same amount of transaction costs, which we believe to be realistic enough. At the same time, it allows the learned trading policy to be compared with the solution of Gârleanu and Pedersen, 2013. For these reasons, we choose the parameter κ_2 so that the amount of quadratic trading costs equals the sum of the linear and the quadratic part in that work according to Eq. 4.7. We simply neglect the absolute part of the cost function in Eq. 4.7 by posing $\kappa_1 = 0$. Then we create an equally spaced interval, $[1e + 2, 1e + 5]$, of $N = 10000$ possible trades a_t . Considering $\kappa_1 = 2.89E-4$ and, $\kappa_2 = 7.91E-4$ as in Kolm and Ritter, 2014, we compute the corresponding costs $c_t(a_t)$ through Eq. 4.7. To approximate the same amount of transaction cost with the quadratic part of the function only, we use OLS to estimate the cost parameter so that $c_t(a_t) = \kappa_2 a_t^2 / (0.01pV)$. The resulting parameter is $\kappa_2 = 0.0009 = 9$ bps. Figure 4.1 presents the different cost functions according to the model of (Kyle and Obizhaeva, 2016) and our fitted cost function.

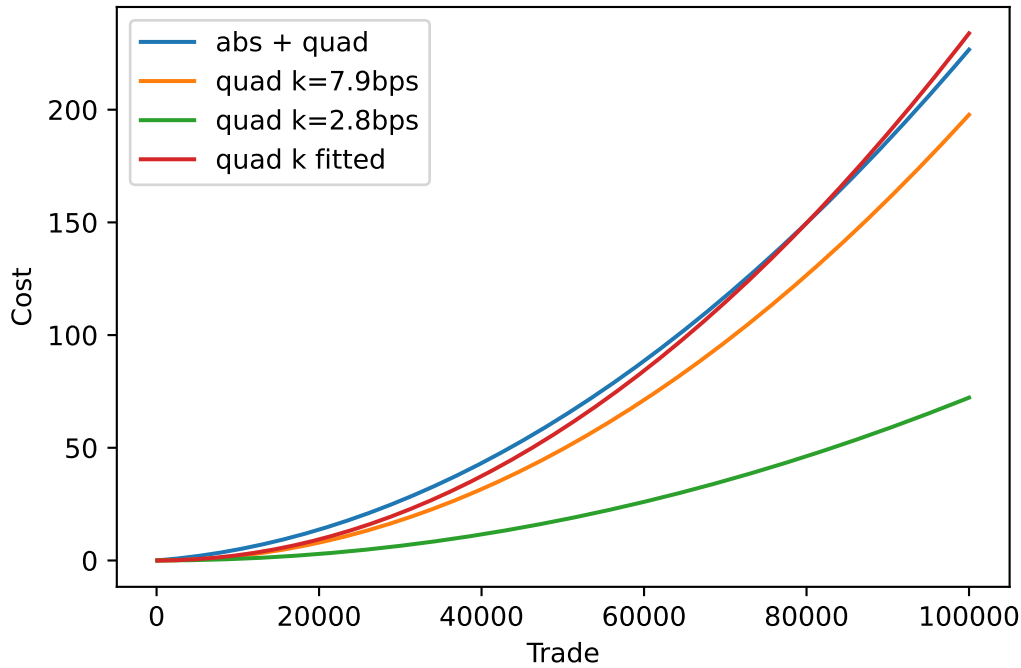


Figure 4.1: The blue represents the transaction cost function used in Kolm and Ritter, 2014. The green and the orange accounts only for the quadratic part with different parameter κ_2 . The red line is the quadratic cost function with the fitted κ_2 , which best approximates the full cost function.

4.4 The residual RL approach: empirical results

This section introduces the residual RL approach and presents its performance compared with the model-free approach. As stated in Section 4.3, we use the solution provided in Gârleanu and Pedersen, 2013 as a benchmark to reach, and we motivate the need of the RL agent to point toward such a solution in solving the outlined problem under the MPO setting.

4.4.1 Residual RL approach

We know from the empirical experiments presented in Chapter 3 that the portfolio problem described in Eq. 4.5 can be solved through the solution

4. Reinforcement learning for trading in a market with frictions: a residual approach

of Gârleanu and Pedersen, 2013. The application of dynamic programming for MPO problems is crucial to account for transaction costs that may burn the trader’s profit due to large portfolio rebalancing. On the other hand, the Markowitz portfolio theory (Markowitz, 1952) solves the same problem in an SPO setting as

$$h_t^M = (\gamma \Sigma)^{-1} r_t, \quad (4.8)$$

where Σ is the covariance matrix of the assets. One can see the benchmark solution of Gârleanu and Pedersen, 2013 as an extension of Markowitz’s solution in an MPO, where it becomes necessary to consider the opportunity cost of entering a trade and then modify the position according to the decay of the alpha returns over time. If transaction costs are neglected, as in an SPO setting, a strong signal for specific assets may result in a considerable exposure for the trader. However, suppose the term structure of the alphas quickly decays to zero. In that case, the trader will find trouble in the subsequent time steps trying to liquidate its position that became unprofitable and incurring consistent transaction costs. Considering the decay of the signal help improve the management of position sizes and avoid the erosion of profitable trades.

The intuition behind the residual approach is to leverage prior knowledge of the problem by exactly exploiting the optimality of the Markowitz solution in Eq. 4.8 in the absence of transaction costs. Since the trader acts in a dynamic setting where transaction costs play a relevant role, we reduce the problem by learning the nonlinear impact these costs have on the profitability of the Markowitz trades. Knowing that Eq. 4.8 optimizes the mean-variance part of the reward function in Eq. 4.6, RL corrects the solution to reflect the impact of the market frictions. Therefore, we express the quantity to trade at each time step as a residual part

$$h_t = h_{t-1}^M (1 - a_t) \quad (4.9)$$

where a_t is a real-numbered action between 0 and 1 that helps express the percentage of the Markowitz portfolio, which is better not to buy for the current time t to contain the cost and follow the optimal path. The aim of the RL agent is then to learn how to trade a residual version of the Markowitz portfolio in Eq. 4.8 to match the solution in Eq. 3.5. Indeed, we know from Gârleanu and Pedersen, 2013 that the trader needs always to slow down the trade with respect to what is suggested by the Markowitz solution. It is equivalent to saying that the agent learns to move towards the Markowitz

portfolio at a rate defined by the decay of the signal. If the decay is too fast, the agent will slow down the trade because the profitable opportunity is not supposed to be durable. On the contrary, if the decay is slow, the agent trades more closely to Markowitz since the profitable opportunity is long-lasting. The impact of this whole mechanism is lowering the price paid to rebalance the portfolio at every period with respect to an approach that does not account for the effect of the frictions, which we know from (Guéant, 2013; Kyle and Obizhaeva, 2016; Patzelt and Bouchaud, 2018) to be highly nonlinear. Hence, the important role of the transaction costs in discovering a profitable opportunity and their functional form makes RL with nonlinear approximators an appealing framework for solving these classes of problems.

The following sections provide some results of the synthetic experiment performed to show the approach’s effectiveness. We will present the simplest case of the single asset level.

4.4.2 Simulation results

We show the first proof of concept of the residual RL approach by assuming to be a fund manager that already developed a robust model to forecast excess returns. The forecasting step represents a complex problem on its own, so we focus on using the generated forecast to trade. Simulating an alpha term structure with an exponential decay over the horizon of the simulation resembles the capacity of the forecasting model of the manager to produce valuable alpha at a distant point in time. As already anticipated, we consider the simulated exponential decay is a suitable structure to represent the shrinking of the excess returns through time.

The generation of synthetic data allows focusing on testing the RL approach in a controlled environment with a benchmark solution. It will enable understanding when the RL agent retrieves the benchmark solution and assesses the possible causes when it is not the case. The generated alpha term structure is then used as the input signal of the RL agent that needs to determine how to trade according to the given information.

For our synthetic experiments, we use the proximal policy optimization (PPO) algorithm (Schulman et al., 2017), which is a policy-based deep RL algorithm renowned for its performance in a variety of control problems. It is specifically suitable for continuous control so that it is easily comparable to the continuous-time solution provided by the benchmark. For more details about the PPO algorithm, we refer to Section 3.7.

4. Reinforcement learning for trading in a market with frictions: a residual approach

We train the PPO algorithm episodically, and we simulate alpha term structure with starting level r_0 , a half-life of decay τ , and a fixed-length T . For the experiment results proposed in the section these parameters are respectively $r_0 = 30$ bps, $\tau = 240$ time steps and $T = 1200$ time steps. Considering the simulation period as a month, we have approximately $1200/21 = 57.14 \sim 57$ time steps per trading day. One can divide each day of such a month into $D = 57$ time steps to better interpret the time horizon. We express the alpha decay in these terms to help choose the daily volatility to account for in our financial model, which can be scaled linearly with the time steps.

We run the algorithm for 3000 episodes by adopting two approaches: the residual RL approach described in the previous subsection and the standard model-free approach. In the former, the agent chooses the action according to Eq. 4.9, while in the latter the action is chosen as a real-value in a given interval. For each in-sample episode, we keep track of the cumulative episodic reward, which is the sum of Eq. 4.6 for all the time steps T in an episode, and we present the training progresses in Figure 4.2. The policy network takes as input, i.e., the agent's state, the excess return r_t , the current portfolio holding h_t , and the time remaining to the end of the forecasting horizon $T - t$. We appreciate that the residual RL approach converges faster to a stable solution equal to the benchmark one. On the other hand, the model-free approach struggles to find a solution close to it and requires more training episodes to reach in-sample performance stability.

4. Reinforcement learning for trading in a market with frictions: a residual approach

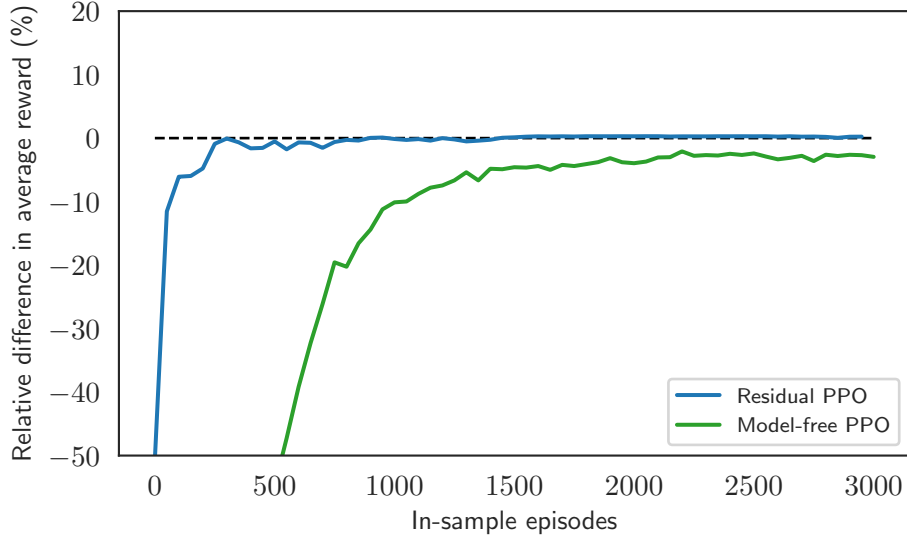


Figure 4.2: convergence to the benchmark solution for the residual RL and the model-free approach trained using PPO with alpha term structure samples. For each approach, we train several agents and calculate their average cumulative reward in-sample to choose the best performing one. Then a rolling average of the relative difference in reward with respect to the benchmark is computed to compare them.

Recalling that r_0 , τ are respectively the starting point and the halflife of decay of the alpha term structure, the relative difference in average rewards presented in Figure 4.2 is obtained as

$$\begin{aligned}
 \frac{E_{\text{PPO}}[G^{r_0, \kappa}] - E_{\text{bnch}}[G^{r_0, \kappa}]}{E_{\text{bnch}}[G^{r_0, \kappa}]} &\approx \frac{\frac{1}{N} \sum_i G_{\text{PPO}}^{r_0, \kappa}(i) - \frac{1}{N} \sum_i G_{\text{bnch}}^{r_0, \kappa}(i)}{\frac{1}{N} \sum_i G_{\text{bnch}}^{r_0, \kappa}(i)} \quad (4.10) \\
 &= \frac{\sum_i (G_{\text{PPO}}^{r_0, \kappa}(i) - G_{\text{bnch}}^{r_0, \kappa}(i))}{\sum_i G_{\text{bnch}}^{r_0, \kappa}(i)}
 \end{aligned}$$

where i denotes a particular path, $E_{\text{PPO}}[G^{r_0, \kappa}]$ is the average cumulative reward G , i.e., the sum of Eq. 4.6 for all the time steps, obtained by the PPO strategy over an alpha term structure characterized by a specific r_0 and τ . The figure highlights the effectiveness of selecting a reasonable starting point and learning the remaining part of the trading solution, i.e., the residual.

4. Reinforcement learning for trading in a market with frictions: a residual approach

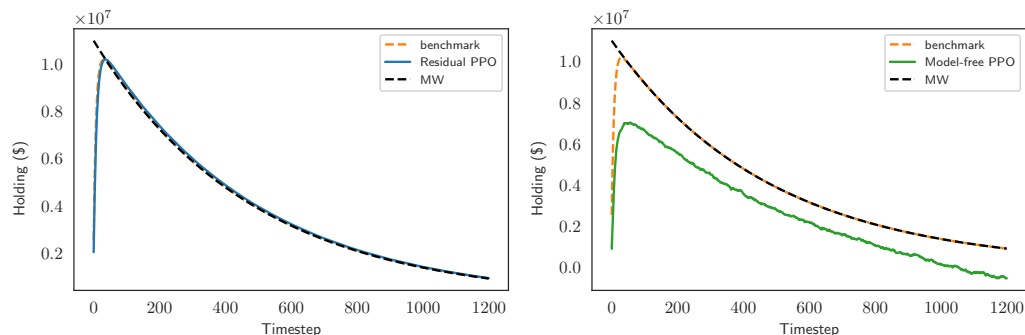


Figure 4.3: The left panel shows holding path obtained by the residual RL algorithm at the end of the training convergence of Figure 4.2, whereas the right panel shows the same for the model-free approach.

We remark then that as a measure of in-sample performance, we use the difference in average reward relative to the reward benchmark instead of using the average relative difference in reward, as in the last part of Eq. 4.10. The former is more efficient for our purpose and is not distorted by possible negative rewards obtained by the PPO algorithm over some alpha term structure paths.

Once we acknowledge that the residual RL approach converges faster in our controlled environment, we pick the best performing trained agents that we use to produce Figure 4.2, and we test them out-of-sample on a different alpha term structure. Figure 4.3 shows the trading paths realized by the residual RL and the model-free together with the benchmark and the Markowitz paths, where the last one represents the starting point solution for the residual RL. The converged and stable solution learned by the residual PPO can closely track the benchmark holding path, trading towards a peak and then slowly moving out of that position until the end of the generated trajectory. On the other hand, the struggle of the model-free approach in converging to a solution close to the benchmark is evident when looking at its generated trading path. The model-free algorithm trades less towards the peak and move out from the position too quickly, resulting in suboptimal performance.

One can look at the same result from the perspective of the learned policies in Figure 4.4. The figures present the traded amount by the learned strategy in correspondence to different levels of the alpha term structure when the other two parameters of the model, the current holding and the

4. Reinforcement learning for trading in a market with frictions: a residual approach

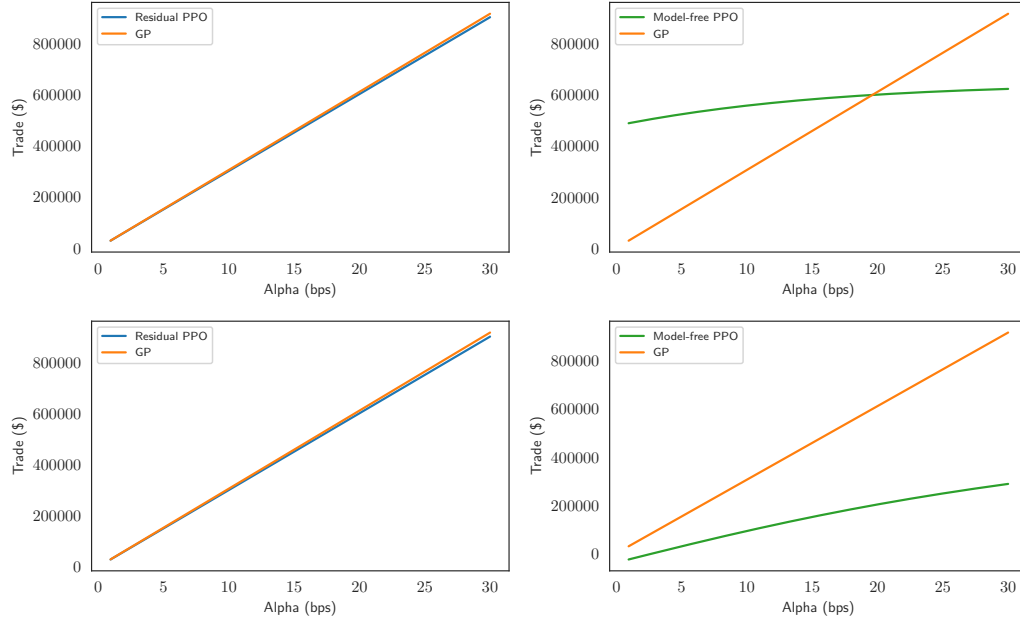


Figure 4.4: Trade size (in dollars) conditional on a current holding of zero in the risk security and remaining time to end equal to 1100 time steps (top panel) and 100 time steps (bottom panel) for the residual PPO approach and the benchmark solution as a function of the predicted alpha term structure.

time by the end of the simulated alpha term structure, are fixed. While residual RL closely tracks the monotonic benchmark policy for the simulated problem, the model-free cannot learn its slope, trading suboptimally. The result holds both when the traded are carried out with many time steps $T - t$ remaining or at the end of the trading period when $T - t$ is small.

However, only looking at Figure 4.4 does not provide any information about the performance obtained by the learned policies themselves since the evaluation needs to be tested for robustness and carried out over many different generated alpha returns paths. To this end, Figure 4.5 presents the cumulative distribution functions (CDFs) of the cumulative reward on a set of simulated alphas with different starting levels and halflives for the two PPO training methods. We still observe that the CDF of the residual PPO closely tracks the distribution of the benchmark solution. At the same time, the CDF of the model-free PPO presents noticeable differences and performs worse than Markowitz for a considerable part of the distribution. To statistically analyze these distributions, we use the Kolmogorov-Smirnoff

4. Reinforcement learning for trading in a market with frictions: a residual approach

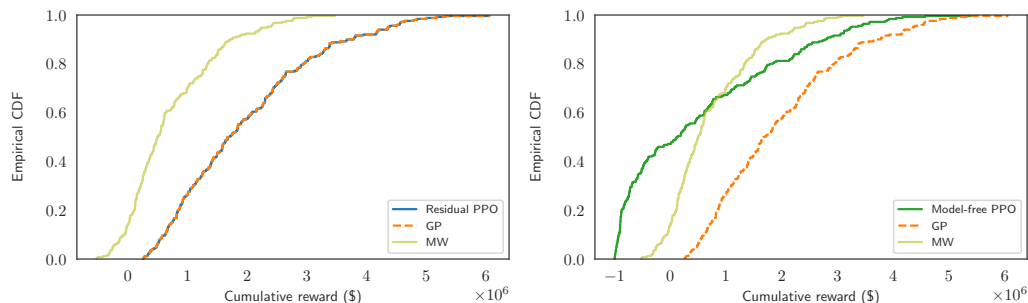


Figure 4.5: The left panel shows the empirical cumulative distribution function for the cumulative reward distribution for the residual RL approach over 250 alpha-term structure simulations of length 1200 time steps each. The same is shown on the right for the model-free approach. The benchmark solution and the Markowitz solution for the same problem are represented in each figure.

(KS) test for comparing the shape of the distribution and the Student’s T-test for comparing their average values with respect to the distribution of the benchmark solution. In the residual PPO case, the values for both tests are close to the unit, and the null hypothesis is not rejected. There is no statistical evidence that the residual RL and the benchmark distributions of the cumulative rewards differ in shape and mean. The opposite holds for the model-free case where the test confirms that the distribution differs, as the reader can appreciate visually.

