

## Todo list



# Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
1.1	The March of the Machines . . . . .	1
1.2	Artificial Intelligence and Finance . . . . .	1
1.3	Structure . . . . .	1
<b>2</b>	<b>Discrete-Time Stochastic Optimal Control</b>	<b>3</b>
2.1	Sequential Decision Problems . . . . .	3
2.2	Markov Decision Processes . . . . .	4
2.3	Policies . . . . .	5
2.4	Risk-Neutral Framework . . . . .	6
2.4.1	Discounted Reward Formulation . . . . .	6
2.4.2	Average Reward Formulation . . . . .	9
2.5	Risk-Sensitive Framework . . . . .	10
2.5.1	Discounted Reward Formulation . . . . .	11
2.5.2	Average Reward Formulation . . . . .	13
2.6	Dynamic Programming Algorithms . . . . .	15
2.6.1	Value Iteration . . . . .	15
2.6.2	Policy Iteration . . . . .	15
<b>3</b>	<b>Reinforcement Learning</b>	<b>17</b>
3.1	The Reinforcement Learning Problem . . . . .	18
3.2	Model-Free RL Methods . . . . .	18
3.2.1	Model Approximation . . . . .	20
3.2.2	Value Approximation . . . . .	20
3.2.3	Policy Approximation . . . . .	21
<b>4</b>	<b>Policy Gradient</b>	<b>23</b>
4.1	Basics of Policy Gradient Methods . . . . .	24
4.1.1	Risk-Neutral Framework . . . . .	25
4.1.2	Risk-Sensitive Framework . . . . .	25
4.2	Finite Differences . . . . .	27

---

4.3	Likelihood Ratio Methods . . . . .	27
4.3.1	Monte Carlo Policy Gradient . . . . .	28
4.3.2	GPOMDP . . . . .	30
4.3.3	Risk-Sensitive Monte Carlo Policy Gradient . . . . .	31
4.3.4	Stochastic Policies . . . . .	31
4.3.5	Policy Gradient with Parameter Exploration . . . . .	32
4.4	Risk-Neutral Policy Gradient Theorem . . . . .	36
4.4.1	Theorem Statement and Proof . . . . .	37
4.4.2	GPOMDP . . . . .	39
4.4.3	Actor-Critic Policy Gradient . . . . .	39
4.4.4	Compatible Function Approximation . . . . .	42
4.4.5	Natural Policy Gradient . . . . .	43
4.5	Risk-Sensitive Policy Gradient Theorem . . . . .	45
4.5.1	Risk-Sensitive Actor-Critic Algorithm . . . . .	47
4.6	Parameter-Based Policy Gradient Theorem . . . . .	48
4.6.1	Risk-Neutral Setting . . . . .	48
4.6.2	Risk-Sensitive Setting . . . . .	48
4.6.3	Natural Policy Gradient . . . . .	48
<b>5</b>	<b>Financial Applications of Reinforcement Learning</b>	<b>49</b>
5.1	Bibliographical Survey . . . . .	49
5.2	Asset Allocation . . . . .	49
5.2.1	Reward Function . . . . .	49
5.2.2	States . . . . .	52
5.2.3	Actions . . . . .	52
<b>6</b>	<b>Numerical Results for the Asset Allocation Problem</b>	<b>55</b>
6.1	Synthetic Risky Asset . . . . .	55
6.1.1	Learning Algorithm Specifications . . . . .	56
6.1.2	Experimental Setup . . . . .	57
6.1.3	Risk-Neutral Framework . . . . .	57
6.1.4	Risk-Sensitive Framework . . . . .	58
6.2	Historic Risky Asset . . . . .	61
6.3	Multiple Synthetic Risky Assets . . . . .	61
6.4	Historic Multiple Risky Assets . . . . .	61
<b>7</b>	<b>Conclusions</b>	<b>63</b>
7.1	Summary . . . . .	63
7.2	Further Developments . . . . .	63

# List of Figures

2.1	Agent-environment interaction in sequential decision problems.	5
2.2	Policy iteration algorithm . . . . .	16
3.1	Solution process for the control problem. . . . .	19
4.1	“Vanilla” policy gradient vs. natural policy gradient . . . . .	44
6.1	Risk-neutral learning process for one synthetic risky asset . . .	57
6.2	Backtest performance with one synthetic risky asset . . . . .	58
6.3	Risk-sensitive learning process for one synthetic risky asset . .	59
6.4	Backtest performance with one synthetic risky asset . . . . .	60



# List of Tables

6.1	Risk-neutral backtest statistics with one synthetic risky asset .	59
6.2	Risk-sensitive backtest statistics with one synthetic risky asset	60





# List of Algorithms

4.1	General setup for a policy gradient algorithm. . . . .	24
4.2	Episodic REINFORCE policy gradient estimate . . . . .	30
4.3	Episodic PGPE algorithm . . . . .	34
4.4	Generic structure for an online actor-critic algorithm. . . . .	40
4.5	TD( $\lambda$ ) policy gradient algorithm. . . . .	42



# Chapter 1

## Introduction

1.1 The March of the Machines

1.2 Artificial Intelligence and Finance

1.3 Structure



## Chapter 2

# Discrete-Time Stochastic Optimal Control

### 2.1 Sequential Decision Problems

In sequential decision problems, an agent interacts with an environment by selecting a series of actions in order to complete a specific task. During this interaction, the agent receives a numerical reward from the environment and his goal is to find the best strategy in order to maximize a certain measure of his performance. The environment evolves stochastically and may be influenced by the interaction with the agent, so that each action taken by the agent may influence the circumstances under which future decisions will be made. Therefore, the agent must balance his desire to obtain a large reward today by acting greedily and the opportunities that will be available in the future. While this setting appears quite simple, it is general enough to encompass a wide range of applications in different fields. A classical example is portfolio management, where an investor must allocate his capital so as to maximize his long-term profits. Another standard example is chess, where two players successively move pieces around the chessboard to checkmate the opponent's king.

The purpose of the following sections is to introduce the notation that will be used in the rest of this work and to recall the fundamental concepts and results of the discrete-time stochastic optimal control theory, which is the standard framework to study sequential decisions problems in mathematical terms. Since our discussion will be far from being comprehensive, we refer the reader to the extensive literature on the subject, such as [12], [64], [11].

## 2.2 Markov Decision Processes

A sequential decision problem under uncertainty can be schematized as in Figure 2.1: at a given time  $t$ , the agent (also known as decision maker or controller) observes the state  $s_t$  of the system (also known as environment) and subsequently performs an action  $a_t$ . Following this action, the agent receives an immediate reward  $r_{t+1}$  (or incurs an immediate cost) and the system evolves to a new state according to a probability distribution which depends on the action selected by the agent. At the subsequent time  $t + 1$ , the agent selects a new action given the new state of the system and the process repeats. This interaction can be modeled rigorously using Markov decision processes.

**Definition 2.2.1** (Markov Decision Process). *A Markov decision process (MDP) is a stochastic dynamical system specified by the tuple  $\langle \mathbb{S}, \mathbb{A}, \mathcal{P}, \mathcal{R}, \gamma \rangle$ , where*

- i)  $(\mathbb{S}, \mathcal{S})$  is a measurable space, called the state space.
- ii)  $(\mathbb{A}, \mathcal{A})$  is a measurable space, called the action space.
- iii)  $\mathcal{P} : \mathbb{S} \times \mathbb{A} \times \mathcal{S} \rightarrow \mathbb{R}$  is a Markov transition kernel, i.e.
  - a) for every  $s \in \mathbb{S}$  and  $a \in \mathbb{A}$ ,  $B \mapsto \mathcal{P}(s, a, B)$  is a probability distribution over  $(\mathbb{S}, \mathcal{S})$ .
  - b) for every  $B \in \mathcal{S}$ ,  $(s, a) \mapsto \mathcal{P}(s, a, B)$  is a measurable function on  $\mathbb{S} \times \mathbb{A}$ .
- iv)  $\mathcal{R} : \mathbb{S} \times \mathbb{A} \rightarrow \mathbb{R}$  is a reward function.
- v)  $\gamma \in (0, 1)$  is a discount factor.

Typically, the state space (and similarly the action space) will be either finite, i.e.  $\mathbb{S} = \{s_1, \dots, s_d\}$ , or continuous, i.e.  $\mathbb{S} \subseteq \mathbb{R}^{D_s}$ . The kernel  $\mathcal{P}$  describes the random evolution of the system: suppose that at time  $t$  the system is in state  $s$  and that the agent takes action  $a$ , then, regardless of the previous history of the system, the probability to find the system in a state belonging to  $B \in \mathcal{S}$  at time  $t + 1$  is given by  $\mathcal{P}(s, a, B)$ , i.e.

$$\mathcal{P}(s, a, B) = \mathbb{P}(S_{t+1} \in B | S_t = s, A_t = a) \quad (2.1)$$

Following this random transition, the agent receives a stochastic reward  $R_{t+1}$ . The reward function  $\mathcal{R}(s, a)$  gives the expected reward obtained when action  $a$  is taken in state  $s$ , i.e.

$$\mathcal{R}(s, a) = \mathbb{E}[R_{t+1} | S_t = s, A_t = a] \quad (2.2)$$

This setting can be easily generalized to the following cases

1. The initial state of the system is a random variable  $S_0 \sim \mu$ .
2. The actions that an agent can select depend on the state of the system.

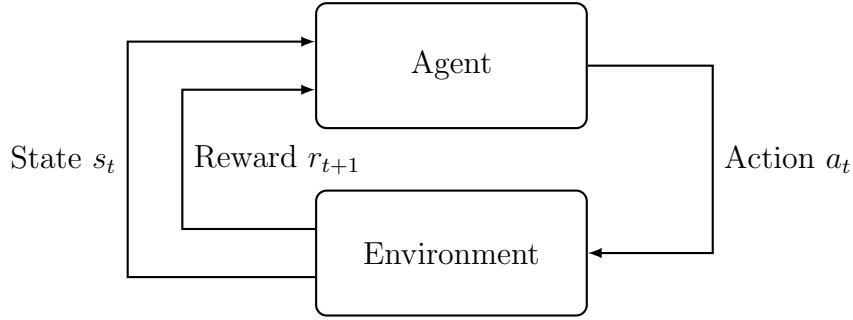


Figure 2.1: Agent-environment interaction in sequential decision problems.

## 2.3 Policies

At any time step, the agent selects his actions according to a certain policy.