The results highlighted in this subsection clarify the importance of simplifying the process of sampling actions when algorithms like PPO are applied to financial trading problems, and the continuous space of actions is enormous. In our first proof-of-concept, the residual PPO allows us to leverage prior information about the financial environment, i.e., knowing that Markowitz is optimal without transaction costs, and significantly speed up the learning process. When model-free cannot even find the solution, residual PPO is robust enough to guarantee an approximate solution close to the GP one.

4.5 Further directions

The synthetic experiments presented in the previous section are the first proof of concept of the adaptation of the model-free reinforcement learning approach to the financial trading problem. Since these models are known

to require a long training runtime and to suffer from instability during the training process, the approach we have presented is beneficial for both issues. However, residual RL is just one of the nuances of model-free RL that still guarantees the flexibility of being model-free while being driven by prior knowledge of the financial environment. The following subsections will explore two directions of research that can improve the model-free approach toward a more efficient and effective training process while maintaining its original flexibility.

4.5.1 Residual approaches

One of the first proposals for modifying the model-free approach is to learn just a residual part of the total solution. This method, which we preliminary explored in our synthetic experiments, changes the RL algorithm at the action level. The agent’s action is no longer expressed as a real-valued amount reflecting the trade to perform, but it is expressed as a difference from a known solution. As already mentioned in the introduction of this section, this approach is helpful in domains such as robotics, where it is difficult to describe an exact model of the environment. However, some of its properties are known and can be used to outperform a fully model-free approach.

The residual approach still guarantees a certain degree of flexibility in choosing the prior knowledge to leverage. We based the agent’s residual action on the Markowitz portfolio in our experiments, but different choices are available in the financial landscape. The residual approach is also suitable for learning a variation of risk-based allocation methods such as the global minimum variance portfolio (Kempf and Memmel, 2006) or the risk parity portfolio (Asness et al., 2012). The adoption of the Markowitz solution as a baseline is motivated by the comparison with the benchmark solution of (Gârleanu and Pedersen, 2013) which modifies the Markowitz solution by accounting for the transaction costs effect.

4.5.2 Regularization approaches

A different approach to improve the model-free RL performance is to use some regularization method. Regularization in RL has not been considered yet at the same level as the deep learning field, although some relevant empirical testing in that direction appears (Liu et al., 2019b). Other works are specific

4. Reinforcement learning for trading in a market with frictions: a residual approach

to either value-based algorithms, such as Farebrother et al., 2018 for Deep-Q-network or Liu et al., 2019a for policy-based methods.

Differently from the residual approach, regularization methods act at the policy level by imposing constraints to the objective function to penalize certain behaviors or by imposing specific training conditions on the policy network, such as the Batch Normalization (Ioffe and Szegedy, 2015), Dropout (Srivastava et al., 2014) or weight decay (Zhang et al., 2018).

Among all of those regularization methods, the possibility of constraining, hence directing the policy towards having specific characteristics, is the most appealing for the financial problem. Constraining the learned policy with a soft penalty is interesting for enforcing specific properties that we know the optimal policy has, such as the monotonicity with respect to the alpha signal as in Figure 4.4.

4.6 Conclusions

In this work, we presented an initial proof-of-concept of the residual RL approach, which aims to leverage prior knowledge of the domain and improve the model-free RL efficiency in solving financial problems. The approach has proven beneficial for learning a known solution in a synthetic financial market.

The work will extend in two directions by proving the effectiveness of the residual method over a large cross-section of assets, either synthetic or real, and by adopting different regularization methods that constrain the learned policy and have already been employed in the classical RL literature.

The former will be challenging because of the runtime required when the cross-section of assets increases. However, the residual approach learns only the nonlinear part of the solution relative to transaction costs, which is usually considered separable and asset-dependent. Therefore, from our next series of experiment we expect the approach to be scalable over a large cross-section of assets since the neural network policy would be able to generalize and learns how to trade assets with different costs level. In other words, there is no need to train a large neural network that accounts for all the different assets, but it will be enough to train a simple network that generalizes over the cross-section.

The latter method will be challenging for devising a proper functional form of the constraint that enforces known characteristics of the optimal

4. Reinforcement learning for trading in a market with frictions: a residual approach

policy without conflicting with the existent training algorithm.

Chapter 5

Reinforcement Learning Policy Recommendation for Interbank Network Stability

5.1 Introduction

The increasingly recurrent and impactful crises affecting the socio-economic system have called for a deep rethinking of the economic theory. Firstly, the literature has made an effort to understand and include in the economic models the sources of contagion. Regardless of the modeling approach used, which ranges from New Keynesian models solved globally or using reduced functional form (see, for instance, Boissay et al., 2016, Gertler et al., 2020, Svensson, 2017) to agent-based models and the most recent network-oriented approaches (see Battiston et al., 2012a; Battiston et al., 2012b, Georg, 2013, Haldane and May, 2011, Upper, 2011, Capponi et al., 2020, Calice et al., 2020), there is a general agreement that identifies interaction and heterogeneity as the drivers of endogenous crises. Moreover, the post-Lehman studies have placed particular emphasis on the propagation of contagion, determining the direction of the attack from financial to real markets and its fuse in the portfolio structure of financial institutions (see Brunnermeier et al., 2012). Many interesting studies, for example, have identified the source of contagion in the asset or liability side of banks' balance sheets. Among them, the effect of the fire-sale price and the (re)payment system between creditors and debtors have proven to be particularly important in generating financial

5. Reinforcement Learning Policy Recommendation for Interbank Network Stability

instability (see Acharya and Yorulmazer, 2008a, Angelini et al., 1996, Dasgupta, 2004 Rochet and Tirole, 1996). In this vein, maturity transformation, sharing risk, herding behavior, and interbank linkages are just some various components able to trigger instability or collapse in financial markets (see Acharya and Yorulmazer, 2008b, Allen and Gale, 2000 and Tedeschi et al., 2021, among the many).

Once the origin of the disease and the channels through which it spreads have been identified, the literature has turned to treatment, that is the identification of the best tools to mitigate financial contagion. Probably spurred on by the words of the former governor of the European Central Bank, Trichet, who said that he felt abandoned by conventional tools during the sovereign debt crisis and advocated the development of alternative instruments borrowed from different disciplines (see Trichet, 2010), the scientific community has heavily focused on the development of new tools aimed at overcoming systemic instability. In this regard, several conventional and non-conventional monetary policies, as well as other alternative tools, have been proposed, but their effects on financial stability are controversial and depend on the overall economic condition (see Goldberg et al., 2020 and Altavilla et al., 2021). A strand of literature, for example, has emphasized the importance of a strict, rule-based and predictable monetary policy to tame systemic risk (see Jiménez et al., 2014 and Taylor, 2011). On another side, instead, different studies have bet on alternative rules, compatible with the underlying economic conditions (see Boissay et al., 2021, De Grauwe, 2011 and Galí, 2015). Unfortunately, the weak empirical evidence, due to the fairly recent development of these alternative techniques, which also include the so-called macro-prudential policies, makes it difficult to prove the supremacy of one approach over the other. While the empirical facts are still uncertain, recent theoretical models have attempted to resolve this “certamen”. An interesting contribution in this direction is the model of Boissay et al., 2021. The authors use a globally solved New Keynesian model with heterogeneous agents to generate endogenous crises. The paper compares two monetary policy instruments, one that follows a strict inflation targeting rule and the other that allows the central bank to curb financial booms and busts. The authors show how the policies that mitigate output fluctuations, by acting on agents’ expectations, help in preventing financial crises. In support of cyclical policies determined by the economic background, there are also many agent-based models (see, Cincotti et al., 2012, Giri et al., 2019 and Riccetti et al., 2018, among the many). The approach of generating complex dynamics in evol-

5. Reinforcement Learning Policy Recommendation for Interbank Network Stability

ing systems is an ideal environment for testing the effect of (un)conventional policies/measures on the financial stability.

Following this last line of research, in this paper we are explicitly interested in understanding the effect that an unconventional and environmentally dependent policy recommendation has on the stability of the interbank system. From the point of view of the functioning of the interbank market, our work follows Berardi and Tedeschi, 2017, where financial institutions establish preferential lending arrangements to insure themselves against the unexpected withdrawal of deposits. Financial connections might change over time via a preferential attachment evolving procedure (see Barabási and Albert, 1999) such that each agent can enter into a lending relationship with others with a probability proportional to a fitness measure. Specifically, the attractiveness of agents is based either on their high supply of liquidity or on their low interest rate. The authors show how the implementation of one or the other strategy generates different architectures of the credit network, which impacts on the spread of systemic risk differently.

The originality of this work with respect to the above-mentioned one concerns the mechanism that drives banks to choose between the two strategies. Where in Berardi and Tedeschi, 2017 the choice is exogenous and fixed, here we introduce a time dependent, reinforcement learning based policy recommendation that directs banks to optimize the entire banking system long term fitness. Specifically, the regulator directs the interbank system towards an optimal strategy that chooses between favoring a high liquidity supply with respect to a low interest rate, by collecting information from the environment. Once the policy recommendation is made public, each bank signals itself in the interbank market according to its optimal level of liquidity supply or interest rate, which are used to establish credit agreements via the above-mentioned preferential attachment mechanism.

Compared to Berardi and Tedeschi, 2017, therefore, the reinforcement learning mechanism allows us both to endogenize and identify the optimal strategy and to model a policy recommendation useful to tame systemic risk. Although this tool is very useful for modeling reward-seeking behavior of agents in complex systems¹ (see(Osoba et al., 2020)), to the best of our knowledge it is barely employed in the agent-based framework. Interesting

¹We refer the reader to Charpentier et al., 2021 and Mosavi et al., 2020 for comprehensive reviews of different use cases of reinforcement learning in financial and economic context.

5. Reinforcement Learning Policy Recommendation for Interbank Network Stability

exceptions are Liu et al., 2018 and Lozano et al., 2007, which uses reinforcement learning to model the credit allocation strategy of financial institutions in the interbank market. Apart from the modeling differences - omitted here - that distinguish us from those works, it is instead important to point out the methodological distinction. Where these works use a tabular reinforcement learning algorithm, as proposed by (Watkins and Dayan, 1992), we use a state-of-the-art deep reinforcement learning algorithm with neural network approximators (Schulman et al., 2017), which describe the complex reward-seeking behavior. While the advantages and disadvantages of these algorithms are well documented and concern issues such as the computational efficiency, the curse of dimensionality and the convergence (Bellman, 1956), what is certain is the better performance of the algorithm with neural network in solving complex problems where the underlying environment not only changes rapidly, but it is also defined by the different forces that relate and compete with each other. These capabilities have already demonstrated effective in solving complex financial and economics problems (see (Du et al., 2020; Jiang et al., 2017; Lin and Beling, 2020; Zhang et al., 2020)).

Without delving into technical details, some clarifications on how the proposed algorithm works should be made. The selected reinforcement learning algorithm optimizes an objective function, which in our context corresponds to the aggregate fitness of the interbank system. The optimization is carried out by training a neural network model. The neural network receives as input variables concerning the economic conditions of the interbank system and return as output the strategy, i.e. the policy recommendation directing the system towards competing on liquidity supply rather than on the interest rate.

One of the criticisms to this family of algorithms concerns the interpretability and measurability of the impact that inputs have on outputs. The output, in fact, often appears as a black box whose determinants remain hidden to the user. To avoid this problem, we act in the following way. Firstly, we limit the choice of inputs to variables easily available to the regulator. To this end we use aggregate systemic variables such as minimum, maximum and average interest rate and liquidity supply of the interbank system. The choice of a limited set of input variables allows us not only to understand their effects in determining the output, but also to model a system with incomplete and asymmetric information (see Bernanke et al., 1999). Secondly, we directly study the impact that each input has on the determination of the output through the SHapley Additive exPlanation (SHAP) framework

5. Reinforcement Learning Policy Recommendation for Interbank Network Stability

(Lundberg and Lee, 2017).

The introduction of the previous reinforcement learning framework into the interbank market model proposed by Berardi and Tedeschi, 2017 allows us to draw some important conclusions about the systemic stability of the system and to determine some policy interventions capable of curbing contagion. Firstly, the proposed algorithm fully endogenizes the evolution of the interbank network, whose architecture therefore changes over time. In this way, we are able to identify that the topology that emerges when the policy recommendation suggests a high supply of liquidity is more resilient in the face of exogenous shocks. Also, at individual level, this policy produces better microeconomic performance. In this circumstance, banks are less heterogeneous, which generates a uniform risk exposure among counter-parties able to favor the resiliency of the system. The negative impact of heterogeneity on systemic stability is in line with various theoretical and empirical studies (see Caccioli et al., 2012, Iori et al., 2006 and Tedeschi et al., 2012). On the other hand, the worse performance of a system dominated by low interest rates reflects the empirical evidence. Indeed, it is well documented that a credit market dominated by "low-for-long" interest rates adversely affects both the banks and the stability of the economy. For financial institutions, low rates might reduce resilience by lowering profitability, and thus their ability to replenish capital after a negative shock. For the system, this strategy would encourage risk-taking which undermines systemic stability (see Bindseil, 2018, for a general overview on the topic). Finally, our results suggest that the policy recommendation implemented via the reinforcement learning is more able to mitigate systemic risk in comparison with alternative tools.

The rest of the work is organized as follows. In Section Section 5.2 we present the functioning of the interbank market, placing particular emphasis on the evolution of the credit network and the implementation of the reinforcement learning algorithm. In Section Section 5.3 we show the results. Specifically, we follow three steps: firstly, we verify performances and robustness of the reinforcement learning algorithm; secondly, we investigate its implication on the interbank network morphology and on the performances of the financial institutions; thirdly, we present the effect on the interbank systemic stability of the policy recommendation. Finally, Section Section 5.4 concludes with some remarks on the achieved results and the provided contribution.

5.2 Model

In this section, we describe the formation and evolution of credit relationships between financial institutions. Due to unexpected deposit movements, banks face liquidity surpluses or shortages, which induce them to enter the interbank market as lenders or borrowers. Each bank own its preferential credit channels and the set of all relationships among institutions defines the interbank network where banks try to fulfill their liquidity needs, establishing bilateral lending agreements. Within the market, banks signal their credit conditions through an attractiveness measure. We model the bank fitness as a combination of a policy recommendation and private information. The first ingredient is a signal obtained via a reinforcement learning mechanism, through which the regulator directs banks to choose the best strategy given the underlying environmental conditions. In particular, the regulator recommends the weight to assign to high liquidity supply rather than to low interest rates, thus directing the competition. The second ingredient is a private signal, based on the bank’s capital structure, consisting of the actual interest rate and credit provision offered. Credit relationships, then, might change over time via a preferential attachment evolving procedure that depends on bank fitnesses.

5.2.1 The interbank market microstructure

We consider a sequential economy operating in discrete time, which is denoted by $t = \{0, 1, 2, \dots, T\}$. At any time t , the system is populated by a large number N of active banks $i, j, k \in \Omega = \{1, \dots, N\}$. Financial institutions interact to each other through credit relationships represented by the set V_t , whose elements are ordered pairs of different banks. Banks (nodes or vertices) and their connections (edges or links) form the interbank network $G_t = (\Omega, V_t)$. The daily balance sheet structure of each bank is defined as

$$L_t^i + C_t^i + R_t^i = D_t^i + E_t^i, \quad (5.1)$$

where assets are on the left-hand side and liabilities on the right-hand one. In particular, L, C and R represent long term assets, liquidity and reserves, while D and E deposits and equity of bank i at time t . Reserves are a portion of deposits, $R_t^i = \hat{r}D_t^i$, where the required reserve rate, \hat{r} , meets the legal requirement of 2%.

5. Reinforcement Learning Policy Recommendation for Interbank Network Stability

At every time period t , deposits are exogenously shocked and the balance sheet in Eq. 5.1 modifies accordingly. Specifically, deposits evolve as

$$D_t^i = D_{t-1}^i(\mu + \omega U(0, 1)), \quad (5.2)$$

with $U(0, 1)$ a uniformly distributed noise between 0 and 1 and μ and ω modelling the expected number of negative shock and thus different market conditions. On the one hand, financial institutions with a negative change in deposits and subject to a complete erosion of their liquidity, $\Delta D_t^i + C_t^i \leq 0$, become potential debtors in the interbank market. On the other hand, banks which suffer a small negative shock or an increase in deposits that raises their liquidity, $\Delta D_t^i + C_t^i > 0$, become potential creditors in the system. Consequently, the respective demand d_t^i and supply s_t^i of liquidity of potential borrowers and lenders are given by

$$\begin{cases} d_t^i = |\Delta D_t^i + C_t^i| \\ s_t^i = \Delta D_t^i + C_t^i. \end{cases}$$

Since we do not assume a walrasian tâtonnement mechanism, the system may endogenously generate a mismatch between credit supply and demand. Moreover, since the interbank network is not fully connected, even at a micro level the demand for liquidity of a borrower bank might not match the credit supply offered by the lender banks connected to it. Specifically, we define the granted loan from a generic lender i to a generic borrower j as $l_t^{i,j} = \min(s_t^i, d_t^j)$. Borrowing banks that are rationed in the interbank market can sell their long-term assets at a fire-sale price as a method of last resort. The amount of loan the borrower has to sell for covering its residual liquidity need is equal to $\Delta L_t^j = \frac{d_t^j - s_t^i}{\rho}$, where ρ is the 'fire-sale' price

At the beginning of the next day, the repayment round takes place. Financial institutions encounter a new deposit movement that increases or decreases their liquidity. On the one hand, lending banks facing a positive (negative) change in deposits remain potential creditors (became potential debtors). On the other hand, borrowing banks face different scenarios depending on whether the deposit shock is positive or negative. Specifically, in the case of a positive shock, it can happen that: i) the change in deposits is sufficient to repay the principal and the interest, or ii) the deposit variation is not sufficient to cope with the loan. In the first case, the debtor can easily meet her obligations, but in the second case, she must sell, at a 'fire-sale'

price, an amount of long term assets sufficient to fully repay the creditor. On the other hand, in the case of a negative shock, banks must sell their long-term assets to pay for previous interbank borrowings and meet the new liquidity needs. All institutions that do not raise enough liquidity to meet their obligations via the fire sale fail, thus creating a bad debt on the lender. The creditor's loss, $B_t^{i,j}$, is simply equal to the granted loan after the liquidation of the debtor assets. Hence, the equity of the bank i obeys the following law of motion:

$$E_t^i = E_{t-1}^i + \sum_j l_{t-1}^{i,j} r_{t-1}^{i,j} - \sum_{j \subseteq \theta_t^i} B_t^{i,j} - (1 - \rho) \hat{L}_t^j, \quad (5.3)$$

where the second term on the right-hand side is the repayment, at the agent-specific interest rate $r^{i,j}$, of the granted loan $l^{i,j}$, the third term is the bad debt of the subset of the bank i clients, θ_t^i , unable to repay their debts back because they go bankrupt and the last term represents fire sales. If the bank has not fulfilled the loan requirements (i.e., if she is unable to repay the principal and interest in full), the lender does not provide credit any longer, forcing her to exit the market. Thus, the borrower exits the market when assets fall short of liabilities, that is $E_t^i < 0$. The failed banks leave the market. The banks exiting in t are replaced in $t + 1$ by new entrants, which are on average smaller than incumbents. So, entrants' size is drawn from a uniform distribution centered around the mode of the size distribution of incumbent banks (see Bartelsman et al., 2005).

5.2.2 Banks microfoundations: the dynamics of lending agreements and trading strategies

In order to meet their liquidity needs, at the beginning of each day, agents meet in the interbank market and sign bilateral potential lending agreements representing the directed links $(i, j) \in V_t$. These agreements can be interpreted as credit lines, which are valid during t , and can be used at the request of the borrower j in case of the lender i available liquidity. The set of all potential lending agreements reproduces the interbank network topology. Let us now explain in detail the mechanism that governs the formation/evolution of credit relationships between financial institutions. We assume that banks are risk neutral agents operating in a perfect competition environment with the purpose of optimizing their expected profit. The bank i expected profit

5. Reinforcement Learning Policy Recommendation for Interbank Network Stability

for a loan provided to j is given by

$$\mathbb{E}[\Pi_t^{i,j}] = p_t^j(r_t^{i,j}c_t^{i,j}) + (1 - p_t^j)(\xi A_t^j - c_t^{i,j}) + \phi A_t^j - \chi A_t^i, \quad (5.4)$$

where p_t^j is the probability that the borrower does not fail, $r_t^{i,j}$ the interest rate asked by the lender i to the borrower j , $c_t^{i,j}$ the maximum amount i is willing to lend to j . Moreover, ξ is the liquidation cost of assets, A_t^j , pledged as collateral, and ϕ and χ the screening costs of creating a credit link that decrease with the debtor dimension and increase with the creditor size (see Dell’Ariccia and Marquez, 2004 and Maudos and De Guevara, 2004, for empirical evidence). Specifically, Eq. 5.4 captures the lender’s expected revenue in the event that the borrower does or does not meet her obligations (first and second term on the right side, respectively), and the opportunity cost of the agreement (last two variables in Eq. 5.4). Moreover, to model a proxy for the debtor’s j survival probability, we apply a heuristic rule. Recalling that the borrower fails if her equity becomes negative, $E_t^j < 0$, the probability of surviving is simply given by the closeness between j ’s equity and the highest net-worth in the system, i.e.

$$p_t^j = \frac{E_t^j}{E_t^{\max}}. \quad (5.5)$$

Finally, the maximum amount that the lender i is willing to lend to j , that is the lending capacity, $c_t^{i,j}$, in Eq. 5.4 is defined as

$$\begin{cases} c_t^{i,j} = (1 - h_t^j)A_t^j > 0, & \text{if } (i, j) \in V_t, \\ c_t^{i,j} = 0 & \text{otherwise,} \end{cases}$$

with $h_t^j \in (0, h_t^{\max})$ to be the borrower haircut, defined as the j ’s leverage, λ_t^j , with respect to the maximum one. Hence $h_t^j = \frac{\lambda_t^j}{\lambda_t^{\max}}$, with $\lambda_t^j = \frac{L_t^j}{E_t^j}$. By setting Eq. 5.4 equal to zero and rewriting it as a function of $r_t^{i,j}$, we get the interbank rate that guarantees zero expected profit:

$$r_t^{i,j} = \frac{\chi A_t^i - \phi A_t^j - (1 - p_t^j)(\xi A_t^j - c_t^{i,j})}{p_t^j c_t^{i,j}}. \quad (5.6)$$

As the reader can verify, in line with the assumption of asymmetric information and costly state verification (see Bernanke et al., 1999), the lender

5. Reinforcement Learning Policy Recommendation for Interbank Network Stability

applies an interest rate that increases with her own size (that is, her assets) and the financial vulnerability of the borrower (that is j 's leverage). This last implication derives from the budget identity (see Eq. 5.1) from which we can derive that $A_t^j = \frac{L_t^j}{\lambda_t^j} + D_t^j$, where $\lambda_t^j = \frac{L_t^j}{E_t^j}$. In addition, the interest rate in Eq. 5.6 is not linearly related to the bank probability of surviving and capacity.

We now have all the elements to describe how traders select their counterpart in the interbank system, i.e. how lending arrangements are formed and evolve. We develop a measure of agent attractiveness to generate an endogenous mechanism of preferential attachment. Specifically, banks try to signal themselves in the interbank market by offering low interest rates or conspicuous supplies of liquidity. Although all agents start from the same initial conditions, as time goes by financial institutions are characterized by heterogeneous levels of their agent-specific variables. In line with this, the fitness of each agent μ_t^i is a combination between her liquidity relative to the highest liquidity provided in the market, C_t^{\max} , and her interest rate compared to the cheapest one, r_t^{\min} , i.e.

$$\mu_t^i = \eta_t \left(\frac{C_t^i}{C_t^{\max}} \right) + (1 - \eta_t) \left(\frac{r_t^{\min}}{r_t^i} \right). \quad (5.7)$$

The parameter η_t reflects a policy recommendation at time t , addressing the choice of the banking sector towards one of two possible strategies. On the one hand, η approaching zero identifies an interbank system moving towards the cheapest interest rates, on the other hand, η close to one highlights a liquidity-based system. We refer the reader to the Subsection Section 5.3.1 for a detailed explanation on the policy recommendation evolution: one of the main contribution of our work is to assume η_t endogenously evolving over time through a reinforcement learning mechanism, modeling the will of the regulator to address the banking system toward the best credit strategy for system stability.

Coming back to the interbank network, in our model credit links are directional because they are created and deleted by the agent j who looks for a loan and points to the agent i that provides credit. The loan flows in the opposite direction. In general local interaction models, the agent interacts directly with a finite number of counter-parties in the population. The set of nodes with whom a single node is linked is referred to as its neighborhoods. In our model, the number of out-going links is constrained to be a small number

5. Reinforcement Learning Policy Recommendation for Interbank Network Stability

\hat{d} , thus borrowers can only get loans from \hat{d} lenders. With this assumption of network sparsity the topology is always locally tree-like, avoiding loops that would preclude us from fully understanding the impact that the network architecture has on economic dynamics such as systemic risk, failures and liquidity diffusion.

At time $t = 0$, each bank j starts having \hat{d} random outgoing links (i.e. potential borrowing positions), and possibly with some incoming links from other agents (i.e. potential lending position). At the beginning of each period, links are rewired in the following way. For any outgoing link i , each borrower j randomly selects a new bank k . Comparing the fitness of the new financial institution with the one of its previous lender i , the borrower j cuts her old link with i and creates a new one with k according to the probability

$$P_t^j = \frac{1}{1 + e^{-\beta(\mu_t^k - \mu_t^i)}}, \quad (5.8)$$

or keep its previous link with probability $1 - P_t^j$. The proposed mechanism for reviewing credit agreements ensures that the most attractive lenders get the highest number of borrowers (i.e. incoming links) and, consequently, earn the highest profits. Nevertheless, the degree of randomness included in the algorithm guarantees that some links with very high performing agents may be cut in favor of less attractive creditors. The amount of randomness is regulated by β and has a double purpose: from a practical point of view, it prevents the system from being centralized around a single financial hub; from a theoretical perspective, it allows us to model incomplete information and bounded rationality.

The evolution of the banking system: determining the policy recommendation

As anticipated in the previous section, we use the reinforcement learning paradigm to move the parameter η_t and obtain an optimal policy recommendation in the described banking system. The aim of reinforcement learning is to solve a decision-making problem in which the timing of costs and benefits is relevant. In an interbank market that follows the specified dynamics for the creation of lending agreements, reinforcement learning can help in determining the policy recommendation that better identifies the optimal attachment strategy to follow in Eq. 5.7, even when partial information

5. Reinforcement Learning Policy Recommendation for Interbank Network Stability

about the system is provided. Hereafter, we refer to the reinforcement learning algorithm simply as the learning algorithm. A Markov Decision Process (MDP) is the mathematical formalism under which the reinforcement learning problem is usually defined. An MDP consists of a set of possible states $S_t \in \mathcal{S}$, a set of possible actions $A_t \in \mathcal{A}$ and a transition probability $P[S_{t+1} = s' \mid S_t = s, A_t = a]$. At each time t , a learning agent that is in state S_t , takes an action A_t and receives a reward $R_{t+1}(S_t, A_t, S_{t+1}) \in \mathbb{R}$ from the environment before moving to the next state S_{t+1} . We define the agent strategy $\pi : \mathcal{S} \times \mathcal{A} \mapsto [0, 1]$ as the conditional probability $\pi(a \mid s)$ of taking the action $A_t = a$ being in the state $S_t = s$. The reinforcement learning problem is the stochastic control problem of maximizing the expected discounted cumulative reward

$$\mathbb{E}_\pi \left[\sum_{t=0}^{\infty} \gamma^t R_{t+1}(S_t, A_t, S_{t+1}) \right], \quad (5.9)$$

where $\gamma \in [0, 1)$ is a discount factor and the expectation is w.r.t. the sequence of states and actions reached following the strategy π .

In our MDP, the sequential economy in which the banking system operates plays the role of the environment. Banks interact with the environment by changing their credit lines: each day they can adapt their attachment strategy between liquidity supply and interest rate discount, which is regulated through Eq. 5.7, with the choice of η_t , playing the role of the action A_t . We assume the agent is the system as a whole rather than the single bank and that the optimal strategy is realized at the system level, i.e. that the regulator directs financial institutions towards the correct combination of the two strategies. This assumption has a twofold purpose. On the one hand, it helps us to model a system with incomplete/asymmetric information, where the central bank has a richer information set than the single economic actor (see, for instance, Hoff and Stiglitz, 1990 and Thakor, 2020). On the other hand, it allows us to incorporate economic policy, seen as the optimal indication that the regulator gives to the system with the aim of reducing the interbank market vulnerability (see Trichet, 2010, for a global overview)². The state S_t includes information on both the liq-

²Considering η as a system variable allows us to reduce the mathematical and computational complexity of the problem and to study the behavior of the banking system as a whole. Making η bank specific leads towards multiagent reinforcement learning applications (Buşoniu et al., 2010) which consider agents that compete with each other and are out-of-the scope of the present paper.

5. Reinforcement Learning Policy Recommendation for Interbank Network Stability

uidity C_t and the interest rate r_t distributions of the banking system. In particular $S_t = (C_t^{\max}, C_t^{\min}, r_t^{\max}, C_t^{\text{avg}}, r_t^{\min}, r_t^{\text{avg}})$, where $x_t^{\max} = \max_{i \in \Omega} x_t^i$, $x_t^{\min} = \min_{i \in \Omega} x_t^i$, $x_t^{\text{avg}} = \sum_{i=1}^N x_t^i / N$, being x the variable of interest. We believe that this setting of the state space is realistic enough to model the regulator partial information about the banking system: it would be difficult and costly to retrieve detailed and specific data on all the banks included in the system at each time step. It is easier indeed to gather information about the best and the worst liquidity provider in the interbank network, as much as average estimates of the entire market.

Finally, the reward function we consider is the system’s total fitness

$$R_t(S_t, A_t, S_{t+1}) = \sum_{i=1}^N \mu_t^i \quad (5.10)$$

and the problem in Eq. 5.9 becomes a maximization of the discounted cumulative bank’s total fitness. From the definition of bank fitness, this means not only to guarantee a better flow of liquidity through the banking system, but also an efficient allocation at a more convenient interest rate. The learning algorithm operates in a model-free setting because it only receives partial information on the relevant variables of the system, while it has no knowledge of the internal dynamics (i.e. transition probability) with which the balance sheets of the banks moves and lending agreements are generated. This piece of information has to be inferred through the sequence of states, actions and rewards during the learning process.

5.2.3 The optimization algorithm: Proximal Policy Optimization

The optimization problem in Eq. 5.9 can be solved using a policy gradient algorithm like the Proximal Policy Optimization (PPO) (Schulman et al., 2017). A policy gradient algorithm directly parametrizes the optimal strategy $\pi_\theta = \pi(a \mid s; \theta)$, for example using a multilayer neural network with parameters θ . The optimization problem is approximately solved by computing the gradient of the cumulative fitness of the system $J(\theta) = \sum_{t=0}^{\infty} \gamma^t R_{t+1}(S_t, A_t, S_{t+1}; \pi_\theta)$ and then carrying out gradient ascent updates according to

$$\theta_{t+1} = \theta_t + \alpha \nabla_{\theta} J(\theta_t), \quad (5.11)$$

5. Reinforcement Learning Policy Recommendation for Interbank Network Stability

where α is a scalar learning rate. The policy gradient theorem (Marbach and Tsitsiklis, 2001; Sutton et al., 1999) provides an analytical expression for the gradient of $J(\theta)$ as

$$\begin{aligned}\nabla_{\theta} J(\theta) &= \mathbb{E}_{\pi_{\theta}} \left[\frac{\nabla_{\theta} \pi(A_t | S_t; \theta)}{\pi(A_t | S_t; \theta)} Q_{\pi_{\theta}}(S_t, A_t) \right] \\ &= \mathbb{E}_{\pi_{\theta}} [\nabla_{\theta} \log \pi(A_t | S_t; \theta) Q_{\pi_{\theta}}(S_t, A_t)],\end{aligned}\tag{5.12}$$

where the expectation, with respect to (S_t, A_t) , is taken along a trajectory (episode) that occurs adopting the strategy π_{θ} and the action-value function

$$Q_{\pi}(s, a) \equiv \mathbb{E} \left[\sum_{k=0}^{\infty} \rho^k R_{t+1+k} \mid S_t = s, A_t = a, \pi \right],\tag{5.13}$$

represents the long-term reward associated with the action a taken in the state s if the strategy π is followed hereafter. It can be proven that it is possible to modify the action value function $Q_{\pi}(s, a)$ in (5.12) by subtracting a baseline that reduces the variance of the empirical average along the episode, while keeping the mean unchanged. A popular baseline choice is the state-value function,

$$V_{\pi}(s) \equiv \mathbb{E} \left[\sum_{k=0}^{\infty} \rho^k R_{t+1+k} \mid S_t = s, \pi \right],\tag{5.14}$$

which reflects the long-term reward starting from the state s if the strategy π is adopted onwards. The gradient thus can be rewritten as

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\pi_{\theta}} [\nabla_{\theta} \log \pi(A_t | S_t; \theta) \mathbb{A}_{\pi_{\theta}}(S_t, A_t)]\tag{5.15}$$

where

$$\mathbb{A}_{\pi}(s, a) \equiv Q_{\pi}(s, a) - V_{\pi}(s),\tag{5.16}$$

is called advantage function and can be interpreted as the gain obtained by choosing a specific value of a in a given state with respect to its average value for the strategy π .

Different policy gradient algorithms derive from the way the advantage function is estimated. In PPO, the advantage estimator $\mathbb{A}(s, a; \psi)$ is parameterized by another neural network with parameters ψ . This approach is known as actor-critic: the actor is represented by the policy estimator

$\pi(a|s; \theta)$ that outputs a probability for each possible value of $a \in \mathcal{A}$, which the learning algorithm uses to sample actions, while the critic is the advantage function estimator $\mathbb{A}(s, a; \psi)$ whose output is a single scalar value. The two neural networks interact during the learning process: the critic drives the updates of the actor, which successively collects new sample sequences that will be used to update the critic and again evaluated by it for new updates. The PPO algorithm can therefore be described by the extended objective function

$$J^{\text{PPO}}(\theta, \psi) = J(\theta) - c_1 L^{\text{AF}}(\psi) + c_2 H(\pi(a | s; \theta)) \quad (5.17)$$

where the second term is a loss between the advantage function estimator $\mathbb{A}(s, a; \psi)$ and a target \mathbb{A}^{targ} , represented by the cumulative sum of discounted reward, needed to train the critic neural network, and the last term represents an entropy bonus to guarantee an adequate level of exploration. Details about the specific choice of the target together with additional information about the general algorithm implementation are given in the App Section 5.6. In what follows, PPO can be generally referred to as the learning algorithm.

5.3 Simulation Results

In this section, we perform numerical experiments to test the capability of the learning algorithm to identify an optimal strategy for selecting η and trading-off the two competing ways of establishing credit relationships. In this respect, we analyze the effects of the η dynamics on agents' economic performances, the interbank network topology and its resilience in the face of exogenous shocks. Finally, we study the effect of the policy recommendation obtained through reinforcement learning in controlling credit crunch phenomena and mitigating systemic risk.

The results provided in the following subsections are obtained from simulated tests which shares some choices for the parameter involved in the dynamic simulation of the system. The number of Monte Carlo simulations performed is $M = 200$ and each simulation is $T = 1000$ periods long. We simulate a system with $N = 50$ banks whose out-degree is $\hat{d} = 1$, so each bank can obtain at most one outgoing link at each time step, while can have many possible incoming links. Each bank is subjected to an initial probability of being isolated, which is set at 0.25. The parameters of the screening costs χ and ϕ that enters in Eq. 5.6 are set respectively at 0.015 and 0.025,

5. Reinforcement Learning Policy Recommendation for Interbank Network Stability

while the liquidation cost of collateral ξ is 0.3. The parameters μ and ω shifting the uniformly distributed noise that shocks the bank deposits are set at 0.7 and 0.55. All the banks starts with the same initial interest rate equal to 2% and are endowed with the same initial balance sheet $C_0 = 30$, $L_0 = 120$, $D_0 = 135$ and $E_0 = 15$. The price of fire sale $\rho = 0.3$ and the intensity for breaking the connection between banks $\beta = 5$ in Eq. 5.8 are other parameters common to all the agents in the network. In the App. Section 5.5 we check the robustness of our qualitative results by changing some key parameters. Specifically, we vary the intensity of choice, β , from 0 to 40 with steps of 2; the fire-sale price, ρ , from 0.1 to 0.5 with steps of 0.1 and, finally, the parameter ω regarding the volatility shock on bank deposit. We have then studied the moments of the distributions of the statistics of interest. Results confirm that our findings are robust to some variation of the banking system simulation.

The PPO algorithm parametrizes a discrete strategy function so that the learning algorithm can choose the value of η among a finite set of actions $\mathcal{A} = \{0, 1\}$ ³. The choice of a discrete strategy with respect to a continuous strategy on the interval $(0, 1)$ is motivated by the willing to learn a strategy that clearly discriminates between an interest rate strategy ($\eta = 0$) and a liquidity strategy ($\eta = 1$). Learning a continuous strategy that can pick any real value in the given interval is known to be more difficult and make it harder to interpret whether the learned dynamic strategy is pointing towards an objective rather than the other.

5.3.1 Training the PPO algorithm

As a first step in our numerical analysis, we evaluate the performance of the strategy learned by the PPO algorithm. We train four PPO instances on $E_{in} = 1000$ consecutive episodes, which are independent simulations of the banking system. The PPO instances differ for the random seed used to initialize the neural networks and to train them using a stochastic gradient descent approach. Multiple concurrent training of different instances is needed in order to provide an average performance together with a confi-

³Under the same setting, we also trained PPO instances that are allowed to pick discrete values between 0 and 1 as a possible action. The distributions of the η provided by these discrete strategies present a negligible amount of non-extreme values in the interval $(0, 1)$ and are similar to those visualized in the left-hand side of Figure 5.1. We resort then to the binary case, which is more interesting for our analysis.

5. Reinforcement Learning Policy Recommendation for Interbank Network Stability

dence interval that highlights the robustness of the learning process. Each training episode consists of a simulation of the banking system for T periods that allows the learning algorithm to collect samples of data with which it can perform updates of the model parameters. During the learning phase, we evaluate the learning progress of each instance at several intermediate steps. We fix the weights of the neural networks that parametrize the η public signal and perform $E_{out} = 5$ out-of-sample test episodes before carrying on the training process to assess the learned behavior up to that point. We refer to the App. Section 5.6 for the technical difference between an in-sample and an out-of-sample test episode.

We compare the PPO performance with respect to a dynamic random baseline that picks the value of η according to a Bernoulli distribution with parameter equal to 0.5. This random policy that chooses between 0 and 1 with equal probability represents a meaningful benchmark, as we observe in the left-hand side of Figure 5.1, where the values of η in both scenarios are identically distributed over the M performed Monte Carlo simulations. The Kolmogorov-Smirnoff test statistically confirms up to the 1% confidence level that the distribution of the η values generated by the selected⁴ PPO instance is not significantly different from the one of the random baseline. The right-hand side of Figure 5.1 summarizes the results of the learning process, where the average cumulative fitness of the system in Eq. 5.10 is represented on the y -axis. Every PPO instance is tested $E_{out} = 5$ times using Monte Carlo simulations of length T . We notice that the performance metric is always greater for PPO than for the random recommendation, signaling that the banks in the system generated by the PPO signal tend to be more attractive for the borrowers by exhibiting a higher aggregated fitness through time. Moving η randomly causes banks to be less attractive for the borrowers in their interbank market. This result implies that the PPO instances learn to choose the value of η by leveraging the information available about the system, without changing the distribution of the values with respect to the random case. The learning procedure allows us to discover when it is exactly convenient to pick a side in this trade-off. A further comparison with some fixed signals is provided in the App. Section 5.5. However, fixing the η for all time steps has an evolutionary impact on the system which has already been

⁴It is common in reinforcement learning applications to train different instances of the same algorithm and then select the best performing one over some out-of-sample tests (Andrychowicz et al., 2020)

5. Reinforcement Learning Policy Recommendation for Interbank Network Stability

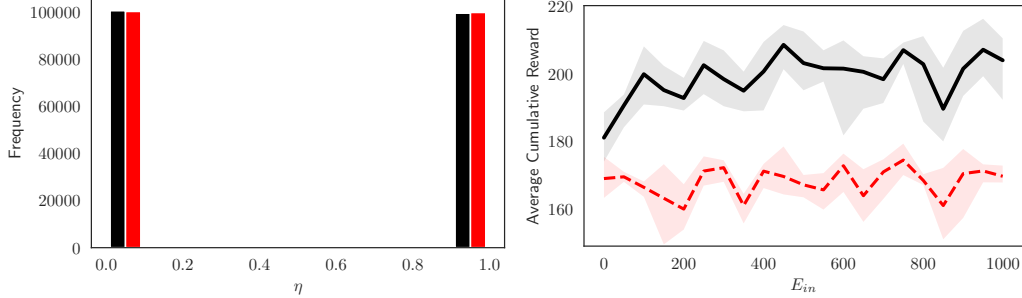


Figure 5.1: The left panel shows the discrete distributions of the η values selected respectively by PPO (in black) and by a Bernoulli distribution (in red) with parameter equal to 0.5 over M Monte Carlo simulations of the system. The right panel shows the average cumulative fitness of the system as a function of the number of training episodes for the trained PPO instances (in solid black) and the Bernoulli distribution of η (in dashed red) with the corresponding confidence intervals.

studied by Berardi and Tedeschi, 2017 and it is not centered on studying the effect of any η that changes through time.