**Definition 2.3.1** (Policy). *A policy is a function  $\pi : \mathbb{S} \times \mathcal{A} \rightarrow \mathbb{R}$  such that*

- i) for every  $s \in \mathbb{S}$ ,  $C \mapsto \pi(s, C)$  is a probability distribution over  $(\mathbb{A}, \mathcal{A})$ .*
- ii) for every  $C \in \mathcal{A}$ ,  $s \mapsto \pi(s, C)$  is a measurable function.*

Intuitively, a policy represents a stochastic mapping from the current state of the system to actions. Deterministic policies are a particular case of this general definition. We assumed that the agent's policy is stationary and only depends on the current state of the system. We might in fact consider more general policies that depends on the whole history of the system. However, as we will see, we can always find an optimal policy that depends only on the current state, so that our definition is not restrictive. A policy  $\pi$  and an initial state  $s_0 \in \mathbb{S}$  determine a random state-action-reward sequence  $\{(S_t, A_t, R_{t+1})\}_{t \geq 0}$  with values on  $\mathbb{S} \times \mathbb{A} \times \mathbb{R}$  following the mechanism described above.

**Definition 2.3.2** (History). *Given an initial state  $s_0 \in \mathbb{S}$  and a policy  $\pi$ , a history (or equivalently trajectory or roll-out) of the system is a random sequence  $H_\pi = \{(S_t, A_t)\}_{t \geq 0}$  with values in  $\mathbb{S} \times \mathbb{A}$ , defined on some probability space  $(\Omega, \mathcal{F}, \mathbb{P})$ , such that for  $t = 0, 1, \dots$*

$$\begin{cases} S_0 = s_0 \\ A_t \sim \pi(S_t, \cdot) \\ S_{t+1} \sim \mathcal{P}(S_t, A_t, \cdot) \end{cases}$$

*we will denote by  $(\mathbb{H}, \mathcal{H})$  the measurable space of all possible histories.*

Moreover, we observe that

- i) the state sequence  $\{S_t\}_{t \geq 0}$  is a Markov process  $\langle \mathbb{S}, \mathcal{P}_\pi \rangle$

ii) the state-reward sequence  $\{(S_t, R_t)\}_{t \geq 0}$  is a Markov reward process  $\langle \mathbb{S}, \mathcal{P}_\pi, \mathcal{R}_\pi, \gamma \rangle$  where we denoted

$$\begin{aligned}\mathcal{P}_\pi(s, s') &= \int_{\mathbb{A}} \pi(s, a) \mathcal{P}(s, a, s') da \\ \mathcal{R}_\pi(s) &= \int_{\mathbb{A}} \pi(s, a) \mathcal{R}(s, a) da\end{aligned}$$

In stochastic optimal control, the goal of the agent is to find a policy that maximizes a measure of the agent's long-term performance. In the next sections we discuss some objective functions that are commonly used in the infinite horizon framework.

## 2.4 Risk-Neutral Framework

In the risk-neutral setting, the agent is only interested in maximizing his reward, without considering the risk he needs to take on to achieve it. In an infinite horizon task, the agent's performance is typically measured either as the total discounted reward or as the average reward obtained at each time step. These two approaches are radically different both from a theoretical and an algorithmic point of view. For this reason, they will be always treated separately.

### 2.4.1 Discounted Reward Formulation

In the discounted reward formulation, the agent's performance is measured as the expected return obtained following a specific policy.

**Definition 2.4.1** (Return). *The return is the total discounted reward obtained by the agent starting from  $t$*

$$G_t = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}$$

where  $0 < \gamma < 1$  is the discount factor.

In some domains, such as economics, discounting can be used to represent interest earned on rewards, so that an action that generates an immediate reward will be preferred over one that generates the same reward some steps into the future. Discounting thus models the trade-off between immediate and delayed reward: if  $\gamma = 0$  the agent selects his actions in a myopic way,



while if  $\gamma \rightarrow 1$  he acts in a far-sighted manner. There are other possible reasons for discounting future rewards. The first is because it is mathematically convenient, as it avoids infinite returns and it solves many convergence issues. Another interpretation is that it models the uncertainty about the future, which may not be fully represented. Indeed, the discount factor could be seen as the probability that the world does not stop at a given time step. Since the return is stochastic, we consider its expected value.

**Definition 2.4.2** (State-Value Function). *The state-value function  $V_\pi : \mathbb{S} \rightarrow \mathbb{R}$  is the expected return that can be obtained starting from a state and following policy  $\pi$*

$$V_\pi(s) = \mathbb{E}_\pi [G_t | S_t = s] \quad (2.3)$$

where the subscript in  $\mathbb{E}_\pi$  indicates that all the actions are selected according to policy  $\pi$ . The state-value function measures how good it is for the agent to be in a given state and follow a certain policy. Similarly, we can introduce an action-value function that measures how good it is for the agent to be in a state, take a certain action and then follow the policy.

**Definition 2.4.3** (Action-Value Function). *The action-value function  $Q_\pi : \mathbb{S} \times \mathbb{A} \rightarrow \mathbb{R}$  is the expected return that can be obtained starting from a state, taking an action and then following policy  $\pi$*

$$Q_\pi(s, a) = \mathbb{E}_\pi [G_t | S_t = s, A_t = a] \quad (2.4)$$

We have the following relationship between  $V_\pi$  and  $Q_\pi$

$$V_\pi(s) = \int_{\mathbb{A}} \pi(s, a) Q_\pi(s, a) da \quad (2.5)$$

Almost all reinforcement learning algorithms are designed to estimate these value functions and are typically based on the Bellman equations.

**Proposition 2.4.1** (Bellman Expectation Equations).

$$V_\pi(s) = \mathcal{R}_\pi(s) + \gamma T_\pi V_\pi(s) \quad (2.6)$$

$$Q_\pi(s, a) = \mathcal{R}(s, a) + \gamma T_a V_\pi(s) \quad (2.7)$$

where we denoted by  $T_a$  (resp.  $T_\pi$ ) the transition operator for action  $a$  (resp. for policy  $\pi$ )

$$T_a F(s) = \mathbb{E} [F(S_{t+1}) | S_t = s, A_t = a] = \int_{\mathbb{S}} \mathcal{P}(s, a, s') F(s') ds'$$

$$T_\pi F(s) = \mathbb{E}_\pi [F(S_{t+1}) | S_t = s] = \int_{\mathbb{A}} \pi(s, a) \int_{\mathbb{S}} \mathcal{P}(s, a, s') F(s') ds' da$$

If we introduce the Bellman expectation operator  $B_\pi$ , defined as

$$B_\pi V_\pi(s) = \mathcal{R}_\pi(s) + \gamma T_\pi V_\pi(s)$$

Then Eq. (2.6) can be written as a fixed-point equation

$$V_\pi(s) = B_\pi V_\pi(s)$$

which, under some simple assumptions on the reward functions, admits a unique solution by the contraction mapping theorem. A similar argument applies to Eq. (2.7). The agent's goal is to select a policy  $\pi_*$  that maximizes his expected return in all possible states. Such a policy is called *optimal*.

**Definition 2.4.4** (Optimal State-Value Function). *The optimal state-value function  $V_* : \mathbb{S} \rightarrow \mathbb{R}$  is the largest expected return that can be obtained starting from a state*

$$V_*(s) = \sup_{\pi} V_\pi(s) \quad (2.8)$$

**Definition 2.4.5** (Optimal Action-Value Function). *The optimal action-value function  $Q_* : \mathbb{S} \times \mathbb{A} \rightarrow \mathbb{R}$  is the largest expected return that can be obtained starting from a state and taking an action*

$$Q_*(s, a) = \sup_{\pi} Q_\pi(s, a) \quad (2.9)$$

The optimal value functions satisfy the following Bellman equations.

**Proposition 2.4.2** (Bellman Optimality Equations).

$$V_*(s) = \sup_a Q_*(s, a) = \sup_a \{\mathcal{R}(s, a) + \gamma T_a V_*(s)\} \quad (2.10)$$

$$\begin{aligned} Q_*(s, a) &= \mathcal{R}(s, a) + \gamma T_a V_*(s) \\ &= \mathcal{R}(s, a) + \gamma \int_{\mathbb{S}} \mathcal{P}(s, a, s') \sup_{a'} Q_*(s', a') ds' \end{aligned} \quad (2.11)$$

Again, these two equations are fixed-point equations and the existence and uniqueness of a solution is guaranteed, under some technical assumptions, by the contraction mapping theorem. Starting from the optimal value functions, we can easily derived an optimal policy. Let us define a partial ordering in the policy space

$$\pi \succeq \pi' \Leftrightarrow V_\pi(s) \geq V_{\pi'}(s) \quad \forall s \in \mathbb{S} \quad (2.12)$$

Then the optimal policy  $\pi_* \succeq \pi, \forall \pi$ . We have the following result

**Theorem 2.4.1** (Optimal Policy). *For any Markov decision process,*

- i) *It exists an optimal policy  $\pi_*$  such that  $\pi_* \succeq \pi, \forall \pi$ .*
- ii)  *$V_{\pi_*}(s) = V_*(s)$ .*
- iii)  *$Q_{\pi_*}(s, a) = Q_*(s, a)$ .*

An optimal policy can be found by acting greedily with respect to  $Q_*$ , i.e. in each state  $s$  the agent selects the action that maximizes the action-value function

$$a_* = \arg \sup_a Q_*(s, a) \quad (2.13)$$

We see that this policy is deterministic and only depends on the current state of the system.

## 2.4.2 Average Reward Formulation

Most of the research in RL has studied a problem formulation where agents maximize the cumulative sum of rewards. However, this approach cannot handle infinite horizon tasks, where there are no absorbing goal states, without discounting future rewards. Clearly, discounting is only necessary in cyclical tasks, where the cumulative reward sum can be unbounded. More natural long-term measure of optimality exists for such cyclical tasks, based on maximizing the average reward per action. For a more in-depth presentation, the reader may again refer to the extensive literature on the subject, such as [3], [49] and the references therein. In the average reward setting, also known as long-run reward or ergodic reward, the goal of the agent is to find a policy that maximizes the expected reward per step.

**Definition 2.4.6** (Average Reward). *The average reward  $\rho_\pi$  associated to a policy  $\pi$  is defined as*

$$\begin{aligned} \rho_\pi &= \lim_{T \rightarrow \infty} \frac{1}{T} \mathbb{E}_\pi \left[ \sum_{t=0}^{T-1} R_{t+1} \right] \\ &= \mathbb{E}_{\substack{S \sim d_\pi \\ A \sim \pi}} [\mathcal{R}(S, A)] \\ &= \int_{\mathbb{S}} d_\pi(s) \int_{\mathbb{A}} \pi(s, a) \mathcal{R}(s, a) da ds \end{aligned} \quad (2.14)$$

where  $d_\pi$  is the stationary distribution of the Markov process induced by  $\pi$ .

The agent aims to find an *average optimal* policy

$$\pi_* = \arg \sup_{\pi} \rho_\pi \quad (2.15)$$

In this setting, we introduce the *average adjusted* value and action-value functions.

**Definition 2.4.7** (Average Adjusted State-Value Function). *The average adjusted state-value function  $V_\pi : \mathbb{S} \rightarrow \mathbb{R}$  is the expected residual return that can be obtained starting from a state and following policy  $\pi$*

$$V_\pi(s) = \mathbb{E}_\pi \left[ \sum_{t=0}^{\infty} (R_{t+1} - \rho_\pi) \middle| S_0 = s \right] \quad (2.16)$$

The term  $V_\pi(s)$  is usually referred to as the *bias* value, or the *relative* value, since it represents the relative difference in total reward gained starting from a state  $s$  as opposed to a generic state.  $\rho_\pi$  serves as a baseline that allows to avoid divergence in the value function definition.

**Definition 2.4.8** (Average Adjusted Action-Value Function). *The average adjusted action-value function  $Q_\pi : \mathbb{S} \times \mathbb{A} \rightarrow \mathbb{R}$  is the expected residual return that can be obtained starting from a state, taking an action and then following policy  $\pi$*

$$Q_\pi(s, a) = \mathbb{E}_\pi \left[ \sum_{t=0}^{\infty} (R_{t+1} - \rho_\pi) \middle| S_0 = s, A_0 = a \right] \quad (2.17)$$

We have the following relation between the state-value function and the action-value function

$$V_\pi(s) = \int_{\mathbb{A}} \pi(s, a) Q_\pi(s, a) \quad (2.18)$$

The value functions satisfy the following Bellman equation

**Proposition 2.4.3** (Bellman Expectation Equations).

$$V_\pi(s) = \mathcal{R}_\pi(s) - \rho_\pi + T_\pi V_\pi(s) \quad (2.19)$$

$$Q_\pi(s, a) = \mathcal{R}(s, a) - \rho_\pi + T_a V_\pi(s) \quad (2.20)$$

Again, by introducing opportune Bellman operators, these equations can be rewritten as fixed-point equations. In the discrete case, where the transition operator correspond to matrices, these Bellman equations become linear systems that can be solved to obtain the value functions.

## 2.5 Risk-Sensitive Framework

In many application, in addition to maximizing the average reward, the agent may want to control risk by minimizing some measure of variability in rewards. In the risk-sensitive framework, the goal of the agent is to find the

policy that optimally solves the trade-off between reward and risk. Although risk-sensitive sequential decision-making has a long history in operations research and finance, it has only recently grabbed the attention of the machine learning community. Hence, the literatures offers many reference which approach the risk-sensitive control problem from the traditional stochastic optimal control perspective. On the other hand, there are only few references that attack the problem in the reinforcement learning setting. Again, we can consider the discounted formulation or the average formulation.

### 2.5.1 Discounted Reward Formulation

A standard way to measure the risk associated with a policy  $\pi$  is the variance of the total discounted reward obtained starting from a given state  $s$

$$\Lambda_\pi(s) = \text{Var}_\pi(G_t | S_t = s) \quad (2.21)$$

This approach is the one considered in [74].

**Definition 2.5.1** (Square State-Value Function). *The square state-value function  $U_\pi(s)$  is the second moment of returns that can be obtained starting from a state  $s$  and following policy  $\pi$*

$$U_\pi(s) = \mathbb{E}_\pi[G_t^2 | S_t = s] \quad (2.22)$$

Then, the variance can be decomposed as

$$\Lambda_\pi(s) = U_\pi(s) - V_\pi(s)^2 \quad (2.23)$$

We would like to obtain a Bellman expectation equation for  $U_\pi(s)$ , in order to piggyback on the discussion about the standard state-value function. Let us introduce the square reward function

**Definition 2.5.2** (Square Reward Function). *The square reward function  $\mathcal{M}(s, a)$  is the second moment of the reward that can be obtained in state  $s$  when taking action  $a$*

$$\mathcal{M}(s, a) = \mathbb{E}_\pi[R_{t+1}^2 | S_t = s, A_t = a] \quad (2.24)$$

Let us also define

$$\mathcal{M}_\pi(s) = \mathbb{E}_\pi[R_{t+1}^2 | S_t = s] = \int_{\mathbb{A}} \pi(s, a) \mathcal{M}(s, a) da \quad (2.25)$$

Then, it is easy to show that the square state-value function satisfies the following Bellman equation

**Proposition 2.5.1** (Bellman Expectation Equation).

$$U_\pi(s) = \mathcal{K}_\pi(s) + \gamma^2 T_\pi U_\pi(s) \quad (2.26)$$

where

$$\mathcal{K}_\pi(s) = \mathcal{M}_\pi(s) + 2\gamma \mathbb{E}_\pi [R_{t+1} G_{t+1} | S_t = s] \quad (2.27)$$

This equation is analogous to the Bellman expectation equation for the state-value function, with a synthetic reward function  $\mathcal{K}_\pi(s)$  and a discount factor  $\gamma^2$ . However, this equation is completely useless in practice because the second term in Eq. (2.27), which represents the influence of the first day reward on all the future rewards, is very hard, if not impossible, to estimate. In [80] and [63], the authors implicitly assume that the reward  $R_{t+1}$  is conditionally independent from the future rewards given the current state. Exploiting this assumption, it is easy to show that the return variance satisfies the following Bellman expectation equation

**Proposition 2.5.2** (Bellman Expectation Equation).

$$\Lambda_\pi(s) = \mathcal{V}_\pi(s) + \gamma^2 \text{Var}_\pi (V_\pi(S_{t+1}) | S_t = s) + \gamma^2 T_\pi \Lambda_\pi(s) \quad (2.28)$$

where  $\mathcal{V}_\pi(s)$  denotes the conditional variance of the reward

$$\mathcal{V}_\pi(s) = \text{Var}_\pi (R_{t+1} | S_t = s) \quad (2.29)$$

*Proof.*

$$\begin{aligned} \Lambda_\pi(s) &= \text{Var}_\pi (G_t | S_t = s) \\ &= \text{Var}_\pi (R_{t+1} + \gamma G_{t+1} | S_t = s) \\ &= \mathcal{V}_\pi(s) + 2\gamma \text{Cov}_\pi (R_{t+1}, G_{t+1} | S_t = s) + \gamma^2 \text{Var}_\pi (G_{t+1} | S_t = s) \end{aligned}$$

Under the independence hypothesis, the second term is null. Applying the law of total variance we obtain

$$\begin{aligned} \text{Var}_\pi (G_{t+1} | S_t = s) &= \mathbb{E}_\pi [\text{Var}_\pi (G_{t+1} | S_{t+1}) | S_t = s] + \text{Var}_\pi (\mathbb{E}_\pi [G_{t+1} | S_{t+1}] | S_t = s) \\ &= T_\pi \Lambda_\pi(s) + \text{Var}_\pi (V_\pi(S_{t+1}) | S_t = s) \end{aligned}$$

Plugging it into the first equality yields the result

$$\Lambda_\pi(s) = \text{Var}_\pi (R_{t+1} | S_t = s) + \text{Var}_\pi (V_\pi(S_{t+1}) | S_t = s) + \gamma^2 T_\pi \Lambda_\pi(s)$$

□

Again, even though this result is appealing from a theoretical point of view, it cannot be exploited to derive a reinforcement learning algorithm. Indeed, the variance of the state-value function would be extremely difficult or even impossible to estimate. For this reason, we will not develop further the discounted reward formulation in the risk-sensitive setting and we will focus instead on the average reward formulation.

## 2.5.2 Average Reward Formulation

In [63], the authors consider the long-run variance as a measure of the risk associated to a policy  $\pi$

**Definition 2.5.3** (Long-Run Variance). *The long-run variance  $\Lambda_\pi$  under policy  $\pi$  is defined as*

$$\Lambda_\pi = \lim_{T \rightarrow \infty} \frac{1}{T} \mathbb{E}_\pi \left[ \sum_{t=0}^{T-1} (R_{t+1} - \rho_\pi)^2 \right] \quad (2.30)$$

The long-run variance can be decomposed as follows

$$\Lambda_\pi = \eta_\pi - \rho_\pi^2 \quad (2.31)$$

where  $\eta_\pi$  is the average square reward per step

**Definition 2.5.4** (Average Square Reward).

$$\begin{aligned} \eta_\pi &= \lim_{T \rightarrow \infty} \frac{1}{T} \mathbb{E}_\pi \left[ \sum_{t=0}^{T-1} R_{t+1}^2 \right] \\ &= \mathbb{E}_{\substack{S \sim d_\pi \\ A \sim \pi}} [\mathcal{M}(S, A)] \\ &= \int_{\mathbb{S}} d_\pi(s) \int_{\mathbb{A}} \pi(s, a) \mathcal{M}(s, a) \end{aligned} \quad (2.32)$$

where we denoted by  $\mathcal{M}(s, a)$  the square reward function

$$\mathcal{M}(s, a) = \mathbb{E} [R_{t+1}^2 | S_t = s, A_t = a] \quad (2.33)$$

As before, we introduce the residual state-value and action-value functions associated with the square reward under policy  $\pi$

**Definition 2.5.5** (Average Adjusted Square State-Value Function). *The average adjusted square state-value function  $U_\pi : \mathbb{S} \rightarrow \mathbb{R}$  is the expected square residual return that can be obtained starting from a state and following policy  $\pi$*

$$U_\pi(s) = \mathbb{E}_\pi \left[ \sum_{t=0}^{\infty} (R_{t+1}^2 - \eta_\pi) \mid S_0 = s \right] \quad (2.34)$$

**Definition 2.5.6** (Average Adjusted Square Action-Value Function). *The average adjusted square action-value function  $Q_\pi : \mathbb{S} \times \mathbb{A} \rightarrow \mathbb{R}$  is the expected square residual return that can be obtained starting from a state, taking an action and then following policy  $\pi$*

$$W_\pi(s, a) = \mathbb{E}_\pi \left[ \sum_{t=0}^{\infty} (R_{t+1}^2 - \eta_\pi) \mid S_0 = s, A_0 = a \right] \quad (2.35)$$