In order to shed light on the decisions taken by the best performing trained PPO instance, we use the SHapley Additive exPlanation (SHAP)⁵ framework (see Lundberg and Lee, 2017, Shapley, 2016). This approach allows explaining a complex nonlinear model like a neural network by shedding light on the contribution of each input feature to the output formation. For each input vector $x \in \mathbb{R}^K$ and a model f , the SHAP value $\phi_i(f, x)$, $i = 1, \dots, K$ quantifies the effect (in a sense, the importance) on the output $f(x)$ of the i -th feature. To compute this effect one measures, for any subset $S \subseteq \{1, \dots, K\}$, the effect of adding/removing the i -th feature to the set, i.e. $f_{S \cup \{i\}}(x) - f_S(x)$. The SHAP value is defined as the weighted average

$$\phi_i(f, x) = \sum_{S \subseteq \{1, \dots, K\} \setminus \{i\}} \frac{|S|!(K - |S| - 1)!}{K!} [f_{S \cup \{i\}}(x) - f_S(x)], \quad (5.18)$$

where the weights ensure that $\sum_i \phi_i = f(x)$.

Figure 5.2 shows the magnitude of the Shapley values for the policy recommendation learned by the best performing PPO instance, referred to the two possible outcomes $\eta = 0$ and $\eta = 1$. The left-hand side shows that high

⁵For the implementation we use the Python package linked to Lundberg and Lee, 2017

5. Reinforcement Learning Policy Recommendation for Interbank Network Stability

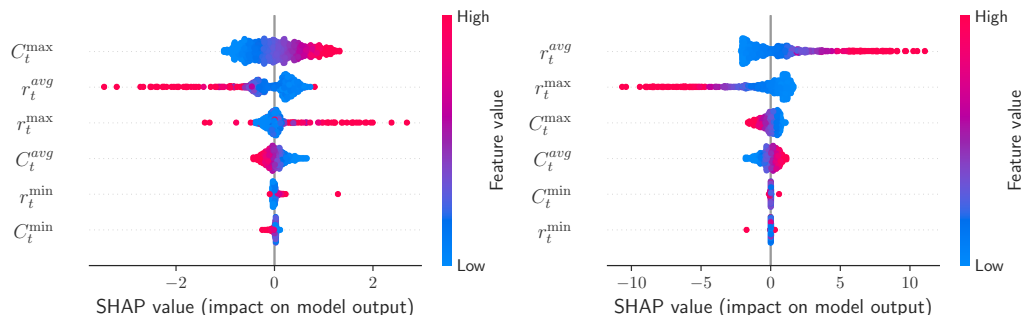


Figure 5.2: SHAP values relative to the strategy outputs $\eta = 0$ (left panel) and $\eta = 1$ (right panel). The cloud of colored dots for each input variable expresses the importance and the correlation with respect to the model output. Features are ordered on the y -axis by relevance, so the first on the top influences the most the model output.

values for the maximum liquidity available in the system tends to favor the choice of a η based on the interest rate. Also, a low average interest rate and a high maximum interest rate point to the choice of $\eta = 0$. The right-hand side shows an opposite input relevance, with a dominant role for high values of the average interest rate and low values of the maximum interest rate. The two figures show that the trained learning algorithm chooses one of the two signals by looking at the main characteristics of the opposite one. When it chooses $\eta = 0$, it is more interested to know if there are participants in the network which are large, while when it chooses $\eta = 1$, it looks for homogeneity of interest rate, that is a common feature obtained by always playing towards the interest rate. The learning algorithm suggests a switch towards the other competing recommendation in order to avoid extreme cases in which a disadvantage of one or the other choice exacerbates, i.e. a huge financial institution that gather all the demand of the borrowers when $\eta = 1$ could not be sustainable in the long term, so the algorithms suggests switching to the other option. On the other hand, a majority of banks of medium size that offer medium rates when $\eta = 0$ could not be able to gather enough liquidity to deal with deposit shocks and would be better to resort to the opposite signal.

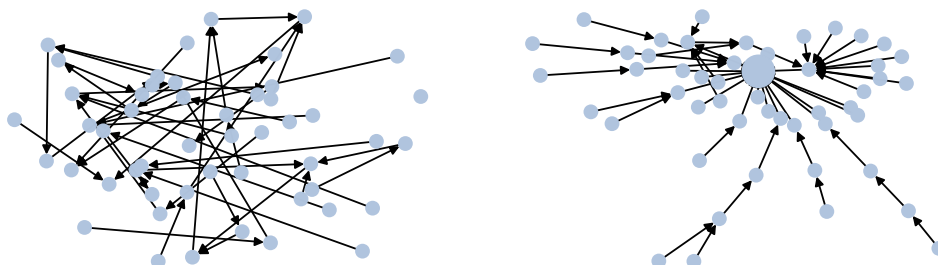


Figure 5.3: Network configuration at time $t=0$ (left side), and $t=800$ (right side).

5.3.2 Micro and macro consequences of the policy recommendation

In this subsection, we deal with the implications that the dynamics of the η parameter has on the interbank network morphology and on resulting performances of the financial institutions. Finally, we study the effects that the emerging network topology has on the stability of the market.

Topology and evolution of the interbank network

Before starting the analysis, it is worth remembering the dynamics of η , that appears in the banks' fitness (see Eq. 5.7), determines the probability of creating credit links in the system as shown in Eq. 5.8. Therefore, it is appropriate to begin the analysis by describing the topology of the interbank network.

In Figure 5.3 we plot the configuration of the endogenous interbank network at two different time steps of a single simulation of the system. As the reader can appreciate, the market configuration goes through different phases, ranging from a random topology with isolated agents to a highly centralized architecture where a few hubs compete for credit supply. A more detailed analysis on the evolution of the interbank network architecture over time can be found in the left-hand side of Figure 5.4, where we show the time series of network degree centrality

$$C_t^{\text{Net}} = \frac{\sum_i (k_t^{\text{max}} - k_t^i)}{N(N-1) - |V_t|}, \quad (5.19)$$

5. Reinforcement Learning Policy Recommendation for Interbank Network Stability

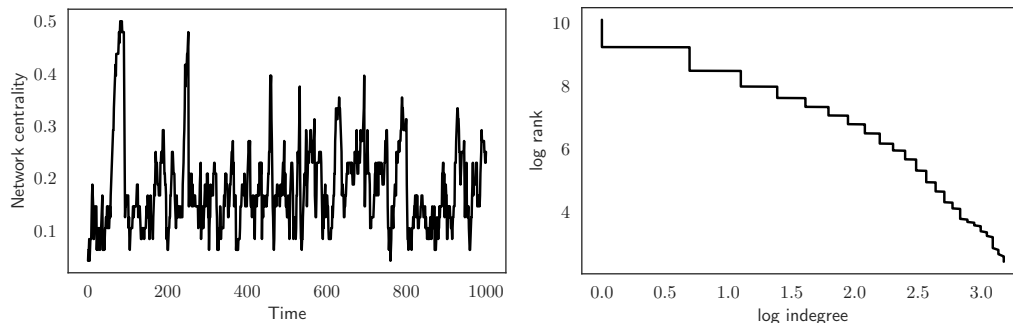
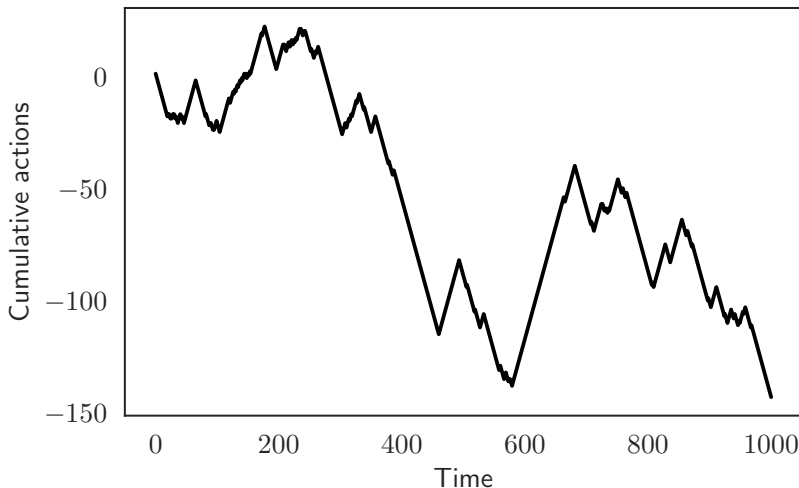


Figure 5.4: Time series of interbank network centrality (left side). The decumulative distribution (DDF) of the in-degree (right side).

where N is the number of banks, $|V_t|$ is the total number of incoming links in the system, k_t^i is the number of incoming links for the i -th bank and k_t^{\max} is the number of incoming links hold by the hub of the network.

The dynamics of network centrality shows how the morphology of the credit market evolves over time, going from periods in which the network is decentralized and made of many small components to periods in which more than 45% of banks are connected to a single hub. In addition, the topology of the emerging network as a whole is different from that of the random graph, where the in-degree distribution decays exponentially. Similar to real credit networks, in our system some banks are found to have a disproportionately large number of incoming links while others have very few (see Iori and Mantegna, 2018, for a survey of the relevant literature). This result is shown in the right-hand side of Figure 5.4 where we plot the decumulative distribution function of the in-degree. As the reader can observe, this distribution is in keeping with that of scale-free networks and displays a 'fat tail'.

5. Reinforcement Learning Policy Recommendation for Interbank Network Stability



y_t	b_0	b_1
Centrality	0.1830 ^{*****} (599.06)	-0.0047 ^{*****} (-11.8013)
Density	0.1042 ^{*****} (317.08)	-0.0070 ^{*****} (-16.01)
Diameter	10.14 ^{*****} (1104.02)	0.22 ^{*****} (17.10)
Components	1.48 ^{*****} (656.45)	0.0095 ^{*****} (3.04)
Avg nodes per components	39.50 ^{*****} (940.20)	-0.21 ^{*****} (-3.57)

*** $p < 0.01$, ** $p < 0.05$, * $p < 0.1$

Figure 5.5: Top Panel: Time series of η cumulative values over the simulation. Bottom Panel: Estimated results with the respective T-test in brackets for Eq. 5.20. b_0 is the estimated mean value of y when $\eta = 1$ and b_1 the deviation from this mean value when $\eta = 0$. Data are obtained through M Monte Carlo simulations of the system.

To conclude the analysis on the interbank market architecture, we deal with the effect of the η parameter on the credit network topology. In the top panel of Figure 5.5 we plot a single realization of the cumulative value

5. Reinforcement Learning Policy Recommendation for Interbank Network Stability

of η through time. The figure shows how the reinforcement learning algorithm generates a time evolution in the choice of policy recommendations. Specifically, increasing (decreasing) values in the curve correspond to a signal that directs the system towards high liquidity supply (low interest rate), i.e. $\eta = 1$ (i.e. $\eta = 0$).

The effect of the signal in shaping the topology of the interbank network is, instead, shown in the lower panel of Figure 5.5, where we estimate a categorical regression model

$$y_t = b_0 + b_1(1 - \eta_t), \quad (5.20)$$

where b_0 is the estimated mean value assumed by the dependent variable y when $\eta = 1$ and $b_0 + b_1$ is the mean when $\eta = 0$. As shown in the bottom panel of Figure 5.5, when the system selects low interest rates, the interbank network is less centralized, more sparse and with a larger diameter. Moreover, the graph is fragmented into many scarcely-populated islands.

Having described the architecture of the interbank network, let us now examine its evolution over time. It is worth remembering that banks signal in the market their attractiveness μ according to the recommendation received from the regulator, i.e. whether to compete more on low interest rates, $\eta = 0$, or on high liquidity supply, $\eta = 1$. While the regulator's signal is market specific, liquidity supplies and interest rates (based on Eq.6) are bank specific variables. This mechanism creates competition among financial institutions for credit allocation. The war in the granting of credit, modeled through the possibility of redefining lending agreements via Eq. 5.8 is shown in Figure 5.6.

5. Reinforcement Learning Policy Recommendation for Interbank Network Stability

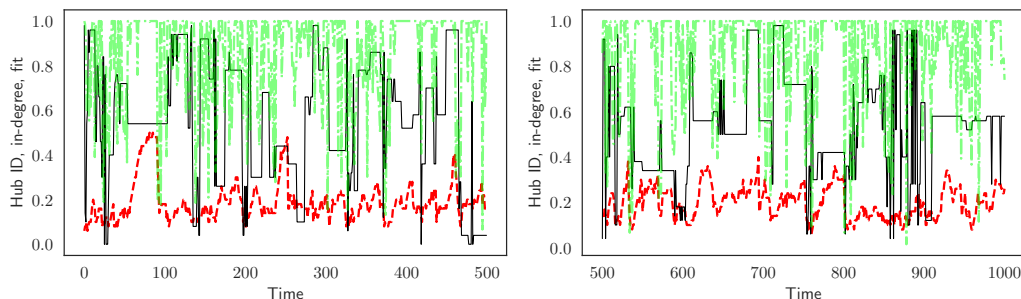


Figure 5.6: Time series of the evolution of the most connected lender (hub) along the time T . The black solid line identifies the normalized hub id, the red dashed line her number of clients (incoming links) and the green dotted line the hub's fitness. Colors are available on the web side version.

The black solid, red dashed and green dotted lines represent the normalized id of the lender with the highest number of clients (i.e. the hub), her incoming links (i.e. number of clients) and her fitness, respectively. As the reader can appreciate, the simulation presents periods of hub stability and periods of alternation and competition between different hubs. When the hub succeeds in standing out from her competitors and signaling a significantly higher fitness (i.e. the green dotted line approaches the unit), she is able to attract numerous clients as shown by her high number of incoming links. However, the attractiveness of the hub may work against her. A large portfolio of customers increases the likelihood that some of them may fail. This either decreases the attractiveness of the hub herself⁶ or even causes her failure. In any case, the drop in the agent's fitness gradually reduces her number of clients and makes other lenders more attractive. These agents can replace the unsuccessful hub and so become in turn the most appealing lenders.

⁶The reduction of the hub's fitness due to one of her clients' failure works in the following way. On the one hand, when the fitness uses a strategy based on a low interest rate, the client's approach to the bankruptcy threshold increases both the borrower's financial fragility and probability of bankruptcy. Both these effects produce an increase in the lending interest rate, which makes the hub less attractive (see Eq. 5.6). On the other hand, when μ moves towards a high liquidity supply, the borrower's bad debt is absorbed by the lender's net-worth. The fall in the latter causes a parallel reduction in the hub liquidity, as shown by the balance sheet identity (see Eq. 5.1).

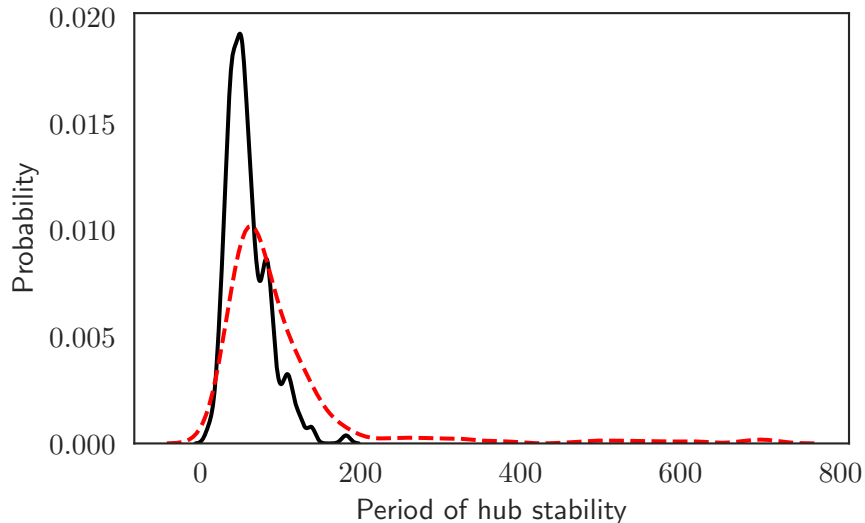


Figure 5.7: Density distributions over M Monte Carlo simulations of maximum period of hub stability in which the strategy doesn't change. The black solid and red dashed lines show $\eta = 0$ and $\eta = 1$, respectively.

Micro consequences of the reinforcement learning policy

In this subsection, we investigate how the dynamics of η affects the performance of the hub and other financial institutions.

In Figure 5.7 we show how the choice between a low interest rate and a high liquidity supply strategy affects the hub longevity. The figure shows the distribution, over M simulations, of the maximum period of hub stability in which the strategy doesn't change, respectively for $\eta = 0$ (black) or $\eta = 1$ (red). As the figure shows, the hub is in general more stable if the regulator recommends a high liquidity supply (red dashed line in Fig Figure 5.7). Moreover, also at a micro level, we show that $\eta = 1$ seems to produce better individual performances. This result is shown in the top panel of Figure 5.8, where we report the effect of the two possible values of η on some key individual variables.

5. Reinforcement Learning Policy Recommendation for Interbank Network Stability

y_t	b_0	b_1
Liquidity	2960.34 ^{*****} (1062.34)	331.42 ^{*****} (82.59)
Equity	888.96 ^{*****} (1171.31)	-110.72 ^{*****} (-135.79)
Leverage	0.01673 ^{*****} (435.43)	0.000175 ^{*****} (3.23)
Rationing	0.33 ^{*****} (71.80)	0.28 ^{*****} (38.31)
Bad debt	36.05 ^{*****} (446.46)	2.19 ^{*****} (19.49)
Failed banks	3.14 ^{*****} (520.17)	0.35 ^{*****} (41.00)

*** $p < 0.01$, ** $p < 0.05$, * $p < 0.1$

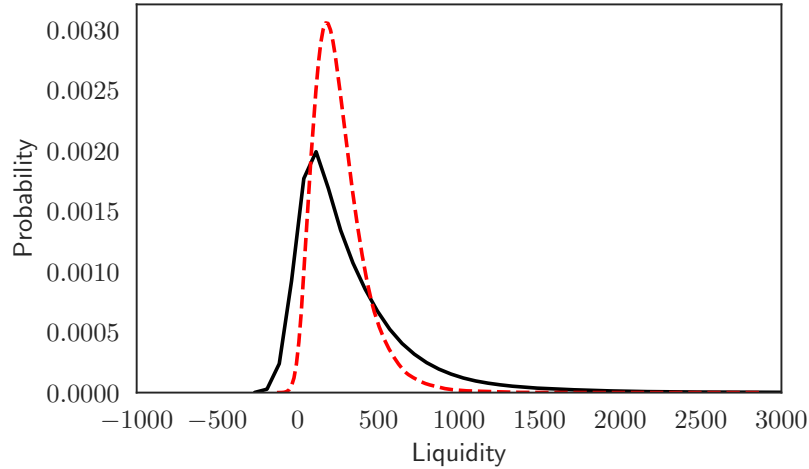


Figure 5.8: Top Panel: Estimated results with the respective T-test in brackets for Eq. 5.20. b_0 is the estimated mean value of y when $\eta = 1$ and b_1 the deviation from this mean value when $\eta = 0$. Data are obtained through M Monte Carlo simulations of the system. Bottom Panel: Density distributions of aggregated liquidity over times and over M Monte Carlo simulations. The black solid and red dashed lines show $\eta = 0$ and $\eta = 1$, respectively.

Specifically our results, estimated via the categorical regression model in Eq. 5.20, shows that a signal that directs the system towards an abundant supply of liquidity (i.e. $\eta = 1$) produces better results on controlling leverage, rationing, bad debt and bankruptcies. Moreover, in line with the hypothesis that the failure of banks occurs as net-worth falls below a minimum threshold, the equity is higher in the case of $\eta = 1$.

The result on the liquidity is, however, less intuitive. In fact, the system that competes on the interest rate level is significantly more liquid than the one adopting a high liquidity, with an average liquidity value of 3291 in the case of $\eta = 0$ and 2960 in the opposite case. The reason for the apparent better performance on liquidity in the case of $\eta = 0$ lies in the competition that arises among banks when they use interest rates. As clarified by Eq. 5.6, the financial institutions applying the lowest interest rates are the smallest ones. This implies that the biggest banks are less attractive to borrowers because they charge higher rates. The system, therefore, excludes these economic agents from trading, while it encourages small institutions to provide liquidity. This mechanism of selection has a twofold effect. On the one hand, it generates a strong heterogeneity between lenders and borrowers. Creditors, which are much smaller than debtors, are overwhelmed in the event of their clients' bankruptcy. On the other hand, the exclusion from the exchanges of the largest institutions leaves a lot of unallocated liquidity in the system. The first effect, i.e. agents' heterogeneity, determines the worst performances under $\eta = 0$, while the second effect, i.e. exclusion, determines the highest level of unallocated liquidity in the system. In contrast, a signal that directs the system towards an abundant liquidity supply produces a more homogeneous distribution among the liquidity of banks, as shown in the bottom panel of Figure 5.8. This homogeneity between economic agents generates a uniform risk exposure among counter-parties, which favors the resiliency of the system in front of shocks. This result is in line with other studies showing that agents' heterogeneity is a leading force in generating propagation of systematic failure (see, for instance, Caccioli et al., 2012, Berardi and Tedeschi, 2017, Iori et al., 2006, Lenzu and Tedeschi, 2012 and Tedeschi et al., 2012).

Systemic impact of the network

To conclude the section, we combine the results on network topology and individual performance as a function of η in order to capture the overall effect of the interbank architecture on systemic stability. To this end, in Tab.

5. Reinforcement Learning Policy Recommendation for Interbank Network Stability

Indep. Variable	Dep. Variable		
	Rationing	Failed banks	Leverage
Net centrality	-0.25 ^{***} (-6.04)	-2.09 ^{***} (-41.82)	-0.016 ^{***} (-55.69)
Density	-1.32 ^{***} (-52.64)	-9.14 ^{***} (-254.08)	-0.051 ^{***} (-235.33)
Diameter	0.011 ^{***} (8.95)	0.032 ^{***} (21.69)	0.0002 ^{***} (27.02)
Components	0.029 ^{***} (5.28)	0.020 ^{***} (3.09)	0.0004 ^{***} (10.57)
Avg nodes per comp	-0.0011 ^{***} (-4.18)	-0.0022 ^{***} (-6.97)	-0.00002 ^{***} (-13.19)

*** $p < 0.01$, ** $p < 0.05$, * $p < 0.1$

Table 5.1: Regression results between indicators of the interbank stability and network measures. T-stats for each coefficient are provided in parentheses. Data are obtained through M Monte Carlo simulations of the system.

Table 5.1 we report the results of a linear regression estimated through ordinary least squares where the independent variables are some measures of the interbank network topology and dependent variables are some indicators of the market systemic stability. In line with what has been observed so far, when the network tends to be centralized, i.e. denser towards the hub and with a smaller diameter, the risk of contagion decreases, i.e. bankruptcies, rationing and leverage are reduced. Obviously, this architecture corresponds to a graph composed of a few highly populated components. It is worth noting that this topology emerges when the interbank system is oriented towards an abundant supply of liquidity, which generates a certain homogeneity among agents able to compensate for the imbalance between lenders and borrowers present in the case of $\eta = 0$. In this respect, a clarification is important: $\eta = 1$ is not the absolute best signal. This is the best strategy given the individual and aggregate conditions of the system at the time of the choice. In fact, the algorithm is designed to identify one or the other recommendation as optimal on the basis of the underlying environmental conditions. The robustness of this observation is shown in Sec. Section 5.3.3 and in the App.

Section 5.5. In the former, we show that the system governed by a regulator that directs the choice via the implemented reinforcement learning algorithm outperforms a system based on a random selection between the two signals. In the latter, we demonstrate the better performances of the reinforcement learning with respect to keeping constant the two values of η .

5.3.3 The reinforcement learning based recommendation for taming systemic risk

In this subsection, we study the effect on the interbank systemic stability of the policy recommendation obtained through the reinforcement learning mechanism solved by the PPO algorithm.

Specifically, we answer the following question: how would the interbank system perform in terms of aggregate resiliency when the regulator direct financial institutions as a whole to choose the optimal strategy between competing on low interest rate, $\eta = 0$, or on high liquidity, $\eta = 1$? Again, we compare the effects of the learned strategy on the market stability with those of a random strategy.

A common finding in several theoretical and empirical works is that the interbank market works properly when credit flows efficiently through the system, thus ensuring it against liquidity shocks (see, for instance, Allen and Gale, 2000; Carlin et al., 2007; Freixas et al., 2000).

5. Reinforcement Learning Policy Recommendation for Interbank Network Stability

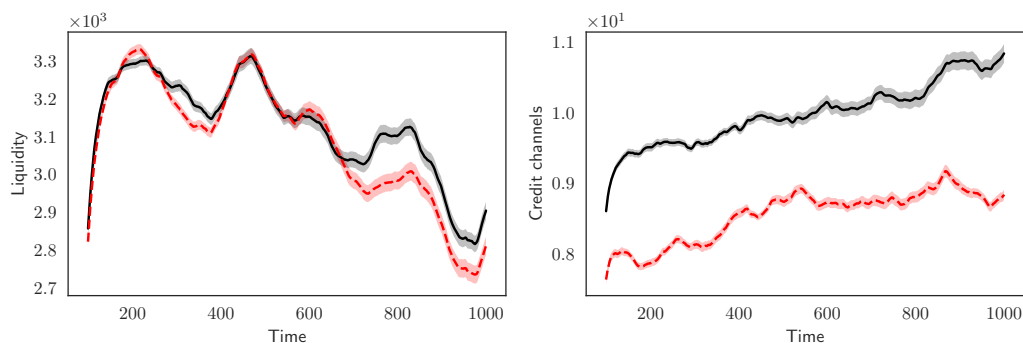


Figure 5.9: Liquidity of the system (left panel) and number of credit channels (right panel). Black solid and red dashed lines refer to the best performing reinforcement learning optimal strategy and to the random strategy, respectively. The curves reproduce the mean and the standard deviation over M simulations of the system and a rolling window of 100 time steps.

Starting from this consideration and recalling the severity of liquidity crises, we show in Figure 5.9 (left side) the effectiveness of the implemented reinforcement learning strategy in spreading liquidity through the system. In the figure, once selected the best performing learned strategy as shown in the right-hand side of Figure 5.1, the aggregated average liquidity of M simulations over a rolling windows of 100 time steps is shown through time. Although in some time periods the learned strategy strongly competes with the random one, its supremacy becomes evident from the time step 700 on-wards. In addition, the average liquidity, over all time periods and simulations, of the learned strategy is statistically higher than the one obtained with the random strategy (i.e. 3129.98 (std. 1.5128) vs 3091.51 (std. 4.4258), respectively).

A possible explanation for this phenomenon can be seen in the right-hand side of Figure 5.9, where we plot the active credit links in the two frameworks. As the reader can appreciate, the number of activated credit channels is higher when the system follows the learned strategy with respect to the case of random strategy, and this guarantees a higher circulation of liquidity in the system. In detail, the average number of credit channels, over time and simulations, in the first scenario is 9.9823 (std. 0.4321), while in the second case is 8.5464 (std. 0.3596). On the whole, this result reveals the ability of the reinforcement learning optimal policy to design an interbank network

5. Reinforcement Learning Policy Recommendation for Interbank Network Stability

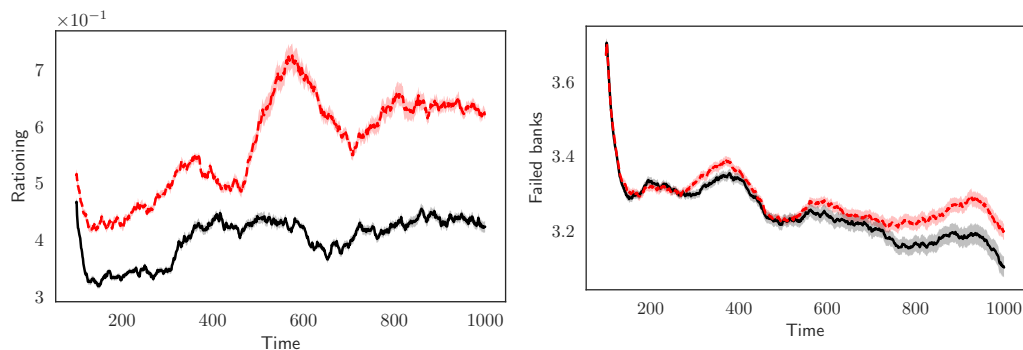


Figure 5.10: Rationing of the system (left panel) and number of failed banks (right panel). Black solid and red dashed lines refer to the best performing reinforcement learning optimal strategy and to the random strategy, respectively. The curves reproduce the mean and the standard deviation over M simulations of the system and a rolling window of 100 time steps.

architecture promoting an efficient credit allocation and, therefore, reducing liquidity shortage phenomena. As a consequence, the emerging topology of the credit network effectively controls rationing and avoids failures due to credit crunch phenomena, as shown in Figure 5.10, left and right panel, respectively.

The robustness of the two latter results is confirmed by the average values of these variables over all times and simulations. Specifically, the mean and standard deviation of the rationing in the case of the learned strategy (resp. random strategy) are 0.4024 and 0.0375 (resp. 0.5671 and 0.08465), while the mean and standard deviation of the number of failed banks in the case of the learned policy (resp. random policy) are 3.2101 and 0.0410 (resp. 3.2931 and 0.0423).

It is important to note the ability of the reinforcement learning mechanism to generate an interbank network whose architecture is resilient in the face of financial attacks. This characteristic provides on the one hand, an additional monetary policy tool that can be implemented in times of economic adversity, and on the other hand, enriches the vast literature that emphasizes the importance of credit network architecture in dealing with systemic shocks (see Grilli et al., 2017, for a survey of the relevant literature).

We conclude this section by analyzing the effect of the reinforcement learning optimal policy on the market financial (in)stability. The approach followed here in explaining the materialization of financial frictions is very

5. Reinforcement Learning Policy Recommendation for Interbank Network Stability

close in spirit to the Minskyan financial instability hypothesis, and therefore uses the leverage of banks as the main indicator (see Minsky, 1964). In our stylized market, the link between leverage and systemic instability works as follows. Given our naive banks' balance-sheet (see Eq. 5.1), the leverage is defined as assets on equity. Moreover, credit costs (i.e. interest rates) are strongly positively affected by the leverage (see Eq. 5.6). When a lender grants a loan to a bank with a low probability of surviving (i.e. an over-leveraged borrower) she charges higher interest rate via the financial accelerator. This, in turn, exacerbates the financial condition of the borrower herself, pushing her towards the bankruptcy state. If one or more borrowers are not able to pay back their loans, even the lenders' equity is affected by bad debts. Therefore, lenders decrease their credit supply and increase the borrowers' rationing. In this way, the profit margin of borrowers decreases and a new round of failures may occur. The leverage dynamics when the system follows the reinforcement learning recommended policy and in the random case are shown in the left-hand side of Figure 5.11. The figure highlights two important features. First, the recommended learned policy keep the leverage below the values obtained with the random policy. Specifically, the average leverage, over time and simulations, in the first scenario is 1.59 (std. 0.042), while in the second case is 1.69 (std. 0.031). Second, the leverage fluctuates over time, thus recalling the different phases of lending suggested by Minsky. In fact, there are periods when financial institutions grant more loans without considering the overall financial fragility. However, it can happen that banks underestimate their credit risk, making the system more vulnerable when default materializes. This ambiguous effect of the leverage, first positive and then negative, on interbank stability, is clearly shown in the right-hand side of Figure 5.11, where the correlation wave between bankruptcies and agents' leverage first decreases from lag $\tau = -21$ up to $\tau = -11$, then increases from $\tau = -8$ up to $\tau = 9$, and finally, returns to decrease from $\tau = 15$.

5. Reinforcement Learning Policy Recommendation for Interbank Network Stability

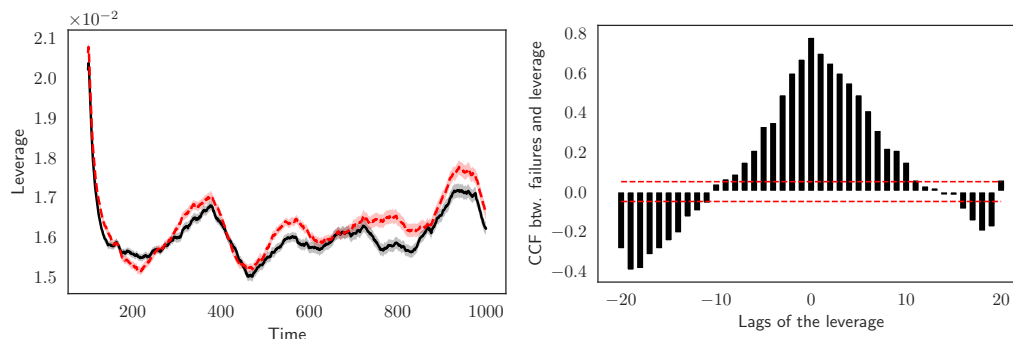


Figure 5.11: Left side: Leverage of the system. Black solid and red dashed lines refer to the best performing reinforcement learning optimal strategy and to the random strategy, respectively. The curves reproduce the mean and the standard deviation over M simulations of the system and a rolling window of 100 time steps. Right side: Average correlation between number of bankruptcies and lagged leverage, at a 1% confidence level.

5.4 Concluding remarks

In this work, we have shown the effects of a policy recommendation obtained through a reinforcement learning mechanism in an artificial interbank market. Specifically, we assume that the financial institutions receive a signal from the regulator regarding the best strategy to adopt for the creation of their lending agreements. Depending on the underlying economic conditions, the signal directs the system towards the provision of a high liquidity supply or a low interest rate. The use of a reinforcement learning approach to provide this public signal has proven to be effective since the method exploits the available information and redirects the system towards an efficient flow of liquidity, when compared to other different static and dynamic policies. Moreover, through the use of the SHAP framework, that dissects the contribution of each piece of information to the recommended policy, we have been able to interpret what is the main input that drives the choice of the policy. We have acknowledged that the occurrence of one or the other circumstance (liquidity vs interest rate) generates important consequences affecting either the agents' performances and either the topology and resiliency of the interbank network. Specifically, when the signal directs the system towards an abundant liquidity provision, the interbank network, composed by a few populated communities, is more centralized and dense towards hub banks than in

the low interest rate scenario. This network architecture is accompanied by better individual performances and higher resilience of the system in the face of exogenous shocks. Our results have shown that the better general conditions underlying this signal are due to the homogeneity between lenders and borrowers, which generates a uniform risk exposure among counter-parties able to favor the resiliency of the system.

Leaving aside the results on the comparison between the two signals, in the second part of the paper we have analyzed the general effect of the policy recommendation implemented via the reinforcement learning procedure. Our results have shown how systemic risk is mitigated by such a tool and how this outperforms other alternative policy instruments.

5.5 Appendix: A sensitivity analysis on model parameters

In this appendix, we investigate the performances of the learning algorithm by varying some key parameters. The first investigated parameter, β , governs the network topology (see Grilli et al., 2014, for a mathematical explanation). As the intensity of choice increases, the interbank architecture ranges from a random configuration to a star one. The effect of the network topology on the interbank system is studied by changing β from 0 to 40 with steps of 2. The second parameter we consider is fire sale price ρ . An increase in ρ has an impact on both lenders and borrowers. On the one hand, it compensates the losses that lenders incur due to the failure of their clients (see Eq. 5.4). On the other hand, a higher fire-sale increases the likelihood that the borrower, rationed in the interbank market, can face deposit repayments. Here we vary the fire-sale price, ρ , from 0.1 to 0.5 with steps of 0.1. Thirdly, we modify the skewness of the distribution of the random shock affecting the bank deposit at the beginning of each period. Recalling the equation for the deposit movements as $D_t^i = D_{t-1}^i(\mu + \omega U(0, 1))$, we remark that it allows us to reproduce bearish and bullish market periods. The uniformly distributed noise component can be shifted towards more negative or positive shocks at convenience in order to represent different market situations. Having fixed $\mu = 0.7$ in our simulations, we let ω vary from 0.52 to 0.6 with steps of 0.02, which corresponds respectively to a highly negative skewed and to a perfectly symmetrical shock distribution. This specific choice of the noise structure

5. Reinforcement Learning Policy Recommendation for Interbank Network Stability

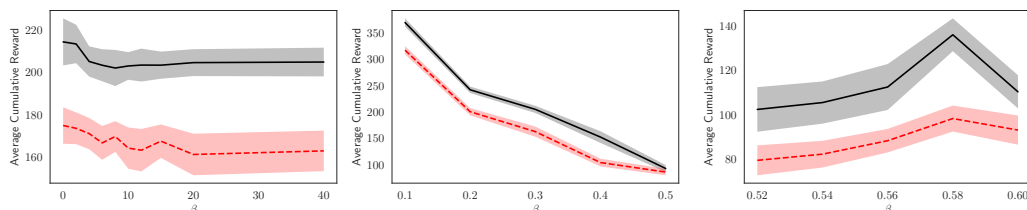


Figure 5.12: Average cumulative fitness of the system as a function of changes in β , ρ and ω , in the first, second and third panel, respectively. The reinforcement learning algorithm is in solid black, while the random strategy is in dashed red.

allows us to perturb the system with shocks that are mainly negative and generate situations where more banks ask for cash in the interbank market.

The last part of this appendix is dedicated to the comparison between η evolving via the reinforcement learning mechanism and the two fixed strategies, i.e. a constant η equal to 0 and 1, respectively.

In all these experiments, we run our model 100 times for different values of the initial seed, generating the pseudo-random numbers over a time span of $T = 1000$ periods. Moreover, all the agents' initialization parameters, except for the variations studied here, coincide with those presented in Sec. Section 5.3.

Let us begin the analysis by focusing on the implications that the three parameters variations have on the model's results. Each variation of a parameter represents a different configuration of the banking system, which is used to test the different strategies over M simulation. The cumulative reward of these simulations is then averaged to obtain the mean values and the respective confidence interval for the reinforcement learning strategy and the random strategy. Figure 5.12 shows the average cumulative reward over the M simulation as a function of a single parameter variation. We notice that the performance of the reinforcement learning algorithm solved with the PPO procedure is still superior with respect to the random strategy for all the three sensitivity cases presented. Therefore, we can conclude that the effect analysis in the main paper still holds if one modifies some characteristics of the underlying financial system.

In Figure 5.13 we show the sensitivity of the average values, over all the M simulations and a rolling window of 100 time steps, of relevant quantities

5. Reinforcement Learning Policy Recommendation for Interbank Network Stability

at the systemic level with respect to the three parameters described above⁷. The reinforcement learning strategy consistently outperforms the random strategy over all parameters and variables considered.

In the first column of Figure 5.13, we show the effects that the intensity of choice, β , has on the systemic variables. When β increases from 0 to 40, the liquidity and the credit channels increase up to $\beta = 10$ and then stabilize. The underlying reason for this dynamic is as follows: a β value greater than or equal to 10 generates a stable topology in the interbank network, which makes the investigated values insensitive to further changes in the parameter. Similar to the trend of the previous variables is the dynamics of the leverage, which increases with β but at a decreasing rate. Indeed, the more liquidity is available in the system, the more exchange of loans between banks happens. Finally, an increasing β causes the amount of rationing of the system to decrease, while the failures of the agent happen to be stable over the period. In the second column of the Figure 5.13, we focus on the effects produced by a variation in the fire-sale price. An increase of ρ protects both lenders and borrowers from losses, and it is beneficial when looking at the liquidity up to $\rho = 0.3$. From that level, borrowers do not enter the interbank market very frequently, because they can cover their needs by selling their long-term assets at a satisfactory price. This is reflected also in the amount of rationing and failures that decrease when ρ is above 0.3. The leverage immediately decreases with ρ , because the increase in liquidity of the system is more than compensated by the increase in equity, since lenders are usually repaid by borrowers and do not lose parts of their equity. Finally, in the last column of the figure, the impact of the deposits' motion is investigated. The increase of the ω parameter causes an increase in liquidity, since the shocks become gradually less and less negative. This explains also the decrease in the leverage and the rationing, because banks are less negatively impacted by the deposit shock and, consequently, need to gather less money from the market. For the same reason, the amount of credit channels decreases with a more symmetric shock distribution, while the amount of failures is substantially stable, except for a higher variability when ω describes a highly asymmetric shock.