The following relation between square state-value function and the square action-value holds

$$U_\pi(s) = \int_{\mathbb{A}} \pi(s, a) W_\pi(s, a) \quad (2.36)$$

The average adjusted square value functions satisfy the following Bellman equations

**Proposition 2.5.3** (Bellman Expectation Equations).

$$U_\pi(s) = \mathcal{M}_\pi(s) - \eta_\pi + T_\pi U_\pi(s) \quad (2.37)$$

$$\begin{aligned} W_\pi(s, a) &= \mathcal{M}(s, a) - \eta_\pi + T_a U_\pi(s) \\ &= \mathcal{M}(s, a) - \eta_\pi + \int_{\mathbb{S}} \mathcal{P}(s, a, s') \int_{\mathbb{A}} \pi(s', a') W_\pi(s', a') da' ds' \end{aligned} \quad (2.38)$$

In the risk-sensitive setting, the agent wants to find a policy that solves the following mean-variance optimization problem

$$\begin{cases} \max_{\pi} \rho_\pi \\ \text{subject to } \Lambda_\pi \leq \alpha \end{cases} \quad (2.39)$$

for a given  $\alpha > 0$ . Using the Lagrangian relaxation procedure, we can recast (2.39) to the following unconstrained problem

$$\max_{\lambda} \min_{\pi} L(\pi, \lambda) = -\rho_\pi + \lambda(\Lambda_\pi - \alpha) \quad (2.40)$$

Alternatively, the agent may want to optimize the Sharpe ratio, a risk-sensitive performance measure commonly used in finance

$$\text{Sh}_\pi = \frac{\rho_\pi}{\sqrt{\Lambda_\pi}} \quad (2.41)$$



## 2.6 Dynamic Programming Algorithms

The Bellman equations provide the basis for the *value iteration* and *policy iteration* algorithms [77], which simply consist in iteratively applying the Bellman operators to an approximation of the value functions. As we will see, these methods require knowledge of the MDP dynamics and typically they are only applicable to finite state MDPs, for which the value functions can be represented as vectors and the Bellman operators becomes matrices. However, these simple algorithm provide some useful insight that can be exploited to develop more advanced reinforcement learning algorithm

### 2.6.1 Value Iteration

*Value iteration* is an iterative method to compute the optimal state-value function  $V_*(s)$ . Starting from an arbitrary function  $V_0(s)$ , the algorithm iteratively updates the approximation by applying the Bellman optimality operator in a fixed-point scheme

$$V_{k+1}(s) = B_* V_k(s) = \sup_a \{ \mathcal{R}(s, a) + \gamma T_a V_*(s) \}$$

This algorithm is guaranteed to converge to the optimal value function  $V_*$  by the contraction mapping theorem. Let us notice that the transition operator  $T_a$  requires knowledge of the MDP dynamics, which is not available in the typical reinforcement learning framework. Moreover, the update involves an optimization with respect to all possible actions which can be carried out efficiently only in the case of finite action spaces. A drawback of this algorithm is that it does not provide an explicit representation of the optimal policy, which in many control problems is what we are looking for.

### 2.6.2 Policy Iteration

Policy iteration is an iterative method to approximate both the optimal state-value function  $V_*$  and the optimal policy  $\pi_*$ . This algorithm alternates an evaluation step, in which the current policy is evaluated using the state-value function, and an improvement step, in which the policy is improved by acting greedily with respect to the action-value function computed in the evaluation step. In the standard version of policy iteration, the state-value function for the current policy is evaluated starting from an arbitrary function  $V_0(s)$  and iteratively applying the Bellman expectation operator in a fixed-point iteration scheme

$$V_{k+1}(s) = B_{\pi_n} V_k(s)$$

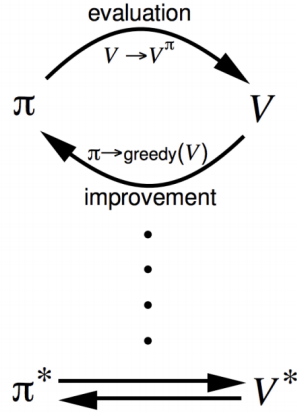


Figure 2.2: Policy iteration algorithm [75].

This algorithm is guaranteed to converge to  $V_{\pi_n}$ . The new policy is then computed as

$$\pi_{n+1} = \text{greedy}(V_{\pi_n})$$

and we go back to the evaluation step for this new policy. The algorithm is guaranteed to converge to  $V_*$  and  $\pi_*$ . Let us notice that it is not necessary to perfectly evaluate to policy  $\pi_n$  before performing the improvement step and the evaluation procedure can be stopped before convergence. This algorithm suffers from the same problem as above. However it provides the basic structure for most of the value-based reinforcement learning methods. In particular, in *generalized policy iteration* the evaluation step is performed in an arbitrary way which does not necessarily employ the Bellman operator. This scheme is illustrated in Figure 2.2,

## Chapter 3

# Reinforcement Learning

The standard way to solve Markov decision processes is through dynamic programming, which simply consists in solving the Bellman fixed-point equations discussed in the previous chapter. Following this approach, the problem of finding the optimal policy is transformed into the problem of finding the optimal value function. However, apart from the simplest cases where the MDP has a limited number of states and actions, dynamic programming becomes computationally infeasible. Moreover, this approach requires complete knowledge of the Markov transition kernel and of the reward function, which in many real-world applications might be unknown or too complex to use. Reinforcement Learning (RL) is a subfield of Machine Learning which aims to turn the infeasible dynamic programming methods into practical algorithms that can be applied to large-scale problems. RL algorithms are based on two key ideas: the first is to use samples to compactly represent the unknown dynamics of the controlled system. The second idea is to use powerful function approximation methods to compactly estimate value functions and policies in high-dimensional state and action spaces. In the following sections, we present the reinforcement learning problem in more depth and discuss how it relates to the standard discrete-time stochastic optimal control theory and to the other subfields of machine learning, such as supervised learning. In particular, we will see how RL extends ideas from optimal control theory and stochastic approximation to address the broader goal of artificial intelligence. This quick overview of RL is propaedeutic to the following chapters, where we will present in more detail a particular class of algorithms, called policy gradient methods, which are well suited for continuous action spaces. For a more thorough presentation, the reader may consult [75], [77] or [85].

### 3.1 The Reinforcement Learning Problem

Reinforcement Learning (RL) is a general class of algorithms in the field of machine learning that allows an agent to learn how to behave in a stochastic and possibly unknown environment, where the only feedback consists of a scalar reward signal. In order to maximize the long-run reward, the agent must learn which actions are the most profitable by trial-and-error. Therefore, RL algorithms can be seen as computational methods to solve Markov decision processes by directly interacting with the environment, for which a model may or may not be available. Trial-and-error search and a delayed reward signal can be seen as the most characteristic features of reinforcement learning.

Compared to supervised learning, one of the main branches of machine learning, the feedback the learner receives is much less. In supervised learning, the agent is provided with examples of the correct or expected behavior by a knowledgeable external supervisor and his goal is to learn how to replicate these examples as well as possible and possibly generalize this knowledge to new examples. In reinforcement learning, the agent only receives a numerical reward that only gives a partial feedback of the goodness of his actions. Therefore, this feedback system is evaluative rather than instructive and it is much more difficult for the agent to learn how to behave in uncharted territory without any external guidance.

This particular framework generates some challenges that are not present in other kinds of learning. The first one is the trade-off between exploration and exploitation. In order to maximize his reward, an agent would greedily select actions that he has already tried in the past and found to be effective in producing rewards. However, to find these actions, he has to test actions that haven't been selected before in order to evaluate their potential. Clearly, this might result in worse performance in the short-term because the actions might be suboptimal. However, without trying them, the agent might not be able to find possible improvements. Thus, an agent must exploit what he already knows to obtain rewards, but he also needs to explore to select better actions in the future. A second challenge is the credit assignment problem. Since rewards might be delayed in time, it will be difficult for the agent to understand which actions are mostly responsible for the outcome.

### 3.2 Model-Free RL Methods

In Section 2.6, we discussed the policy iteration method for computing an optimal policy for an MDP in a finite state and action spaces. This algorithm belongs to the class of *model-based* methods, since it requires perfect

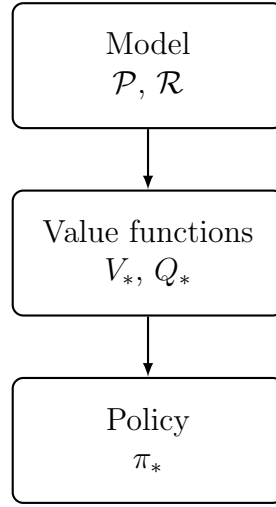


Figure 3.1: Solution process for the control problem.

knowledge of the Markov transition kernel and reward function, i.e. the model of the MDP. RL is primarily concerned with how to obtain an optimal policy when such a model is not available. In this section, we discuss some classes of *model-free* methods which do not rely on the MDP model. The lack of a model generates the need to sample the MDP to gather statistical knowledge about this unknown model. In the control setting, the goal is to approximate the optimal policy, which depends on the optimal value, which in turn depends on the model of the MDP, as shown in Figure 3.1. Indeed, we have already seen that a policy which is greedy with respect to the optimal action-value function  $Q_*$ , i.e. such that

$$\int_{\mathbb{A}} \pi_*(s, a) Q_*(s, a) da = \sup_a Q_*(s, a) \quad (3.1)$$

is optimal. Therefore, we can derive the following three methodologies that differ in which part of the solution process is approximated

- i) *Model-approximation* algorithms approximate the MDP model and compute an estimate of the optimal policy by dynamic programming.
- ii) *Value-approximation* algorithms use samples to directly approximate  $V_*$  or  $Q_*$ , from which an estimate of  $\pi_*$  can be derived by acting greedily.
- iii) *Policy-approximation* algorithms directly try to estimate the optimal policy.

It should be noticed that these approaches are not mutually exclusive and can be combined to derive hybrid algorithms. In the following sections we discuss these three classes of algorithm in more detail, following closely [85].

### 3.2.1 Model Approximation

Model-approximation algorithms approximate the MDP model and compute an optimal policy by dynamic programming, using the techniques discussed in the previous chapter. Since  $\mathbb{S}$ ,  $\mathbb{A}$  and  $\gamma$  are assumed to be known, these methods are based on learning an approximation of the Markov transition kernel  $\mathcal{P}$  and the reward function  $\mathcal{R}$ . Thanks to the Markov property, these quantities only depend on the current state and action, so that their approximation corresponds to a density estimation problem and a regression problem respectively, which are fairly standard supervised learning problems. Learning the model may not be trivial, but it is in general easier than learning the value of a policy or optimizing the policy directly. The major drawback of model-based algorithms in continuous-state MDPs is that, even if the model is available, it is in general infeasible to compute the value functions by dynamic programming and to extract an optimal policy for all states by acting greedily. Alternatively, a transition model estimate may be used to generate sample trajectories from the MDP, which can then be used to estimate the value function or directly improve the policy. However, the value function and the policy estimated using these samples can only be as accurate as the learned model, so that in many cases it may be easier to directly approximate the value function or the policy using the methods described below.

### 3.2.2 Value Approximation

Value-approximation algorithms use samples from the MDP to approximate  $V_*$  or  $Q_*$  directly and then derive an estimate of the optimal policy by acting greedily with respect to  $Q_*$ . Typically, when the state and action spaces are large, the value functions are estimated using a parametric function approximator, whose parameters are iteratively updated given the observed samples. Many reinforcement learning algorithms fall into this category and they can be distinguished based on whether they are on-policy or off-policy and whether they update the value function estimates online or offline. *On-policy* algorithms approximate the state-value function  $V_\pi$  or the action-value function  $Q_\pi$  from samples of the MDP obtained by following the same policy  $\pi$  to be evaluated. Although the optimal policy  $\pi_*$  is initially unknown, such algorithms can eventually approximate the optimal value functions  $V_*$  or  $Q_*$  from which an approximation of the optimal policy can be derived. On the other hand, *off-policy* algorithms can learn the value of a policy different from the one used for obtaining the MDP samples. *Online* algorithms adapt their value approximation after each observed sample while *Offline* algorithms operate on batches of samples. Online algorithms typically re-

quire less computation per sample but their convergence is slower. Online on-policy algorithms include temporal-difference (TD) algorithms, such as TD-learning, Sarsa. Offline on-policy algorithms include least-squares approaches, such least-squares temporal difference (LSTD), least-squares policy evaluation (LSPE) and least-squares policy iteration (LSPI). The most known model-free online off-policy algorithm is Q-learning.

### 3.2.3 Policy Approximation

Policy-approximation algorithms only store a policy and update this policy to maximize a given performance measure and eventually converge to the optimal policy. Since these algorithms only use a policy estimate and don't rely on a value function approximation, they are also referred to as *direct policy-search* or *actor-only* algorithms. Algorithms that store both a policy and a value function are commonly known as *actor-critic* algorithms [76], [43]. Policy gradient algorithms, the most studied policy-approximation methods, will be discussed in much more detail in the next chapter.





# Chapter 4

## Policy Gradient

Policy gradient methods directly store and iteratively improve a parametric approximation of the optimal policy. This policy is commonly referred to as an *actor* and methods that directly approximate the policy without exploiting an approximation of the optimal value function are called *actor-only*. Algorithms that combine an approximation of the optimal policy with an approximation of the value function are commonly called *actor-critic* methods. As discussed in Chapter 2, an optimal policy can be obtained by simply acting greedily with respect to the optimal action-value function. However, in large or continuous action spaces, this leads to a complex optimization problem that is computationally expensive to solve. Therefore, it can be beneficial to store an explicit estimation of the optimal policy from which we can select actions. Policy gradient methods have other advantages compared to standard value-based approaches. In many applications, a good policy has a more compact representation than the value function so that it may be easier to approximate. Moreover, the policy parameterization can be chosen so as to be relevant for the task and to directly incorporate prior knowledge. Another advantage is that these methods can learn stochastic policies and not only deterministic ones, which might be useful in multiplayer frameworks or in partially observable environments. Finally, these methods have better convergence properties and are guaranteed to converge at least to a local optimum, which may be good enough in practice. On the other hand, policy gradient methods are typically characterized by a large variance which may hinder the converge speed.

In the following sections, we present the basics of policy gradient methods closely following the thorough overviews provided in [62]. Then, we discuss some state-of-the-art policy gradient algorithms and we propose some original extensions to these methods.

## 4.1 Basics of Policy Gradient Methods

In policy gradient methods, the optimal policy is approximated using a parametrized policy  $\pi : \mathbb{S} \times \mathcal{A} \times \Theta \rightarrow \mathbb{R}$  such that, given a parameter vector  $\theta \in \Theta \subseteq \mathbb{R}^{D_\theta}$ ,  $\pi(s, B; \theta) = \pi_\theta(s, B)$  gives the probability of selecting an action in  $B \in \mathcal{A}$  when the system is in state  $s \in \mathbb{S}$ . The general goal of policy optimization in reinforcement learning is to optimize the policy parameters  $\theta \in \Theta$  so as to maximize a certain objective function  $J : \Theta \rightarrow \mathbb{R}$

$$\theta^* = \arg \max_{\theta \in \Theta} J(\theta)$$

In the following, we will focus on gradient-based and model-free methods that exploit the sequential structure of the reinforcement learning problem. The idea of policy gradient algorithms is to update the policy parameters using the gradient ascent direction of the objective function

$$\theta^{k+1} = \theta^k + \alpha_k \nabla_{\theta} J(\theta^k) \quad (4.1)$$

where  $\{\alpha_k\}_{k \geq 0}$  is a sequence of learning rates. Typically, the gradient of the objective function is not known and its approximation is the key component of every policy gradient algorithm, a general setup of which is schematized in Algorithm 4.1. It is a well-know result from stochastic optimization [44] that, if the gradient estimate is unbiased and the learning rates satisfy the *Robbins-Monro conditions*

$$\sum_{k=0}^{\infty} \alpha_k = \infty \quad \sum_{k=0}^{\infty} \alpha_k^2 < \infty \quad (4.2)$$

the learning process is guaranteed to converge at least to a local optimum of the objective function. Before describing various methods to approximate the gradient, let us give an overview of some standard objective functions and the situation in which they are commonly used.

---

**Algorithm 4.1** General setup for a policy gradient algorithm.

---

**Input:** Initial parameters  $\theta_0$ , learning rate  $\{\alpha_k\}$

**Output:** Approximation of the optimal policy  $\pi_*$

- 1: **repeat**
  - 2:     Approximate policy gradient  $\hat{g} \approx \nabla_{\theta} J(\theta^k)$
  - 3:     Update parameters using gradient ascent  $\theta^{k+1} = \theta^k + \alpha_k \hat{g}$
  - 4:      $k \leftarrow k + 1$
  - 5: **until** converged
-

### 4.1.1 Risk-Neutral Framework

In the risk-neutral setting, the agent is only interested in maximizing his reward, without considering the risk he needs to take on to achieve it. In an episodic environment where the system always starts from an initial state  $s_0$ , the typical objective function is the start value.

**Definition 4.1.1** (Start Value). *The start value is the expected return that can be obtained starting from the start state  $s_0 \in \mathbb{S}$  and following policy  $\pi_\theta$*

$$J_{start}(\theta) = V_{\pi_\theta}(s_0) = \mathbb{E}_{\pi_\theta} [G_0 \mid S_0 = s_0] \quad (4.3)$$

In a continuing environment, where no terminal state exists and the task might go on forever, it is common to use either the average value or the average reward per time step.

**Definition 4.1.2** (Average Value). *The average value is the expected value that can be obtained following policy  $\pi_\theta$*

$$J_{avV}(\theta) = \mathbb{E}_{S \sim d^\theta} [V_{\pi_\theta}(S)] = \int_{\mathbb{S}} d^\theta(s) V_{\pi_\theta}(s) ds \quad (4.4)$$

where  $d^\theta$  is the stationary distribution of the Markov chain induced by  $\pi_\theta$ .

**Definition 4.1.3** (Average Reward per Time Step). *The average reward per time step is the expected reward that can be obtained over a single time step by following policy  $\pi_\theta$*

$$J_{avR}(\theta) = \rho(\theta) = \mathbb{E}_{\substack{S \sim d^\theta \\ A \sim \pi_\theta}} [\mathcal{R}(S, A)] = \int_{\mathbb{S}} d^\theta(s) \int_{\mathbb{A}} \pi_\theta(s, a) \mathcal{R}(s, a) da ds \quad (4.5)$$

where  $d^\theta$  is the stationary distribution of the Markov chain induced by  $\pi_\theta$ .

### 4.1.2 Risk-Sensitive Framework

In the risk-sensitive framework, in addition to maximizing his reward, the agent also wants to control the risk he takes on to achieve it. Thus, it is necessary to introduce a function  $\Lambda : \Theta \rightarrow \mathbb{R}$  such that  $\Lambda(\theta)$  measures the risk associated with the policy  $\pi_\theta$ . In an episodic framework where the system always starts from the same initial state  $s_0$ , the variance of the total return can be used [80]

**Definition 4.1.4** (Start Variance). *The start variance is the variance of the return that can be obtained starting from the start state  $s_0 \in \mathbb{S}$  and following policy  $\pi_\theta$*

$$\Lambda_{start}(\theta) = \text{Var}_{\pi_\theta} (G_0 \mid S_0 = s_0) = U_{\pi_\theta}(s_0) - V_{\pi_\theta}(s_0)^2 \quad (4.6)$$

In a continuing environment, we might consider the long-run variance [63] defined in Section 2.5

$$\Lambda(\theta) = \lim_{T \rightarrow \infty} \frac{1}{T} \mathbb{E}_{\pi_\theta} \left[ \sum_{t=0}^{T-1} (R_{t+1} - \rho(\theta))^2 \right] = \eta(\theta) - \rho(\theta)^2 \quad (4.7)$$

In order to formalize this trade-off between reward and risk from a mathematical point of view, we borrow from the financial literature two standard risk-sensitive performance measures: the mean-variance criterion [50] and the Sharpe ratio [72].

**Definition 4.1.5** (Mean-Variance Criterion). *The mean-variance criterion is defined as*

$$J_{MV}(\theta) = J(\theta) - \chi \Lambda(\theta) \quad (4.8)$$

where  $\chi \in \mathbb{R}$  is a constant that controls the reward-risk trade-off.