⁷We refer the reader to Sec. Section 5.3.3 for a detailed explanation on the implementation of Figure 5.13

5. Reinforcement Learning Policy Recommendation for Interbank Network Stability

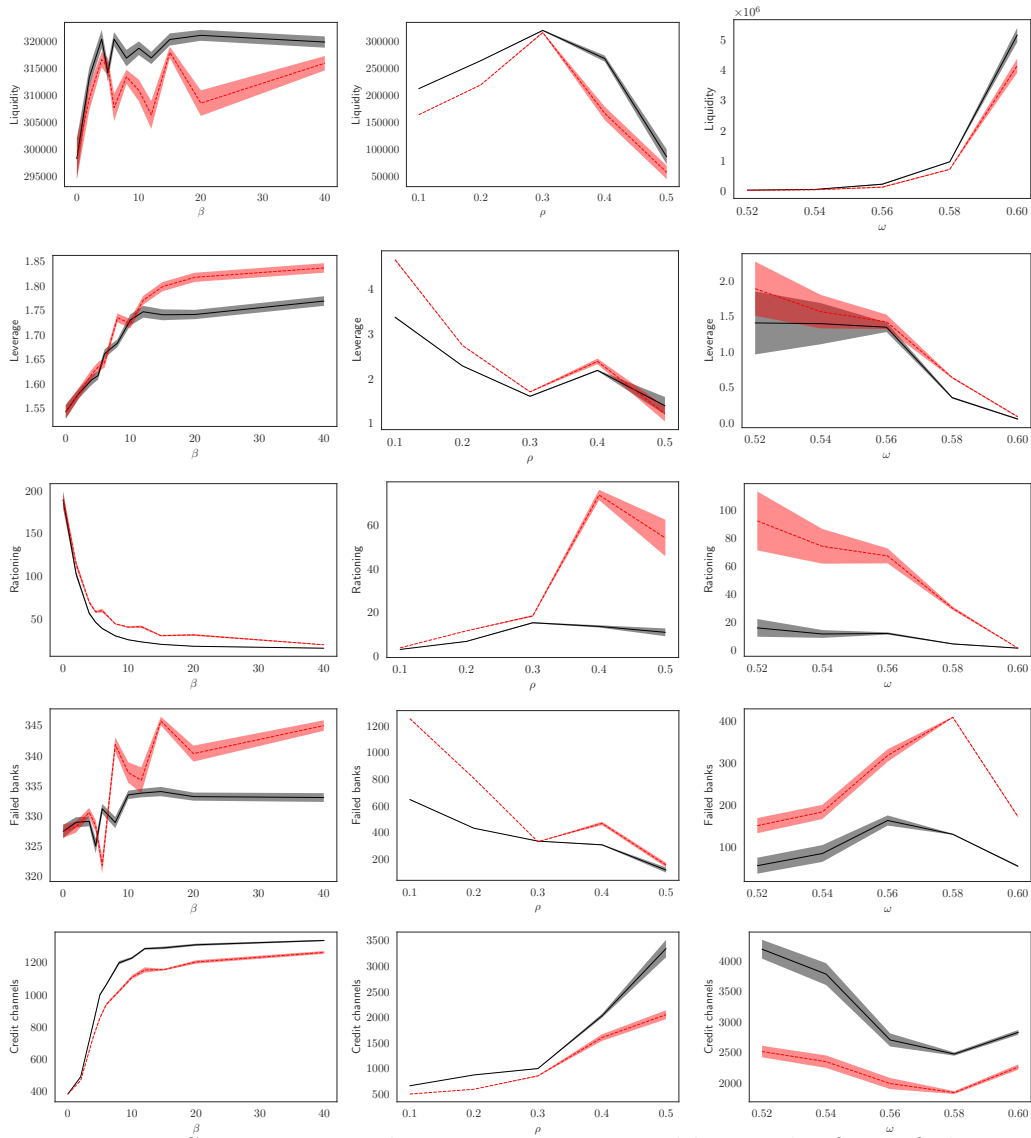


Figure 5.13: Sensitivity analysis on system variables in the face of changes in β , ρ and ω , in the first, second and third columns, respectively. The reinforcement learning algorithm is in solid black, while the random strategy is in dashed red.

5. Reinforcement Learning Policy Recommendation for Interbank Network Stability

In the last part of our analysis, we compare the results of the η obtained with the reinforcement learning algorithm with those of the two strategies, $\eta = 0$ and $\eta = 1$, kept fixed during the evolution of the system. This experiment allows us to verify the resilience of our simulated system with respect to the one obtained by implementing a fixed mechanism of parameter choice as in Berardi and Tedeschi, 2017.

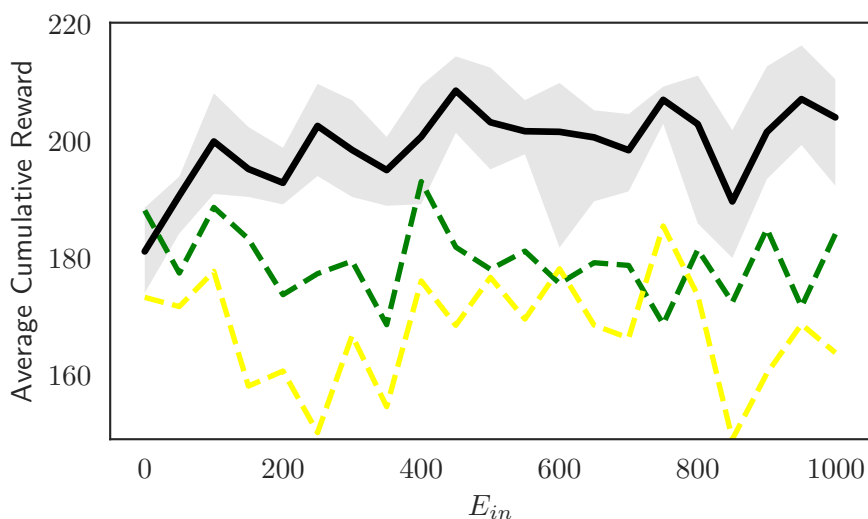


Figure 5.14: Average cumulative fitness of the system as a function of the number of training episodes for the trained PPO instances with the corresponding confidence intervals (in solid black), and the fixed strategies $\eta = 0$ (in dashed green) and $\eta = 1$ (in dashed yellow).

As already outlined in the main part of this work, the benchmark strategy adopted for comparing our results is a random one given the similarity between the two cases in the distribution of η as discussed above (see Figure 5.1, left panel). In Figure 5.14 we show that the reinforcement learning algorithm outperforms the two fixed signals in terms of aggregated fitness of the system. This result indicates that a dynamic selection of the η by looking at the available information allows for more attractive banks in the system than maintaining a fixed η . Finally, in Tab. Table 5.2 we show the results on the economic performance of the three different signals.

It is worth noting that, when the regulator adopts an η evolving through reinforcement learning, the system is more liquid and absorbs shocks better than in the other two cases (as shown by the lower leverage and lower

5. Reinforcement Learning Policy Recommendation for Interbank Network Stability

Ave. value	$\eta = 0.0$	$\eta = 1.0$	$\eta = RL$
Liquidity	302405.63 (173.84)	280389.40 (146.69)	312998.02 (151.28)
Leverage	1.75 (0.054)	1.72 (0.064)	1.59 (0.042)
Rationing	57.40 (2.90)	48.51 (2.27)	40.24 (3.75)
Failed banks	332.49 (4.86)	325.21 (4.81)	321.01 (4.10)
Credit channels	1032.21 (18.66)	1260.51 (120.67)	998.23 (43.21)

Table 5.2: Average values with standard deviations in parentheses, over times and all M Monte-Carlo simulation, of the aggregated economic variables obtained for $\eta = 0$, $\eta = 1$ and η evolving via the reinforcement learning algorithm, i.e. $\eta = RL$.

rationing and failures associated with endogenous η). Furthermore, in line with what has been described in Sec. Section 5.3.2, the $\eta = 1$ strategy always outperforms the $\eta = 0$ one, with the only known exception for the liquidity.

5.6 Appendix: Algorithms and Hyperparameters

The PPO algorithm is easier to implement than a trust-region method (Schulman et al., 2015a) and easier to tune with respect to of Deep-Q network (DQN) Mnih et al., 2015 or its continuous counterpart (Lillicrap et al., 2015). Our implementation of PPO follows Andrychowicz et al., 2020, which performs a large empirical study of the effect of implementation and parameters choices on the PPO performances. Even if we use the algorithm in a different context than their test bed, we follow the direction of their results in order to tune our hyperparameters.

As described in the main, we implement PPO in an actor-critic setting without shared architectures. When used to parametrize discrete strategies,

5. Reinforcement Learning Policy Recommendation for Interbank Network Stability

policy gradient methods like PPO output a set of logits which are then normalized to get the corresponding probabilities. Then, a greedy strategy selects the action which obtains the maximum probability. Exploration during training is guaranteed by the entropy bonus in the objective function.

The on-policy feature of PPO makes the training process episodic, so that experience is collected by interacting with the environment and then discarded immediately once the strategy has been updated. The on-policy learning appears in principle a more obvious setup for learning, even if it comes with some caveats because it makes the training less sample efficient and more computationally expensive since a new sequence of experiences need to be collected after each update step. In this process, the advantage function is computed before performing the optimization steps, when the discounted sum of returns over the episode can be computed. In order to increase the training efficiency, after one sweep through the collected samples, we compute again the advantage estimator and perform another sweep through the same experience. This trick reduces the computational expense of recollecting experiences and increases the sample efficiency of the training process. Usually we do at most 3 sweeps (epochs) over a set of collected experiences before moving on and collecting a new set.

The gradient descent optimizer is Adam (Kingma and Ba, 2014), which performs a batch update of size 100 with a learning rate equal to 0.005. Since in a reinforcement learning setting the data are not all available at the beginning of the training, we can not normalize our input variables as usual in the preprocessing step of a supervised learning context. Hence, we add a Batch Normalization layer (Ioffe and Szegedy, 2015) before the first hidden layer to normalize the inputs batch by batch and obtain the same effect.

Maximizing the objective function that returns the gradient in Eq. 5.15 is known to be unstable, since updates are not bounded and can move the strategy too far from the local optimum. Similarly to TRPO (Schulman et al., 2015a), PPO optimizes an alternative objective to mitigate the instability

$$J^{\text{CLIP}}(\theta, \psi) = \mathbb{E}_{\pi_{\theta}} \left[\min \left(r(\theta) \hat{\mathbb{A}}(s, a; \psi), \text{clip} \left(r(\theta), 1 - \epsilon, 1 + \epsilon \right) \hat{\mathbb{A}}(s, a; \psi) \right) \right] \quad (5.21)$$

where $r(\theta) = \frac{\pi(A_t|S_t;\theta)}{\pi(A_t|S_t;\theta_{\text{old}})}$ is a ratio indicating the relative probability of an action under the current strategy with respect to the old one. Instead of introducing a hard constraint as in TRPO, the ratio is bounded according to a tolerance level ϵ to limit the magnitude of the updates. The combined

5. Reinforcement Learning Policy Recommendation for Interbank Network Stability

objective function in Eq. 5.17 can be easily optimized by the PyTorch’s automatic differentiation engine, which quickly computes the gradients with respect to the two sets of parameters θ and ψ . The implemented advantage estimator depends on the parameterized value function V_ψ and is a truncated version of the one introduced by (Mnih et al., 2016) for a rollout trajectory (episode) of length T :

$$\hat{\mathbb{A}}_t = \delta_t + (\gamma\tau)\delta_{t+1} + \dots + \dots + (\gamma\tau)^{T-t+1}\delta_{T-1} \quad (5.22)$$

where $\delta_t = r_t + \gamma V_\psi(s_{t+1}) - V_\psi(s_t)$, γ is a discount rate with the same role as ρ in DQN and τ is the exponential weight discount which controls the bias variance trade-off in the advantage estimation. The generalized advantage estimator (GAE) uses a discounted sum of temporal difference residuals.

Chapter 6

Conclusions: main results

This work contributes to the financial literature by considering the application of reinforcement learning to such domain under a different lens. Instead of applying state-of-the-art algorithms to known financial problems, we perform simulations in controlled environments where we can quickly evaluate the advantages and the drawbacks of this class of methods. We present different applications of reinforcement learning from the micro perspective of a trader to the macro view of a policymaker. The frequent use of ad-hoc simulated environments helps our work stand back from other pieces of financial literature where the application is blindly carried out over preprocessed datasets, hence being more prone to biases and false discoveries in their analyses.

The first part of the work focuses on providing basic structures and concepts shared by all machine learning algorithms. Chapter 1 introduces the different machine learning paradigms and elaborates on the neural network as a leading machine learning model employed in this work. Chapter 2 outlines the fundamentals of reinforcement learning and its integration with neural networks in a relatively recent thread of research called reinforcement learning.

In the second part of the work, Chapter 3 tests multiple reinforcement learning algorithms to demonstrate their capability to retrieve known signals. Our results have shown that these methods can capture mean-reverting signals and trade accordingly, although they are sometimes unstable during their training process and require long training runtime.

For such reason, in Chapter 4, we extend the same line of research by adopting a hybrid approach that merges the model-free setting of the rein-

forcement learning algorithm applied in Brini and Tantari, 2021 with prior knowledge of the portfolio problem. Using a residual approach, we let the agent learn a deviation from the renowned Markowitz portfolio solution to capture the friction of the market caused by transaction costs. We prove that the residual approach is beneficial to stabilizing the training process and reaching a faster convergence to benchmark solution.

Our third contribution in Chapter 5 showed the flexibility of the reinforcement learning framework, which is extendable to a wide variety of sequential decision making problems, especially in finance and economics. (Brini et al., 2022) represent a still new application of reinforcement learning, where the agent learns a strategy to improve the flow of liquidity through interbank lending agreements. Modeling a public policy recommendation through a reinforcement learning approach is a new application of machine learning in finance. We found that the learned policy plays a crucial role in mitigating systemic risk with respect to alternative policy instruments.

Bibliography

- Abernethy, J., & Kale, S. (2013). Adaptive market making via online learning. *Advances in Neural Information Processing Systems*, 26.
- Abu-Mostafa, Y. S., Magdon-Ismael, M., & Lin, H.-T. (2012). *Learning from data* (Vol. 4). AMLBook New York, NY, USA:
- Acharya, V. V., & Yorulmazer, T. (2008a). Cash-in-the-market pricing and optimal resolution of bank failures. *The Review of Financial Studies*, 21(6), 2705–2742.
- Acharya, V. V., & Yorulmazer, T. (2008b). Information contagion and bank herding. *Journal of money, credit and Banking*, 40(1), 215–231.
- Agarwal, A., Bartlett, P., & Dama, M. (2010). Optimal allocation strategies for the dark pool problem. *Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics*, 9–16.
- Al-Behadili, H. N. K., Ku-Mahamud, K. R., & Sagban, R. (2018). Rule pruning techniques in the ant-miner classification algorithm and its variants: A review. *2018 IEEE Symposium on Computer Applications & Industrial Electronics (ISCAIE)*, 78–84.
- Albert, M. H., Nowakowski, R. J., & Wolfe, D. (2019). *Lessons in play: An introduction to combinatorial game theory*. CRC Press.
- Allen, F., & Gale, D. (2000). Financial contagion. *Journal of political economy*, 108(1), 1–33.
- Alsabah, H., Capponi, A., Ruiz Lacedelli, O., & Stern, M. (2021). Robo-advising: Learning investors' risk preferences via portfolio choices. *Journal of Financial Econometrics*, 19(2), 369–392.
- Altavilla, C., Lemke, W., Linzert, T., Taping, J., & von Landesberger, J. (2021). Assessing the efficacy, efficiency and potential side effects of the ecb's monetary policy instruments since 2014.

BIBLIOGRAPHY

- Amini, S., Hudson, R., Urquhart, A., & Wang, J. (2021). Nonlinearity everywhere: Implications for empirical finance, technical analysis and value at risk. *The European Journal of Finance*, 1–24.
- Ananthakumar, U., & Sarkar, R. (2017). Application of logistic regression in assessing stock performances. *2017 IEEE 15th Intl Conf on Dependable, Autonomic and Secure Computing, 15th Intl Conf on Pervasive Intelligence and Computing, 3rd Intl Conf on Big Data Intelligence and Computing and Cyber Science and Technology Congress (DASC/PiCom/DataCom/CyberSciTech)*, 1242–1247.
- Andrychowicz, M., Raichuk, A., Stańczyk, P., Orsini, M., Girgin, S., Marinier, R., Hussenot, L., Geist, M., Pietquin, O., Michalski, M., Gelly, S., & Bachem, O. (2020). What matters in on-policy reinforcement learning? a large-scale empirical study.
- Andrychowicz, M., Wolski, F., Ray, A., Schneider, J., Fong, R., Welinder, P., McGrew, B., Tobin, J., Pieter Abbeel, O., & Zaremba, W. (2017). Hindsight experience replay. *Advances in neural information processing systems*, 30.
- Angelini, P., Maresca, G., & Russo, D. (1996). Systemic risk in the netting system. *Journal of Banking & Finance*, 20(5), 853–868.
- Arulkumaran, K., Deisenroth, M. P., Brundage, M., & Bharath, A. A. (2017). A brief survey of deep reinforcement learning. *arXiv preprint arXiv:1708.05866*.
- Asness, C. S., Frazzini, A., & Pedersen, L. H. (2012). Leverage aversion and risk parity. *Financial Analysts Journal*, 68(1), 47–59.
- Asness, C. S., Moskowitz, T. J., & Pedersen, L. (2013). Value and momentum everywhere. *Journal of Finance*, 68(3), 929–985.
- Baird, L. (1995). Residual algorithms: Reinforcement learning with function approximation. *Machine learning proceedings 1995* (pp. 30–37). Elsevier.
- Barabási, A.-L., & Albert, R. (1999). Emergence of scaling in random networks. *science*, 286(5439), 509–512.
- Bartelsman, E., Scarpetta, S., & Schivardi, F. (2005). Comparative analysis of firm demographics and survival: Evidence from micro-level sources in oecd countries. *Industrial and corporate change*, 14(3), 365–391.
- Battiston, S., Gatti, D. D., Gallegati, M., Greenwald, B., & Stiglitz, J. E. (2012a). Default cascades: When does risk diversification increase stability? *Journal of Financial Stability*, 8(3), 138–149.
- Battiston, S., Gatti, D. D., Gallegati, M., Greenwald, B., & Stiglitz, J. E. (2012b). Liaisons dangereuses: Increasing connectivity, risk sharing,

BIBLIOGRAPHY

- and systemic risk. *Journal of economic dynamics and control*, 36(8), 1121–1141.
- Baxter, J., & Bartlett, P. L. (2001). Infinite-horizon policy-gradient estimation. *Journal of Artificial Intelligence Research*, 15, 319–350.
- Bellemare, M. G., Dabney, W., & Munos, R. (2017). A distributional perspective on reinforcement learning. *International Conference on Machine Learning*, 449–458.
- Bellman, R. (1956). Dynamic programming and lagrange multipliers. *Proceedings of the National Academy of Sciences of the United States of America*, 42(10), 767.
- Bellman, R. (1957). A markovian decision process. *Journal of mathematics and mechanics*, 6(5), 679–684.
- Bellman, R. (1966). Dynamic programming. *Science*, 153(3731), 34–37.
- Benhamou, E., Saltiel, D., Ungari, S., & Mukhopadhyay, A. (2020). Bridging the gap between markowitz planning and deep reinforcement learning. *arXiv preprint arXiv:2010.09108*.
- Berardi, S., & Tedeschi, G. (2017). From banks' strategies to financial (in) stability. *International Review of Economics & Finance*, 47, 255–272.
- Bernanke, B. S., Gertler, M., & Gilchrist, S. (1999). The financial accelerator in a quantitative business cycle framework. *Handbook of macroeconomics*, 1, 1341–1393.
- Bertoluzzo, F., & Corazza, M. (2012). Testing different reinforcement learning configurations for financial trading: Introduction and applications. *Procedia Economics and Finance*, 3, 68–77.
- Bertsekas, D. P. (2005). *Dynamic programming and optimal control* (3rd, Vol. 1). Athena Scientific.
- Bertsekas, D. P., & Tsitsiklis, J. N. (1996). *Neuro-dynamic programming*. Athena Scientific.
- Bindseil, U. (2018). *Financial stability implications of a prolonged period of low interest rates*. BIS.
- Black, F., & Scholes, M. (1973). The pricing of options and corporate liabilities. *The Journal of Political Economy*, 81(3), 637–654.
- Boissay, F., Collard, F., Galí, J., & Manea, C. (2021). Monetary policy and endogenous financial crises.
- Boissay, F., Collard, F., & Smets, F. (2016). Booms and banking crises. *Journal of Political Economy*, 124(2), 489–538.