**Definition 4.1.6** (Sharpe Ratio). *The Sharpe ratio is defined as*

$$\text{Sh}(\theta) = \frac{J(\theta)}{\sqrt{\Lambda(\theta)}} \quad (4.9)$$

In a risk-sensitive policy gradient algorithm, we try to approximate the optimal parameters that maximize these objective functions by updating the parameters following the gradient ascent directions. In the average-reward formulation, the ascent directions are given by

$$\nabla_\theta J_{MV}(\theta) = \nabla_\theta \rho(\theta) - \chi \nabla_\theta \Lambda(\theta) \quad (4.10)$$

for the mean-variance criterion, where

$$\nabla_\theta \Lambda(\theta) = \nabla_\theta \eta(\theta) - 2\rho(\theta) \nabla_\theta \rho(\theta) \quad (4.11)$$

and by

$$\nabla_\theta \text{Sh}(\theta) = \frac{\eta(\theta) \nabla_\theta \rho(\theta) - \frac{1}{2} \rho(\theta) \nabla_\theta \eta(\theta)}{\Lambda(\theta) \sqrt{\Lambda(\theta)}} \quad (4.12)$$

for the Sharpe ratio. Hence, in the risk-sensitive framework it is necessary to estimate to different gradients:  $\nabla_\theta \rho(\theta)$  and  $\nabla_\theta \eta(\theta)$ . Let us notice that the equivalent expressions for the episodic setting can be obtained by replacing  $\rho(\theta)$  (resp.  $\eta(\theta)$ ) with  $V_\theta(s_0)$  (resp.  $U_\theta(s_0)$ ). In the next sections, we will discuss several techniques to estimate these gradients.

## 4.2 Finite Differences

The most straightforward way to estimate the objective function gradient consists in replacing the partial derivatives with the corresponding finite differences. In other words, the  $k$ -th gradient component is approximated by

$$\frac{\partial J(\theta)}{\partial \theta_k} \approx \frac{J(\theta + \epsilon e_k) - J(\theta)}{\epsilon} \quad (4.13)$$

where  $\epsilon$  is a small constant and  $e_k \in \mathbb{R}^{D_\theta}$  is the  $k$ -th element of the canonical basis. Hence, to compute the gradient it is sufficient to estimate the objective function for  $D_\theta + 1$  different parameters combinations. Since the objective function is unknown, it must be estimated from sample trajectories simulated following the policies associated to the parameterizations appearing in the finite differences. For this reason, this method is computationally demanding and the gradient estimate obtained in this way is extremely noisy, which may slow down or even prevent the convergence of the algorithm. Still, this approach is sometimes effective and works for arbitrary, even non-differentiable, policies. Moreover, finite differences are often used to check gradient estimates during debugging.

## 4.3 Likelihood Ratio Methods

In this section, we describe several methods based on the *likelihood ratio* technique of stochastic optimization. Let  $Z$  be a random variable with a parametric probability density  $p_\theta$  and assume that  $p_\theta$  is known, explicitly computable and differentiable with respect to the parameters. The likelihood ratio technique is used to compute the following gradient

$$\nabla_\theta \mathbb{E}_{Z \sim p_\theta} [f(Z)] = \nabla_\theta \int p_\theta(z) f(z) dz$$

This type of problem appears in many domains, such as when computing the greeks of a derivative product in computational finance [59]. Under some regularity assumptions on the function  $f$  and on the probability distribution  $p_\theta$ , this gradient can be rewritten in a more amenable way

$$\begin{aligned} \nabla_\theta \mathbb{E}_{Z \sim p_\theta} [f(Z)] &= \int \nabla_\theta p_\theta(z) f(z) dz \\ &= \int p_\theta(z) \frac{\nabla_\theta p_\theta(z)}{p_\theta(z)} f(z) dz \\ &= \int p_\theta(z) \nabla_\theta \log p_\theta(z) f(z) dz \\ &= \mathbb{E}_{Z \sim p_\theta} [\nabla_\theta \log p_\theta(Z) f(Z)] \end{aligned}$$

Thus, the likelihood ratio technique simply consists in the following equality

$$\nabla_{\theta} \mathbb{E}_{Z \sim p_{\theta}} [f(Z)] = \mathbb{E}_{Z \sim p_{\theta}} [\nabla_{\theta} \log p_{\theta}(Z) f(Z)] \quad (4.14)$$

where  $\nabla_{\theta} \log p_{\theta}(Z)$  is usually called likelihood score. This result provides a simple way to approximate the gradient using a Monte Carlo estimate: let  $\{Z^{(m)}\}_{m=1}^M$  be i.i.d. samples from  $p_{\theta}$ , then

$$\nabla_{\theta} \mathbb{E}_{Z \sim p_{\theta}} [f(Z)] \approx \frac{1}{M} \sum_{m=1}^M \nabla_{\theta} \log p_{\theta}(Z^{(m)}) f(Z^{(m)}) \quad (4.15)$$

As we will see in the following sections, this technique provides the basis for some of the most common policy gradient algorithms in the reinforcement learning literature.

### 4.3.1 Monte Carlo Policy Gradient

Let  $h = \{(s_t, a_t)\}_{t \geq 0} \in \mathbb{H}$  be a given trajectory of the MDP and let us denote by  $p_{\theta}(h) = \mathbb{P}_{\pi_{\theta}}(\bar{H} = h)$  the probability of obtaining this trajectory under policy  $\pi_{\theta}$ . Let  $G(h)$  denote the expected return obtained on trajectory  $h$

$$G(h) = \mathbb{E}[G_0 | H = h] = \sum_{t=0}^{\infty} \gamma^t \mathcal{R}(s_t, a_t)$$

Let us consider the average value objective function, which can be rewritten as an expectation over all possible trajectories

$$J_{\text{avV}}(\theta) = \mathbb{E}_{H \sim p_{\theta}} [G(H)]$$

Applying the likelihood ratio technique, we obtain

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{H \sim p_{\theta}} [\nabla_{\theta} \log p_{\theta}(H) G(H)] \quad (4.16)$$

The crucial point is that  $\nabla_{\theta} \log p_{\theta}(H)$  can be computed without knowledge of the transition probability kernel  $\mathcal{P}$ . Indeed, by a recursive application of the Markov property, we have

$$p_{\theta}(h) = \mathbb{P}(S_0 = s_0) \prod_{t=0}^{\infty} \pi_{\theta}(s_t, a_t) \mathcal{P}(s_t, a_t, s_{t+1})$$

from which we obtain

$$\log p_{\theta}(h) = \log \mathbb{P}(S_0 = s_0) + \sum_{t=0}^{\infty} \log \pi_{\theta}(s_t, a_t) + \sum_{t=0}^{\infty} \log \mathcal{P}(s_t, a_t, s_{t+1})$$

Since the only term depending on the parameters  $\theta$  is the policy term,

$$\nabla_{\theta} \log p_{\theta}(H) = \sum_{t=0}^{\infty} \nabla_{\theta} \log \pi_{\theta}(s_t, a_t) \quad (4.17)$$

Therefore, we do not need the transition model to compute the  $\nabla_{\theta} \log p_{\theta}(H)$ . Moreover, it is easy to prove that

$$\mathbb{E}_{H \sim p_{\theta}} [\nabla_{\theta} \log p_{\theta}(H)] = 0 \quad (4.18)$$

Hence, a constant baseline  $b \in \mathbb{R}$  can always be subtracted in Eq. 4.16 without changing the gradient value

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{H \sim p_{\theta}} [\nabla_{\theta} \log p_{\theta}(H)(G(H) - b)] \quad (4.19)$$

Intuitively, the baseline  $b$  should measure the expected return under the policy, so that better (resp. worse) than average returns would produce a positive (resp. negative) gradient. A simple approach is to use a moving average of the returns observed while improving the policy. A more elaborate approach is to compute the baseline that minimizes the variance of the gradient estimator.

In an episodic environment, we can derive an estimate of the gradient by sampling  $M$  trajectories  $h^{(m)} = \{(s_t^{(m)}, a_t^{(m)})\}_{t=0}^{T^{(m)}}$  under policy  $\pi_{\theta}$  and by approximating the expected value via Monte Carlo

$$\hat{g}_{\text{RF}} = \frac{1}{M} \sum_{m=1}^M \left[ \sum_{i=0}^{T^{(m)}} \nabla_{\theta} \log \pi_{\theta}(s_i^{(m)}, a_i^{(m)}) \right] \left[ \sum_{j=0}^{T^{(m)}} \gamma^j r_{j+1}^{(m)} - b \right] \quad (4.20)$$

This method, synthetized in Algorithm 4.2, is known in the literature as the REINFORCE algorithm and is guaranteed to converge to the true gradient at a pace of  $O(M^{-1/2})$ . In practice, we can obtain an approximation of the gradient using only one sample which leads to a stochastic gradient ascent method

$$\hat{g}_{\text{SRF}} = \left[ \sum_{i=0}^T \nabla_{\theta} \log \pi_{\theta}(s_i, a_i) \right] \left[ \sum_{j=0}^T \gamma^j r_{j+1} - b \right] \quad (4.21)$$

This method is very easy and works well on many problems. However, the gradient estimate is characterized by a large variance which can hamper the convergence rate of the algorithm. A first approach to address this issue is to optimally set the baseline to reduce the gradient variance.

**Algorithm 4.2** Episodic REINFORCE policy gradient estimate**Input:** Policy parameterization  $\theta$ , number of trajectories  $M$ **Output:** REINFORCE policy gradient estimate  $\hat{g}_{RF} \approx \nabla_{\theta} J(\theta)$ 

- 1: Sample  $M$  trajectories of the MDP following policy  $\pi_{\theta}$
- 2: For all  $k$ , Compute the optimal baseline  $b_k^*$  according to Eq. (4.22)
- 3: Compute  $\hat{g}_{RF}$  according to Eq. (4.20)

**4.3.1.1 Optimal Baseline**

A standard variance reduction technique, called control variate, consists in setting the baseline so as to minimize the gradient estimate variance. More in detail, the optimal baseline for the  $k$ -th gradient component  $\hat{g}_k$  solves

$$b_k^* = \arg \min_b \text{Var}(\hat{g}_k)$$

It is easy to show that

$$b_k^* = \frac{\mathbb{E}[G(H) (\partial_{\theta_k} \log p_{\theta}(H))^2]}{\mathbb{E}[(\partial_{\theta_k} \log p_{\theta}(H))^2]}$$

which can be approximated by

$$\hat{b}_k^* = \frac{\sum_{m=1}^M \left[ \sum_{i=0}^{T^{(m)}} \partial_{\theta_k} \log \pi_{\theta}(s_i^{(m)}, a_i^{(m)}) \right]^2 \sum_{j=0}^{T^{(m)}} \gamma^j r_{j+1}^{(m)}}{\sum_{m=1}^M \left[ \sum_{i=0}^{T^{(m)}} \partial_{\theta_k} \log \pi_{\theta}(s_i^{(m)}, a_i^{(m)}) \right]^2} \quad (4.22)$$

**4.3.2 GPOMDP**

The Monte Carlo policy gradient estimate is typically characterized by a large variance, which may slow the method's convergence. To improve the estimate, it is sufficient to notice that future actions and past rewards are independent, unless the policy has been changed. Therefore, combining this observation with Eq. (4.18) yields

$$\mathbb{E}_{\pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta}(S_t, A_t) \mathcal{R}(S_s, A_s)] = 0 \quad \forall t > s$$

Plugging this equation in Eq. (4.19) leads to the well-known GPOMDP algorithm [6] for generating an estimate of the gradient of the objective function

$$\hat{g}_{\text{GPOMDP}} = \frac{1}{M} \sum_{m=1}^M \sum_{i=0}^{T^{(m)}} \nabla_{\theta} \log \pi_{\theta}(s_i^{(m)}, a_i^{(m)}) \left( \sum_{j=i}^{T^{(m)}} \gamma^j r_{j+1}^{(m)} - b \right) \quad (4.23)$$



By removing almost half of the cross-products, this estimate typically has a smaller variance than the trivial REINFORCE policy gradient and can be further reduced by computing an optimal baseline. In the following sections we will see how this algorithm can be easily derived from a more general result: the policy gradient theorem.

### 4.3.3 Risk-Sensitive Monte Carlo Policy Gradient

For an episodic environment, the REINFORCE algorithm can be easily extended to the risk-sensitive framework described above. Indeed, it is sufficient to adapt its derivation for the average return to the second-moment of return

$$U(\theta) = \mathbb{E}_{H \sim p_\theta} [G(H)^2]$$

Applying the likelihood ratio technique, we obtain

$$\nabla_\theta U(\theta) = \mathbb{E}_{H \sim p_\theta} [\nabla_\theta \log p_\theta(H) G(H)^2]$$

From which we can easily derive a Monte Carlo estimate. Since the problems we will consider in the next chapters are not episodic, we will not investigate this approach further and we refer the interested reader to the extensive work of Tamar et Al. [80], [78], [81], [79], [23].

### 4.3.4 Stochastic Policies

The likelihood ratio technique is a powerful tool and it allows to derive some powerful policy gradient algorithms. However, this approach only works if the policy is stochastic. In most cases this is not a big problem, since stochastic policies are needed anyway to ensure sufficient exploration of the state-action space. Moreover, stochastic policies might be beneficial in partially observable MDP to avoid state aliasing and in adversarial learning where the Pareto optimal strategies of different players are randomized, such as in the classical game “rock-paper-scissors”. We introduce now two standard examples of stochastic policies for discrete and continuous action spaces respectively.

#### 4.3.4.1 Boltzmann Exploration Policy

In discrete action spaces, the Boltzmann exploration policy, also known as softmax policy, is a common choice. In state  $s \in \mathbb{S}$ , this policy selects an action  $a \in \mathbb{A}$  with probability

$$\pi_\theta(s, a) = \frac{e^{\theta^T \Phi(s, a)}}{\sum_{b \in \mathbb{A}} e^{\theta^T \Phi(s, b)}} \quad (4.24)$$

where  $\Phi(s, a) \in \mathbb{R}^{D_\theta}$  is a given feature vector corresponding to state  $s$  and action  $a$ . the likelihood score for this policy is thus given by

$$\nabla_\theta \log \pi_\theta(s, a) = \Phi(s, a) - \sum_{b \in \mathbb{A}} \pi_\theta(s, b) \Phi(s, b) \quad (4.25)$$

#### 4.3.4.2 Gaussian Exploration Policy

4.3.4.1 In continuous action spaces, a Gaussian exploration policy is commonly used. According to this policy, in a state  $s$ , actions are sampled from a Gaussian distribution with a parametric state-dependent mean  $\mu_\psi(s) \in \mathbb{R}^{D_a}$ , with  $\psi \in \mathbb{R}^{D_\psi}$ , and a covariance matrix  $\Sigma \in \mathbb{R}^{D_a \times D_a}$ . The policy parameters consist of  $\theta = \psi, \Sigma$  and the likelihood score are thus given by

$$\nabla_\psi \log \pi_\theta(s, a) = \left( \frac{\partial \mu_\psi(s)}{\partial \psi} \right)^T \Sigma^{-1} (a - \mu_\psi(s)) \quad (4.26)$$

$$\nabla_\Sigma \log \pi_\theta(s, a) = \frac{1}{2} \left[ \Sigma^{-1} (a - \mu_\psi(s)) (a - \mu_\psi(s))^T \Sigma^{-1} - \Sigma^{-1} \right] \quad (4.27)$$

where  $\frac{\partial \mu_\psi(s)}{\partial \psi}$  denotes the Jacobian matrix of  $\mu_\psi$  with respect to  $\psi$ .

### 4.3.5 Policy Gradient with Parameter Exploration

In Monte Carlo Policy Gradient, trajectories are generated by sampling at each time step an action according to a stochastic policy  $\pi_\theta$  and the objective function gradient is estimated by differentiating the policy with respect to the parameters. However, sampling an action from the policy at each time step leads to a large variance in the sampled histories and therefore in the gradient estimate, which can in turn slow down the convergence of the learning process. To address this issue, in [70] the authors propose the *policy gradient with parameter-based exploration* (PGPE) method, in which the search in the policy space is replaced with a direct search in the model parameter space. We start by presenting the episodic case and we will later extend this approach to the infinite horizon setting.

#### 4.3.5.1 Episodic PGPE

Given an episodic MDP, PGPE considers a deterministic controller  $F : \mathbb{S} \times \Theta \rightarrow \mathbb{A}$  that, given a set of parameters  $\theta \in \Theta \subseteq \mathbb{R}^{D_\theta}$ , maps a state  $s \in \mathbb{S}$  to an action  $a = F(s; \theta) = F_\theta(s) \in \mathbb{A}$ . The policy parameters are drawn from a

probability distribution  $p_\xi$ , with hyperparameters  $\xi \in \Xi \subseteq \mathbb{R}^{D_\xi}$ . Combining these two hypotheses, the agent follows a stochastic policy  $\pi_\xi$  defined by

$$\forall B \in \mathcal{A}, \pi_\xi(s, B) = \pi(s, B; \xi) = \int_{\Theta} p_\xi(\theta) \mathbb{1}_{F_\theta(s) \in B} d\theta$$

The advantage of this approach is that the controller is deterministic and therefore the actions do not need to be sampled at each time step, with a consequent reduction of the gradient estimate variance. Indeed, It is sufficient to sample the parameters  $\theta$  once at the beginning of the episode and then generate an entire trajectory following the deterministic policy  $F_\theta$ . As an additional benefit, the parameter gradient is estimated by direct parameter perturbations, without having to backpropagate any derivatives, which allows to use non-differentiable controllers.

The hyperparameters  $\xi$  will be updated by following the gradient ascent direction of the gradient of the expected reward, which can be rewritten as

$$J(\xi) = \mathbb{E}_{\substack{\theta \sim p_\xi \\ H \sim p_\theta}} [G(H)] = \int_{\Theta} \int_{\mathbb{H}} p_\xi(\theta, h) G(h) dh d\theta \quad (4.28)$$

By remarking that  $h$  is conditionally independent from  $\xi$  given  $\theta$ , so that  $p_\xi(\theta, h) = p_\xi(\theta)p_\theta(h)$ , and applying the likelihood ratio technique, we obtain

$$\nabla_\xi J(\xi) = \mathbb{E}_{\substack{\theta \sim p_\xi \\ H \sim p_\theta}} [\nabla_\xi \log p_\xi(\theta) G(H)] \quad (4.29)$$

Again, we can subtract a constant baseline  $b \in \mathbb{R}$  from the total return

$$\nabla_\xi J(\xi) = \mathbb{E} [\nabla_\xi \log p_\xi(\theta) (G(H) - b)] \quad (4.30)$$

In an episodic environment, the gradient can be approximated via Monte Carlo by first drawing  $M$  samples  $\theta^{(m)} \sim p_\xi$  and then, for each combination of parameters, generating a trajectory  $h^{(m)} = \{(s_t^{(m)}, a_t^{(m)})\}_{t \geq 0}$  where actions are selected according to the deterministic controller  $F_{\theta^{(m)}}$ . This leads to the following estimate

$$\hat{g}_{\text{PGPE}} = \frac{1}{M} \sum_{m=1}^M \nabla_\xi \log p_\xi(\theta^{(m)}) [G(h^{(m)}) - b] \quad (4.31)$$

In order to further reduce the estimate variance, an optimal baseline can be computed similarly to the REINFORCE case [88]. The episodic PGPE algorithm obtained in this way is reported in Algorithm 4.3. In the following paragraphs, we discuss some possible choices for the policy parameters distribution  $p_\xi$ , which is the last component of the algorithm.