BIBLIOGRAPHY

- Bollerslev, T. (1987). A conditionally heteroskedastic time series model for speculative prices and rates of return. *The review of economics and statistics*, 542–547.
- Boyd, S., Busseti, E., Diamond, S., Kahn, R. N., Koh, K., Nystrup, P., & Speth, J. (2017). *Multi-period trading via convex optimization*.
- Boyd, S., & Vandenberghe, L. (2009). *Convex optimization*. Cambridge university press.
- Brini, A., & Tantari, D. (2021). Deep reinforcement trading with predictable returns. *arXiv preprint arXiv:2104.14683*.
- Brini, A., Tedeschi, G., & Tantari, D. (2022). Reinforcement learning policy recommendation for interbank network stability. *arXiv preprint arXiv:2204.07134*.
- Brock, W. A. (2018). Nonlinearity and complex dynamics in economics and finance. *The economy as an evolving complex system* (pp. 77–97). CRC Press.
- Brunnermeier, M. K., Eisenbach, T. M., & Sannikov, Y. (2012). Macroeconomics with financial frictions: A survey.
- Buehler, H., Gonon, L., Teichmann, J., Wood, B., Mohan, B., & Kochems, J. (2019). Deep hedging: Hedging derivatives under generic market frictions using reinforcement learning. *Swiss Finance Institute Research Paper*, (19-80).
- Buşoniu, L., Babuška, R., & De Schutter, B. (2010). Multi-agent reinforcement learning: An overview. *Innovations in multi-agent systems and applications-1*, 183–221.
- Caccioli, F., Catanach, T. A., & Farmer, J. D. (2012). Heterogeneity, correlations and financial contagion. *Advances in Complex Systems*, 15(supp02), 1250058.
- Cai, Y., Judd, K. L., & Xu, R. (2013). *Numerical solution of dynamic portfolio optimization with transaction costs* (Working Paper No. 18709). National Bureau of Economic Research.
- Calice, G., Sala, C., & Tantari, D. (2020). Contingent convertible bonds in financial networks. *arXiv preprint arXiv:2009.00062*.
- Cao, J., Chen, J., Hull, J., & Poulos, Z. (2021). Deep hedging of derivatives using reinforcement learning. *The Journal of Financial Data Science*, 3(1), 10–27.
- Capponi, A., Sun, X., & Yao, D. D. (2020). A dynamic network model of interbank lending—systemic risk and liquidity provisioning. *Mathematics of Operations Research*, 45(3), 1127–1152.

BIBLIOGRAPHY

- Carlin, B. I., Lobo, M. S., & Viswanathan, S. (2007). Episodic liquidity crises: Cooperative and predatory trading. *The Journal of Finance*, 62(5), 2235–2274.
- Çelikyurt, U., & Özekici, S. (2007). Multiperiod portfolio optimization models in stochastic markets using the mean–variance approach. *European Journal of Operational Research*, 179(1), 186–202.
- Chaouki, A., Hardiman, S., Schmidt, C., Sérié, E., & De Lataillade, J. (2020). Deep deterministic portfolio optimization. *The Journal of Finance and Data Science*, 6, 16–30.
- Charpentier, A., Elie, R., & Remlinger, C. (2021). Reinforcement learning in economics and finance. *Computational Economics*, 1–38.
- Charpentier, A., Flachaire, E., & Ly, A. (2018). Econometrics and machine learning. *Economie et Statistique*, 505(1), 147–169.
- Cincotti, S., Raberto, M., & Tegli, A. (2012). Macroprudential policies in an agent-based artificial economy. *Revue de l'OFCE*, (5), 205–234.
- Clevert, D.-A., Unterthiner, T., & Hochreiter, S. (2015). Fast and accurate deep network learning by exponential linear units (elus). *arXiv preprint arXiv:1511.07289*.
- Cobbe, K. W., Hilton, J., Klimov, O., & Schulman, J. (2021). Phasic policy gradient. *International Conference on Machine Learning*, 2020–2027.
- Cohen, A., Qiao, X., Yu, L., Way, E., & Tong, X. (2019). Diverse exploration via conjugate policies for policy gradient methods. *Proceedings of the AAAI Conference on Artificial Intelligence*, 33(01), 3404–3411.
- Cong, L. W., Tang, K., Wang, J., & Zhang, Y. (2021). Alphaportfolio: Direct construction through deep reinforcement learning and interpretable ai. *Available at SSRN 3554486*.
- Cont, R. (2001). Empirical properties of asset returns: Stylized facts and statistical issues. *Quantitative Finance*, 1(2), 223–236.
- Cybenko, G. (1989). Approximation by superpositions of a sigmoidal function. *Mathematics of control, signals and systems*, 2(4), 303–314.
- Dabérius, K., Granat, E., & Karlsson, P. (2019). Deep execution-value and policy based reinforcement learning for trading and beating market benchmarks. *Available at SSRN 3374766*.
- Dabney, W., Rowland, M., Bellemare, M., & Munos, R. (2018). Distributional reinforcement learning with quantile regression. *Proceedings of the AAAI Conference on Artificial Intelligence*, 32(1).

BIBLIOGRAPHY

- Daniel, J. S. R., & Rajendran, C. (2005). A simulation-based genetic algorithm for inventory optimization in a serial supply chain. *International Transactions in Operational Research*, 12(1), 101–127.
- Dasgupta, A. (2004). Financial contagion through capital connections: A model of the origin and spread of bank panics. *Journal of the European Economic Association*, 2(6), 1049–1084.
- De Grauwe, P. (2011). The banking crisis: Causes, consequences and remedies. *Systemic implications of transatlantic regulatory cooperation and competition* (pp. 23–46). World Scientific.
- Dell’Ariccia, G., & Marquez, R. (2004). Information and bank credit allocation. *Journal of financial Economics*, 72(1), 185–214.
- de Lope, J., Maravall, D. et al. (2011). Robust high performance reinforcement learning through weighted k-nearest neighbors. *Neurocomputing*, 74(8), 1251–1259.
- Deng, J., Dong, W., Socher, R., Li, L.-J., Li, K., & Fei-Fei, L. (2009). Imagenet: A large-scale hierarchical image database. *2009 IEEE conference on computer vision and pattern recognition*, 248–255.
- Deng, Y., Bao, F., Kong, Y., Ren, Z., & Dai, Q. (2016). Deep direct reinforcement learning for financial signal representation and trading. *IEEE transactions on neural networks and learning systems*, 28(3), 653–664.
- Devlin, J., Chang, M.-W., Lee, K., & Toutanova, K. (2018). Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*.
- Dixon, M. F., Halperin, I., & Bilokon, P. (2020). *Machine learning in finance*. Springer.
- Du, J., Jin, M., Kolm, P. N., Ritter, G., Wang, Y., & Zhang, B. (2020). Deep reinforcement learning for option replication and hedging. *The Journal of Financial Data Science*, 2(4), 44–57.
- Du, X., Zhai, J., & Lv, K. (2016). Algorithm trading using q-learning and recurrent reinforcement learning. *positions*, 1(1).
- Engle, R. F., & Ferstenberg, R. (2007). Execution risk. *The Journal of Trading*, 2(2), 10–20.
- Engstrom, L., Ilyas, A., Santurkar, S., Tsipras, D., Janoos, F., Rudolph, L., & Madry, A. (2020). Implementation matters in deep policy gradients: A case study on ppo and trpo. *arXiv preprint arXiv:2005.12729*.
- Fadja, A. N., Lamma, E., Riguzzi, F. et al. (2018). Vision inspection with neural networks. *RiCeRcA@ AI* IA*.

BIBLIOGRAPHY

- Farebrother, J., Machado, M. C., & Bowling, M. (2018). Generalization and regularization in dqn. *arXiv preprint arXiv:1810.00123*.
- Freixas, X., Parigi, B. M., & Rochet, J.-C. (2000). Systemic risk, interbank relations, and liquidity provision by the central bank. *Journal of money, credit and banking*, 611–638.
- Friedman, J., Hastie, T., Tibshirani, R. et al. (2001). *The elements of statistical learning* (Vol. 1). Springer series in statistics New York.
- Fujimoto, S., Hoof, H., & Meger, D. (2018). Addressing function approximation error in actor-critic methods. *International conference on machine learning*, 1587–1596.
- Galí, J. (2015). *Monetary policy, inflation, and the business cycle: An introduction to the new keynesian framework and its applications*. Princeton University Press.
- Ganchev, K., Nevmyvaka, Y., Kearns, M., & Vaughan, J. W. (2010). Censored exploration and the dark pool problem. *Communications of the ACM*, 53(5), 99–107.
- Gao, Z., Gao, Y., Hu, Y., Jiang, Z., & Su, J. (2020). Application of deep q-network in portfolio management. *2020 5th IEEE International Conference on Big Data Analytics (ICBDA)*, 268–275.
- Gârleanu, N., & Pedersen, L. H. (2013). Dynamic trading with predictable returns and transaction costs. *The Journal of Finance*, 68(6), 2309–2340.
- Gârleanu, N., Pedersen, L. H., & Poteshman, A. M. (2008). Demand-based option pricing. *The Review of Financial Studies*, 22(10), 4259–4299.
- Georg, C.-P. (2013). The effect of the interbank network structure on contagion and common shocks. *Journal of Banking & Finance*, 37(7), 2216–2228.
- Géron, A. (2019). *Hands-on machine learning with scikit-learn, keras, and tensorflow: Concepts, tools, and techniques to build intelligent systems*. O’Reilly Media.
- Gertler, M., Kiyotaki, N., & Prestipino, A. (2020). A macroeconomic model with financial panics. *The Review of Economic Studies*, 87(1), 240–288.
- Ghadirzadeh, A., Bütepage, J., Maki, A., Kragic, D., & Björkman, M. (2016). A sensorimotor reinforcement learning framework for physical human-robot interaction. *2016 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2682–2688.

BIBLIOGRAPHY

- Giri, F., Riccetti, L., Russo, A., & Gallegati, M. (2019). Monetary policy and large crises in a financial accelerator agent-based model. *Journal of Economic Behavior & Organization*, 157, 42–58.
- Glorot, X., & Bengio, Y. (2010). Understanding the difficulty of training deep feedforward neural networks. *Proceedings of the thirteenth international conference on artificial intelligence and statistics*, 249–256.
- Goldberg, J. E., Klee, E., Prescott, E. S., & Wood, P. R. (2020). Monetary policy strategies and tools: Financial stability considerations.
- Gondzio, J., & Grothey, A. (2007). Solving non-linear portfolio optimization problems with the primal-dual interior point method. *European Journal of Operational Research*, 181(3), 1019–1029.
- Goodell, J. W., Kumar, S., Lim, W. M., & Pattnaik, D. (2021). Artificial intelligence and machine learning in finance: Identifying foundations, themes, and research clusters from bibliometric analysis. *Journal of Behavioral and Experimental Finance*, 100577.
- Goodfellow, I., Bengio, Y., & Courville, A. (2016). *Deep learning*. MIT press.
- Graves, A. (2012). Supervised sequence labelling. *Supervised sequence labelling with recurrent neural networks* (pp. 5–13). Springer.
- Grilli, R., Iori, G., Stamboglis, N., & Tedeschi, G. (2017). A networked economy: A survey on the effect of interaction in credit markets. *Introduction to agent-based economics* (pp. 229–252). Elsevier.
- Grilli, R., Tedeschi, G., & Gallegati, M. (2014). Network approach for detecting macroeconomic instability. *2014 Tenth International Conference on Signal-Image Technology and Internet-Based Systems*, 440–446.
- Grinold, R. (2006). A dynamic model of portfolio management. *Journal of Investment Management*, 4, 5–22.
- Gruenstein, J., Chen, T., Doshi, N., & Agrawal, P. (2021). Residual model learning for microrobot control. *2021 IEEE International Conference on Robotics and Automation (ICRA)*, 7219–7226.
- Guéant, O. (2013). Permanent market impact can be nonlinear. *arXiv preprint arXiv:1305.0413*.
- Guida, T. (2020). Factor investing and asset pricing anomalies. *Machine learning for factor investing* (pp. 13–34). Chapman; Hall/CRC.
- Haarnoja, T., Zhou, A., Hartikainen, K., Tucker, G., Ha, S., Tan, J., Kumar, V., Zhu, H., Gupta, A., Abbeel, P. et al. (2018). Soft actor-critic algorithms and applications. *arXiv preprint arXiv:1812.05905*.
- Haldane, A. G., & May, R. M. (2011). Systemic risk in banking ecosystems. *Nature*, 469(7330), 351–355.

BIBLIOGRAPHY

- Halperin, I. (2020). Qlbs: Q-learner in the black-scholes (-merton) worlds. *The Journal of Derivatives*, 28(1), 99–122.
- Halperin, I., Liu, J., & Zhang, X. (2022). Combining reinforcement learning and inverse reinforcement learning for asset allocation recommendations. *Available at SSRN 4002715*.
- Hambly, B., Xu, R., & Yang, H. (2021). Recent advances in reinforcement learning in finance. *arXiv preprint arXiv:2112.04553*.
- Han, D., Zhang, J., Zhou, Y., Liu, Q., & Yang, N. (2019). Intelligent trader model based on deep reinforcement learning. *International Conference on Web Information Systems and Applications*, 15–21.
- Hasselt, H. (2010). Double q-learning. *Advances in neural information processing systems*, 23.
- He, K., Zhang, X., Ren, S., & Sun, J. (2015). Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. *Proceedings of the IEEE international conference on computer vision*, 1026–1034.
- Heaton, J. B., Polson, N. G., & Witte, J. H. (2017). Deep learning for finance: Deep portfolios. *Applied Stochastic Models in Business and Industry*, 33(1), 3–12.
- Hendricks, D., & Wilcox, D. (2014). A reinforcement learning extension to the almgren-chriss framework for optimal trade execution. *2014 IEEE Conference on Computational Intelligence for Financial Engineering & Economics (CIFEr)*, 457–464.
- Hochreiter, S., & Schmidhuber, J. (1997). Long short-term memory. *Neural computation*, 9(8), 1735–1780.
- Hoff, K., & Stiglitz, J. E. (1990). Introduction: Imperfect information and rural credit markets: Puzzles and policy perspectives. *The world bank economic review*, 4(3), 235–250.
- Hornik, K., Stinchcombe, M., & White, H. (1989). Multilayer feedforward networks are universal approximators. *Neural networks*, 2(5), 359–366.
- Hu, G., Yang, Y., Yi, D., Kittler, J., Christmas, W., Li, S. Z., & Hospedales, T. (2015). When face recognition meets with deep learning: An evaluation of convolutional neural networks for face recognition. *Proceedings of the IEEE international conference on computer vision workshops*, 142–150.

BIBLIOGRAPHY

- Hu, Y.-J., & Lin, S.-J. (2019). Deep reinforcement learning for optimizing finance portfolio management. *2019 Amity International Conference on Artificial Intelligence (AICAI)*, 14–20.
- Hua, Z., Wang, Y., Xu, X., Zhang, B., & Liang, L. (2007). Predicting corporate financial distress based on integration of support vector machine and logistic regression. *Expert Systems with Applications*, *33*(2), 434–440.
- Huge, B. N., & Savine, A. (2020). Differential machine learning. *Available at SSRN 3591734*.
- Ioffe, S., & Szegedy, C. (2015). Batch normalization: Accelerating deep network training by reducing internal covariate shift. *arXiv preprint arXiv:1502.03167*.
- Iori, G., Jafarey, S., & Padilla, F. G. (2006). Systemic risk on the interbank market. *Journal of Economic Behavior & Organization*, *61*(4), 525–542.
- Iori, G., & Mantegna, R. N. (2018). Empirical analyses of networks in finance. *Handbook of computational economics* (pp. 637–685). Elsevier.
- Jiang, Z., Xu, D., & Liang, J. (2017). A deep reinforcement learning framework for the financial portfolio management problem. *arXiv preprint arXiv:1706.10059*.
- Jiménez, G., Ongena, S., Peydró, J.-L., & Saurina, J. (2014). Hazardous times for monetary policy: What do twenty-three million bank loans say about the effects of monetary policy on credit risk-taking? *Econometrica*, *82*(2), 463–505.
- Jin, O., & El-Saawy, H. (2016). Portfolio management using reinforcement learning. *Stanford University*.
- Johannink, T., Bahl, S., Nair, A., Luo, J., Kumar, A., Loskyll, M., Ojea, J. A., Solowjow, E., & Levine, S. (2019). Residual reinforcement learning for robot control. *2019 International Conference on Robotics and Automation (ICRA)*, 6023–6029.
- Kaniel, R., Lin, Z., Pelger, M., & Van Nieuwerburgh, S. (2021). Machine-learning the skill of mutual fund managers. *Available at SSRN 3977883*.
- Kempf, A., & Memmel, C. (2006). Estimating the global minimum variance portfolio. *Schmalenbach Business Review*, *58*(4), 332–348.
- Kim, K.-j. (2003). Financial time series forecasting using support vector machines. *Neurocomputing*, *55*(1-2), 307–319.
- Kingma, D. P., & Ba, J. (2014). Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*.

BIBLIOGRAPHY

- Kohavi, R., Wolpert, D. H. et al. (1996). Bias plus variance decomposition for zero-one loss functions. *ICML*, 96, 275–83.
- Kolm, P., & Ritter, G. (2019a). Modern perspectives on reinforcement learning in finance. *SSRN Electronic Journal*.
- Kolm, P. N., & Maclin, L. (2012). Algorithmic trading, optimal execution, and dynamic portfolios. *The oxford handbook of quantitative asset management*.
- Kolm, P. N., & Ritter, G. (2014). Multiperiod portfolio selection and bayesian dynamic models. *Risk*, 28(3), 50–54.
- Kolm, P. N., & Ritter, G. (2019b). Dynamic replication and hedging: A reinforcement learning approach. *The Journal of Financial Data Science*, 1(1), 159–171.
- Kolm, P. N., Tütüncü, R., & Fabozzi, F. J. (2014). 60 years of portfolio optimization: Practical challenges and current trends. *European Journal of Operational Research*, 234(2), 356–371.
- Konidaris, G., Osentoski, S., & Thomas, P. (2011). Value function approximation in reinforcement learning using the fourier basis. *Twenty-fifth AAAI conference on artificial intelligence*.
- Krizhevsky, A., Sutskever, I., & Hinton, G. E. (2012). Imagenet classification with deep convolutional neural networks. *Advances in neural information processing systems*, 25, 1097–1105.
- Kumar, M., & Thenmozhi, M. (2006). Forecasting stock index movement: A comparison of support vector machines and random forest. *Indian institute of capital markets 9th capital markets conference paper*.
- Kyle, A. S., & Obizhaeva, A. A. (2016). Market microstructure invariance: Empirical hypotheses. *Econometrica*, 84(4), 1345–1404.
- Le, Q. V. (2013). Building high-level features using large scale unsupervised learning. *2013 IEEE international conference on acoustics, speech and signal processing*, 8595–8598.
- LeCun, Y., Bengio, Y. et al. (1995). Convolutional networks for images, speech, and time series. *The handbook of brain theory and neural networks*, 3361(10), 1995.
- LeCun, Y., Boser, B., Denker, J. S., Henderson, D., Howard, R. E., Hubbard, W., & Jackel, L. D. (1989). Backpropagation applied to handwritten zip code recognition. *Neural computation*, 1(4), 541–551.
- Lee, I., Epelman, M. A., Romeijn, H. E., & Smith, R. L. (2017). Simplex algorithm for countable-state discounted markov decision processes. *Operations Research*, 65(4), 1029–1042.

BIBLIOGRAPHY

- Lenzu, S., & Tedeschi, G. (2012). Systemic risk on different interbank network topologies. *Physica A: Statistical Mechanics and its Applications*, 391(18), 4331–4341.
- Leshno, M., Lin, V. Y., Pinkus, A., & Schocken, S. (1993). Multilayer feed-forward networks with a nonpolynomial activation function can approximate any function. *Neural networks*, 6(6), 861–867.
- Levine, S., Finn, C., Darrell, T., & Abbeel, P. (2016). End-to-end training of deep visuomotor policies. *The Journal of Machine Learning Research*, 17(1), 1334–1373.
- Liang, Z., Chen, H., Zhu, J., Jiang, K., & Li, Y. (2018). Adversarial deep reinforcement learning in portfolio management. *arXiv preprint arXiv:1808.09940*.
- Lillicrap, T. P., Hunt, J. J., Pritzel, A., Heess, N., Erez, T., Tassa, Y., Silver, D., & Wierstra, D. (2015). Continuous control with deep reinforcement learning. *arXiv preprint arXiv:1509.02971*.
- Lillo, F., Farmer, J. D., & Mantegna, R. N. (2003). Master curve for price-impact function. *Nature*, 421(6919), 129–130.
- Lim, B., Zohren, S., & Roberts, S. (2019). Enhancing time-series momentum strategies using deep neural networks. *The Journal of Financial Data Science*, 1(4), 19–38.
- Lin, L.-J. (1992). *Reinforcement learning for robots using neural networks*. Carnegie Mellon University.
- Lin, S., & Beling, P. A. (2020). An end-to-end optimal trade execution framework based on proximal policy optimization. *IJCAI*, 4548–4554.
- Linnainmaa, S. (1970). The representation of the cumulative rounding error of an algorithm as a Taylor expansion of the local rounding errors. *Master's Thesis (in Finnish), Univ. Helsinki*, 6–7.
- Liu, A., Mo, C. Y. J., Paddrik, M. E., & Yang, S. Y. (2018). An agent-based approach to interbank market lending decisions and risk implications. *Information*, 9(6), 132.
- Liu, J., Gu, X., & Liu, S. (2019a). Policy optimization reinforcement learning with entropy regularization. *arXiv preprint arXiv:1912.01557*.
- Liu, X.-Y., Yang, H., Chen, Q., Zhang, R., Yang, L., Xiao, B., & Wang, C. D. (2020). Finrl: A deep reinforcement learning library for automated stock trading in quantitative finance. *arXiv preprint arXiv:2011.09607*.
- Liu, Z., Li, X., Kang, B., & Darrell, T. (2019b). Regularization matters in policy optimization—an empirical study on continuous control. *arXiv preprint arXiv:1910.09191*.

BIBLIOGRAPHY

- López de Prado, M. (2019). Beyond econometrics: A roadmap towards financial machine learning. *Available at SSRN 3365282*.
- Lozano, F., Lozano, J., & García Molina, M. (2007). An artificial economy based on reinforcement learning and agent based modeling. *Documentos de Trabajo, Facultad de Economía, Universidad del Rosario*, (18).
- Lundberg, S. M., & Lee, S.-I. (2017). A unified approach to interpreting model predictions. *Proceedings of the 31st international conference on neural information processing systems*, 4768–4777.
- Malliaris, A. G., & Malliaris, M. (2015). What drives gold returns? a decision tree analysis. *Finance Research Letters*, 13, 45–53.
- Mao, H., Alizadeh, M., Menache, I., & Kandula, S. (2016). Resource management with deep reinforcement learning. *Proceedings of the 15th ACM workshop on hot topics in networks*, 50–56.
- Marbach, P., & Tsitsiklis, J. (2001). Simulation-based optimization of markov reward processes. *IEEE Transactions on Automatic Control*, 46(2), 191–209.
- Markowitz, H. (1952). Portfolio selection. *The Journal of Finance*, 7(1), 77–91.
- Maudos, J., & De Guevara, J. F. (2004). Factors explaining the interest margin in the banking sectors of the european union. *Journal of Banking & Finance*, 28(9), 2259–2281.
- McCulloch, W. S., & Pitts, W. (1943). A logical calculus of the ideas immanent in nervous activity. *The bulletin of mathematical biophysics*, 5(4), 115–133.
- Mehdipour Ghazi, M., & Kemal Ekenel, H. (2016). A comprehensive analysis of deep learning based representation for face recognition. *Proceedings of the IEEE conference on computer vision and pattern recognition workshops*, 34–41.
- Mei, X., DeMiguel, V., & Nogales, F. J. (2016). Multiperiod portfolio optimization with multiple risky assets and general transaction costs. *Journal of Banking & Finance*, 69, 108–120.
- Melo, F. S., Meyn, S. P., & Ribeiro, M. I. (2008). An analysis of reinforcement learning with function approximation. *Proceedings of the 25th international conference on Machine learning*, 664–671.
- Merton, R. C. (1969). Lifetime portfolio selection under uncertainty: The continuous-time case. *Review of Economics and statistics*, 51(3), 247–257.

BIBLIOGRAPHY

- Minsky, H. P. (1964). Longer waves in financial relations: Financial factors in the more severe depressions. *The American Economic Review*, 54(3), 324–335.
- Minsky, M., & Papert, S. (1969). Perceptrons.
- Mnih, V., Badia, A. P., Mirza, M., Graves, A., Lillicrap, T. P., Harley, T., Silver, D., & Kavukcuoglu, K. (2016). Asynchronous methods for deep reinforcement learning.
- Mnih, V., Kavukcuoglu, K., Silver, D., Graves, A., Antonoglou, I., Wierstra, D., & Riedmiller, M. (2013). Playing atari with deep reinforcement learning. *arXiv preprint arXiv:1312.5602*.
- Mnih, V., Kavukcuoglu, K., Silver, D., Rusu, A. A., Veness, J., Bellemare, M. G., Graves, A., Riedmiller, M., Fidjeland, A. K., Ostrovski, G. et al. (2015). Human-level control through deep reinforcement learning. *nature*, 518(7540), 529–533.
- Mosavi, A., Faghan, Y., Ghamisi, P., Duan, P., Ardabili, S. F., Salwana, E., & Band, S. S. (2020). Comprehensive review of deep reinforcement learning methods and applications in economics. *Mathematics*, 8(10), 1640.
- Moskowitz, T. J., Ooi, Y. H., & Pedersen, L. H. (2012). Time series momentum. *Journal of Financial Economics*, 104(2), 228–250.
- Noguer i Alonso, M., & Srivastava, S. (2020). Deep reinforcement learning for asset allocation in us equities. *Sonam, Deep Reinforcement Learning for Asset Allocation in US Equities (October 14, 2020)*.
- OECD. (2021). Artificial intelligence, machine learning and big data in finance: Opportunities, challenges and implications for policy makers.
- Osoba, O. A., Vardavas, R., Grana, J., Zutshi, R., & Jaycocks, A. (2020). Policy-focused agent-based modeling using rl behavioral models. *arXiv preprint arXiv:2006.05048*.
- Park, H., Sim, M. K., & Choi, D. G. (2020). An intelligent financial portfolio trading strategy using deep q-learning. *Expert Systems with Applications*, 158, 113573.
- Patzelt, F., & Bouchaud, J.-P. (2018). Universal scaling and nonlinearity of aggregate price impact in financial markets. *Physical Review E*, 97(1), 012304.
- Peters, J., & Bagnell, J. A. (2010). Policy gradient methods. *Scholarpedia*, 5(11), 3698.

BIBLIOGRAPHY

- Radford, A., Metz, L., & Chintala, S. (2015). Unsupervised representation learning with deep convolutional generative adversarial networks. *arXiv preprint arXiv:1511.06434*.
- Raina, R., Madhavan, A., & Ng, A. Y. (2009). Large-scale deep unsupervised learning using graphics processors. *Proceedings of the 26th annual international conference on machine learning*, 873–880.
- Ravi, V., Pradeepkumar, D., & Deb, K. (2017). Financial time series prediction using hybrids of chaos theory, multi-layer perceptron and multi-objective evolutionary algorithms. *Swarm and Evolutionary Computation*, 36, 136–149.
- Riccetti, L., Russo, A., & Gallegati, M. (2018). Financial regulation and endogenous macroeconomic crises. *Macroeconomic Dynamics*, 22(4), 896–930.
- Rochet, J.-C., & Tirole, J. (1996). Interbank lending and systemic risk. *Why Are there So Many Banking Crises?*, 140.
- Rosenblatt, F. (1958). The perceptron: A probabilistic model for information storage and organization in the brain. *Psychological review*, 65(6), 386.
- Ruder, S. (2016). An overview of gradient descent optimization algorithms. *arXiv preprint arXiv:1609.04747*.
- Rumelhart, D. E., Hinton, G. E., & Williams, R. J. (1986). Learning representations by back-propagating errors. *nature*, 323(6088), 533–536.
- Ryll, L., Barton, M. E., Zhang, B. Z., McWaters, R. J., Schizas, E., Hao, R., Bear, K., Preziuso, M., Seger, E., Wardrop, R. et al. (2020). Transforming paradigms: A global ai in financial services survey.
- Samuelson, P. A. (1969). Lifetime portfolio selection by dynamic stochastic programming. *The Review of Economics and Statistics*, 51(3), 239–246.
- Sarakhsi, M. K., Ghomi, S. F., & Karimi, B. (2016). A new hybrid algorithm of scatter search and nelder–mead algorithms to optimize joint economic lot sizing problem. *Journal of Computational and Applied Mathematics*, 292, 387–401.
- Sato, Y. (2019). Model-free reinforcement learning for financial portfolios: A brief survey. *arXiv preprint arXiv:1904.04973*.
- Schneeweis, T., & Spurgin, R. B. (1999). Alpha, alpha... who's got the alpha? *The Journal of Alternative Investments*, 2(3), 83–87.
- Schrittwieser, J., Antonoglou, I., Hubert, T., Simonyan, K., Sifre, L., Schmitt, S., Guez, A., Lockhart, E., Hassabis, D., Graepel, T. et al. (2020).

BIBLIOGRAPHY

- Mastering atari, go, chess and shogi by planning with a learned model. *Nature*, 588(7839), 604–609.
- Schulman, J., Levine, S., Abbeel, P., Jordan, M., & Moritz, P. (2015a). Trust region policy optimization. *International conference on machine learning*, 1889–1897.
- Schulman, J., Moritz, P., Levine, S., Jordan, M., & Abbeel, P. (2015b). High-dimensional continuous control using generalized advantage estimation. *arXiv preprint arXiv:1506.02438*.
- Schulman, J., Wolski, F., Dhariwal, P., Radford, A., & Klimov, O. (2017). Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*.
- Sehgal, N., & Pandey, K. K. (2015). Artificial intelligence methods for oil price forecasting: A review and evaluation. *Energy Systems*, 6(4), 479–506.
- Shah, D., & Xie, Q. (2018). Q-learning with nearest neighbors. *Advances in Neural Information Processing Systems*, 31.
- Shapley, L. S. (2016). *17. a value for n-person games*. Princeton University Press.
- Sharang, A., & Rao, C. (2015). Using machine learning for medium frequency derivative portfolio trading. *arXiv preprint arXiv:1512.06228*.
- Sharpe, W. F. (1994). The sharpe ratio. *Journal of portfolio management*, 21(1), 49–58.
- Silver, D., Hubert, T., Schrittwieser, J., Antonoglou, I., Lai, M., Guez, A., Lanctot, M., Sifre, L., Kumaran, D., Graepel, T. et al. (2017a). Mastering chess and shogi by self-play with a general reinforcement learning algorithm. *arXiv preprint arXiv:1712.01815*.
- Silver, D., Lever, G., Heess, N., Degris, T., Wierstra, D., & Riedmiller, M. (2014). Deterministic policy gradient algorithms. *International conference on machine learning*, 387–395.
- Silver, D., Schrittwieser, J., Simonyan, K., Antonoglou, I., Huang, A., Guez, A., Hubert, T., Baker, L., Lai, M., Bolton, A. et al. (2017b). Mastering the game of go without human knowledge. *nature*, 550(7676), 354–359.
- Silver, D., Singh, S., Precup, D., & Sutton, R. S. (2021). Reward is enough. *Artificial Intelligence*, 299, 103535.
- Silver, T., Allen, K., Tenenbaum, J., & Kaelbling, L. (2018). Residual policy learning. *arXiv preprint arXiv:1812.06298*.
- Spooner, T., Fearnley, J., Savani, R., & Koukorinis, A. (2018). Market making via reinforcement learning. *arXiv preprint arXiv:1804.04216*.

BIBLIOGRAPHY

- Srivastava, N., Hinton, G., Krizhevsky, A., Sutskever, I., & Salakhutdinov, R. (2014). Dropout: A simple way to prevent neural networks from overfitting. *The journal of machine learning research*, 15(1), 1929–1958.
- Sutton, R. S., & Barto, A. G. (2018). *Reinforcement learning: An introduction*. MIT press.
- Sutton, R. S., McAllester, D., Singh, S., & Mansour, Y. (1999). Policy gradient methods for reinforcement learning with function approximation. *Advances in neural information processing systems*, 12.
- Svensson, L. E. (2017). Cost-benefit analysis of leaning against the wind. *Journal of Monetary Economics*, 90, 193–213.
- Szepesvári, C. (2010). Algorithms for reinforcement learning. *Synthesis lectures on artificial intelligence and machine learning*, 4(1), 1–103.
- Taherkhani, A., Cosma, G., & McGinnity, T. M. (2018). Deep-fs: A feature selection algorithm for deep boltzmann machines. *Neurocomputing*, 322, 22–37.
- Taigman, Y., Yang, M., Ranzato, M., & Wolf, L. (2014). Deepface: Closing the gap to human-level performance in face verification. *Proceedings of the IEEE conference on computer vision and pattern recognition*, 1701–1708.
- Taylor, J. B. (2011). Macroeconomic lessons from the great deviation. *NBER macroeconomics annual*, 25(1), 387–395.
- Tedeschi, G., Mazlounian, A., Gallegati, M., & Helbing, D. (2012). Bankruptcy cascades in interbank markets. *PloS one*, 7(12), e52749.
- Tedeschi, G., Vidal-Tomás, D., Delli-Gatti, D., & Gallegati, M. (2021). The macroeconomic effects of default and debt restructuring: An agent based exploration. *International Review of Economics & Finance*, 76, 1146–1163.
- Tesauro, G. et al. (1995). Temporal difference learning and td-gammon. *Communications of the ACM*, 38(3), 58–68.
- Thakor, A. V. (2020). Fintech and banking: What do we know? *Journal of Financial Intermediation*, 41, 100833.
- Trichet, J.-C. (2010). Reflections on the nature of monetary policy non-standard measures and finance theory, opening address at the ecb central banking conference. frankfurt, germany, 18.
- Truong, Q., Nguyen, M., Dang, H., & Mei, B. (2020). Housing price prediction via improved machine learning techniques. *Procedia Computer Science*, 174, 433–442.

BIBLIOGRAPHY

- Tutuncu, R. (2011). Recent advances in portfolio optimization. *The Oxford Handbook of Quantitative Asset Management*.
- Tuyls, P., Hollmann, H. D., Van Lint, J. H., & Tolhuizen, L. (2005). XOR-based visual cryptography schemes. *Designs, Codes and Cryptography*, 37(1), 169–186.
- Uhlenbeck, G. E., & Ornstein, L. S. (1930). On the theory of the brownian motion. *Phys. Rev.*, 36, 823–841.
- Upper, C. (2011). Simulation methods to assess the danger of contagion in interbank markets. *Journal of Financial Stability*, 7(3), 111–125.
- Uther, W. T., & Veloso, M. M. (1998). Tree based discretization for continuous state space reinforcement learning. *Aaai/iaai*, 98, 769–774.
- Van Hasselt, H., Guez, A., & Silver, D. (2016). Deep reinforcement learning with double q-learning. *Proceedings of the AAAI conference on artificial intelligence*, 30(1).
- Vapnik, V., & Chervonenkis, A. Y. (1971). On the uniform convergence of relative frequencies of events to their probabilities. *Measures of Complexity*, 16(2), 11.
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, L., & Polosukhin, I. (2017). Attention is all you need. *Advances in neural information processing systems*, 5998–6008.
- Wang, C., Han, D., Liu, Q., & Luo, S. (2018). A deep learning approach for credit scoring of peer-to-peer lending using attention mechanism lstm. *IEEE Access*, 7, 2161–2168.
- Wang, H., & Yu, S. (2021). Robo-advising: Enhancing investment with inverse optimization and deep reinforcement learning. *2021 20th IEEE International Conference on Machine Learning and Applications (ICMLA)*, 365–372.
- Wang, H., Jiang, Y., & Wang, H. (2009). Stock return prediction based on bagging-decision tree. *2009 IEEE International Conference on Grey Systems and Intelligent Services (GSIS 2009)*, 1575–1580.
- Wang, J., Li, Y., & Cao, Y. (2019). Dynamic portfolio management with reinforcement learning. *arXiv preprint arXiv:1911.11880*.
- Wang, X. (2016). Deep learning in object recognition, detection, and segmentation. *Foundations and Trends in Signal Processing*, 8(4), 217–382.
- Wang, Y., O’Donoghue, B., & Boyd, S. (2015). Approximate dynamic programming via iterated bellman inequalities. *International Journal of Robust and Nonlinear Control*, 25(10), 1472–1496.

BIBLIOGRAPHY

- Wang, Z., Bapst, V., Heess, N., Mnih, V., Munos, R., Kavukcuoglu, K., & de Freitas, N. (2016). Sample efficient actor-critic with experience replay. *arXiv preprint arXiv:1611.01224*.
- Wardell, D. C., & Peterson, G. L. (2006). Fuzzy state aggregation and policy hill climbing for stochastic environments. *International Journal of Computational Intelligence and Applications*, 6(03), 413–428.
- Watkins, C. J., & Dayan, P. (1992). Q-learning. *Machine learning*, 8(3-4), 279–292.
- Wei, H., Wang, Y., Mangu, L., & Decker, K. (2019). Model-based reinforcement learning for predictions and control for limit order books. *arXiv preprint arXiv:1910.03743*.
- Werbos, P. J. (1982). Applications of advances in nonlinear sensitivity analysis. *System modeling and optimization* (pp. 762–770). Springer.
- Werbos, P. J. (1987). Building and understanding adaptive systems: A statistical/numerical approach to factory automation and brain research. *IEEE Transactions on Systems, Man, and Cybernetics*, 17(1), 7–20.
- Wolpert, D. H., & Macready, W. G. (1997). No free lunch theorems for optimization. *IEEE transactions on evolutionary computation*, 1(1), 67–82.
- Xiong, Z., Liu, X.-Y., Zhong, S., Yang, H., & Walid, A. (2018). Practical deep reinforcement learning approach for stock trading. *arXiv preprint arXiv:1811.07522*.
- Xu, B., Wang, N., Chen, T., & Li, M. (2015). Empirical evaluation of rectified activations in convolutional network. *arXiv preprint arXiv:1505.00853*.
- Yu, P., Lee, J. S., Kulyatin, I., Shi, Z., & Dasgupta, S. (2019). Model-based deep reinforcement learning for dynamic portfolio optimization. *arXiv preprint arXiv:1901.08740*.
- Zhang, G., Wang, C., Xu, B., & Grosse, R. (2018). Three mechanisms of weight decay regularization. *arXiv preprint arXiv:1810.12281*.
- Zhang, X., & Saniie, J. (2021). Unsupervised learning for 3d ultrasonic data compression. *2021 IEEE International Ultrasonics Symposium (IUS)*, 1–3.
- Zhang, Z., Zohren, S., & Roberts, S. (2020). Deep reinforcement learning for trading. *The Journal of Financial Data Science*, 2(2), 25–40.
- Zhou, H., Chen, J., Varshney, L. R., & Jagmohan, A. (2020). Nonstationary reinforcement learning with linear function approximation. *arXiv preprint arXiv:2010.04244*.

BIBLIOGRAPHY

- Zhou, L. (2020). Product advertising recommendation in e-commerce based on deep learning and distributed expression. *Electronic Commerce Research*, 20(2), 321–342.
- Zhu, B., Yang, W., Wang, H., & Yuan, Y. (2018). A hybrid deep learning model for consumer credit scoring. *2018 International Conference on Artificial Intelligence and Big Data (ICAIBD)*, 205–208.