**Algorithm 4.3** Episodic PGPE algorithm**Input:** Initial hyperparameters  $\xi_0$ , learning rate  $\{\alpha_k\}$ **Output:** Approximation of the optimal policy  $F_{\xi^*} \approx \pi_*$ 


---

```

1: repeat
2:   for  $m = 1, \dots, M$  do
3:     Sample controller parameters  $\theta^{(m)} \sim p_{\xi^k}$ 
4:     Sample trajectory  $h^{(m)} = \{(s_t^{(m)}, a_t^{(m)})\}_{t \geq 0}$  under policy  $F_{\theta^{(m)}}$ 
5:   end for
6:   Approximate policy gradient  $\nabla_{\xi} J(\xi) \approx \hat{g}_{\text{PGPE}}$  according to Eq. 4.31
7:   Update hyperparameters using gradient ascent  $\xi^{k+1} = \xi^k + \alpha_k \hat{g}_{\text{PGPE}}$ 
8:    $k \leftarrow k + 1$ 
9: until converged

```

---

**Independent Gaussian Parameter Distribution** A simple approach is to assume that all the components of the parameter vector  $\theta$  are independent and normally distributed with mean  $\mu_i$  and variance  $\sigma_i^2$ , i.e.  $\theta_i \sim \mathcal{N}(\mu_i, \sigma_i^2)$ , the gradient with respect to the hyperparameters  $\xi = (\mu_1, \dots, \mu_{D_\theta}, \sigma_1, \dots, \sigma_{D_\theta})^T$  is given by

$$\begin{aligned} \frac{\partial \log p_{\xi}(\theta)}{\partial \mu_i} &= \frac{\theta_i - \mu_i}{\sigma_i^2} \\ \frac{\partial \log p_{\xi}(\theta)}{\partial \sigma_i} &= \frac{(\theta_i - \mu_i)^2 - \sigma_i^2}{\sigma_i^3} \end{aligned} \quad (4.32)$$

Using a constant learning rate  $\alpha_i = \alpha \sigma_i^2$ , the gradient updates takes the following form

$$\begin{aligned} \mu_i^{k+1} &= \mu_i^k + \alpha [G(h) - b] (\theta_i - \mu_i) \\ \sigma_i^{k+1} &= \sigma_i^k + \alpha [G(h) - b] \frac{(\theta_i - \mu_i)^2}{\sigma_i} \end{aligned} \quad (4.33)$$

where  $b$  can be computed as a moving average of the past returns. Intuitively, if  $G(h) > b$  we adjust  $\xi$  so as to increase the probability of  $\theta$  while if  $G(h) < b$  we do the opposite.

**Gaussian Parameter Distribution** A more elaborate approach is to assume a generic dependence among the controller parameters, i.e.  $\theta \sim \mathcal{N}(\mu, \Sigma)$ . However, if we directly used the covariance matrix  $\Sigma$  as an hyperparameter, it would be computationally difficult to enforce that it remains well-defined, i.e. symmetric and semidefinite positive, during the gradient ascent iterations. A simpler approach is to parameterize the distribution

using the Cholesky factor  $\Sigma$ , i.e. the matrix  $C$  such that  $\Sigma = C^T C$ . This choice has two advantages: first,  $C$  makes explicit the  $n(n+1)/2$  independent parameters determining the covariance matrix  $\Sigma$ ; in addition,  $C^T C$  is by construction a well-defined covariance matrix. Hence, the hyperparameters are  $\xi = \{\mu, C\}$ . The likelihood score for this distribution doesn't have a simple expression, but it becomes extremely easy in the Natural version of PGPE which will be discussed in the next sections.

**Symmetric Sampling and Gain Normalization** In some settings, comparing the gain with a baseline can be misleading. In their original work, the authors propose a symmetric sampling technique similar to antithetic variates that further improves the convergence of the method. More in detail, a more robust gradient estimate can be obtained by measuring the difference in reward between two symmetric samples on either side of the current mean. That is, we sample a random perturbation  $\epsilon \sim \mathcal{N}(0, \Sigma)$ , where  $\Sigma = \text{diag}(\sigma_1^2, \dots, \sigma_{D_\theta}^2)$ , and we define  $\theta^+ = \mu + \epsilon$  and  $\theta^- = \mu - \epsilon$ . Denoting by  $G^+$  (resp.  $G^-$ ) the gains obtained on the trajectory associated to  $\theta^+$  (resp.  $\theta^-$ ), the objective function gradient can be approximated with

$$\begin{aligned}\nabla_{\mu_i} J(\xi) &\approx \frac{\epsilon_i (G^+ - G^-)}{2\sigma_i^2} \\ \nabla_{\sigma_i} J(\xi) &\approx \frac{\epsilon_i^2 - \sigma_i^2}{\sigma_i^3} \left( \frac{G^+ + G^-}{2} - b \right)\end{aligned}\tag{4.34}$$

Hence, by choosing  $\alpha_i^k = 2\alpha\sigma_i^2$ , we have the following update rules

$$\begin{aligned}\mu_i^{k+1} &= \mu_i^k + \alpha\epsilon_i (G^+ - G^-) \\ \sigma_i^{k+1} &= \sigma_i^k + \alpha \frac{\epsilon_i^2 - \sigma_i^2}{\sigma_i} (G^+ + G^- - 2b)\end{aligned}\tag{4.35}$$

In addition, the authors propose to normalize the gains in order to make the updates independent of the scale of the rewards. For instance, we could modify the hyperparameters updates as follows

$$\begin{aligned}\mu_i^{k+1} &= \mu_i^k + \alpha\epsilon_i \frac{(G^+ - G^-)}{2m - G^+ - G^-} \\ \sigma_i^{k+1} &= \sigma_i^k + \alpha \frac{\epsilon_i^2 - \sigma_i^2}{\sigma_i} \frac{(G^+ + G^- - 2b)}{m - b}\end{aligned}\tag{4.36}$$

where  $m$  might be the maximum gain the agent can receive, if known, or alternatively the maximum gain achieved so far. Symmetric sampling and gain normalization can drastically improve the gradient estimate quality and consequently the convergence time.

### 4.3.5.2 Infinite Horizon PGPE

While in the episodic PGPE the parameters  $\theta$  are sampled only at the beginning of each episode, in *Infinite Horizon PGPE* (IHPGPE) [69] the parameters and learning are carried out simultaneously, while interacting with the environment. Let  $0 < \varepsilon < 1$  the probability of updating the policy parameters, the parameters  $\theta_t$  can be sampled consecutively as follows

$$p_{\xi}(\theta_{i,t+1}) = \varepsilon \mathcal{N}(\mu_{i,t}, \sigma_{i,t}^2) + (1 - \varepsilon) \delta_{\theta_{i,t}} \quad (4.37)$$

In practice,  $\varepsilon$  should be chosen so that the expected frequency of changing a single parameter is coherent with the typical episode length in the episodic framework. Alternatively, one could sample all the parameters at a certain time step simultaneously

$$p_{\xi}(\theta_{t+1}) = \varepsilon \mathcal{N}(\mu_t, \Sigma_t) + (1 - \varepsilon) \delta_{\theta_t} \quad (4.38)$$

This is equivalent to splitting the state-action space into artificial episodes. However, updating parameters asynchronously changes the policy only slightly thus introducing less noise in the process. Again, parameters can be updated at every time step by gradient ascent

$$\begin{aligned} \mu_{i,t+1} &= \mu_{i,t} + \alpha [G_t(h) - b] (\theta_{i,t} - \mu_{i,t}) \\ \sigma_{i,t+1} &= \sigma_{i,t} + \alpha [G_t(h) - b] \frac{(\theta_{i,t} - \mu_{i,t})^2}{\sigma_{i,t}} \end{aligned} \quad (4.39)$$

Similarly to the episodic case, we can improve the gradient estimate by symmetric sampling and gain normalization.

## 4.4 Risk-Neutral Policy Gradient Theorem

In the previous sections, we saw that the core of a policy gradient algorithm consist in the approximation of the objective function gradient. In this section we present the *policy gradient theorem* [76], which shows that the gradient can be rewritten in a form suitable for estimation from experience aided by an approximate action-value or advantage function. First, we will prove the theorem in the risk-neutral setting for which it was originally proposed. In particular, we will see how the GPOMDP algorithm can be easily derived by this result via a Monte Carlo approximation. Moreover, we will derive a class of actor-critic algorithms [43] that, in addition to a parametric approximation of the policy, also exploit an approximation of the action-value function or of an advantage function to reduce the variance of the gradient

estimate. In particular, we review the powerful idea of compatible function approximation [76] which assures the convergence to a local optimum of the objective function. Finally, we discuss the natural policy gradient idea [40], which forms the basis of many state-of-the-art algorithms. The extension to the risk-sensitive setting will be done in the next section.

#### 4.4.1 Theorem Statement and Proof

**Theorem 4.4.1** (Policy Gradient). *Let  $\pi_\theta$  be a differentiable policy. The policy gradient for the average reward formulation is given by*

$$\nabla_\theta \rho(\theta) = \mathbb{E}_{\substack{S \sim d^\theta \\ A \sim \pi_\theta}} [\nabla_\theta \log \pi_\theta(S, A) Q_{\pi_\theta}(S, A)] \quad (4.40)$$

where  $d^\theta$  is the stationary distribution of the Markov chain induced by  $\pi_\theta$ . The policy gradient for the start value formulation is given by

$$\nabla_\theta J_{\text{start}}(\theta) = \mathbb{E}_{\substack{S \sim d_\gamma^\theta(s_0, \cdot) \\ A \sim \pi_\theta}} [\nabla_\theta \log \pi_\theta(S, A) Q_{\pi_\theta}(S, A)] \quad (4.41)$$

where  $d_\gamma^\theta(s_0, \cdot)$  is the  $\gamma$ -discounted visiting distribution over states starting from the initial state  $s_0$  and following policy  $\pi_\theta$

$$d_\gamma^\theta(s, x) = \sum_{k=0}^{\infty} \gamma^k \mathcal{P}_\theta^{(k)}(s, x) \quad (4.42)$$

*Proof.* We first prove the result for the average-reward formulation and then for the start state formulation. From the basic relation between state-value function and action-value function, we have

$$\begin{aligned} \nabla_\theta V_\theta(s) &= \nabla_\theta \int_{\mathbb{A}} \pi_\theta(s, a) Q_\theta(s, a) da \\ &= \int_{\mathbb{A}} [\nabla_\theta \pi_\theta(s, a) Q_\theta(s, a) + \pi_\theta(s, a) \nabla_\theta Q_\theta(s, a)] da \end{aligned}$$

Using the Bellman expectation equation for  $Q_\theta$

$$\begin{aligned} \nabla_\theta Q_\theta(s, a) &= \nabla_\theta \left[ \mathcal{R}(s, a) - \rho_\theta + \int_{\mathbb{S}} \mathcal{P}(s, a, s') V_\theta(s') ds' \right] \\ &= -\nabla_\theta \rho_\theta + \int_{\mathbb{S}} \mathcal{P}(s, a, s') \nabla_\theta V_\theta(s') ds' \end{aligned}$$

Hence, plugging in the first equation

$$\nabla_\theta V_\theta(s) = \int_{\mathbb{A}} \nabla_\theta \pi_\theta(s, a) Q_\theta(s, a) da - \nabla_\theta \rho_\theta + \int_{\mathbb{A}} \pi_\theta(s, a) \int_{\mathbb{S}} \mathcal{P}(s, a, s') \nabla_\theta V_\theta(s') ds'$$

Integrating both sides with respect to the stationary distribution  $d^\theta$  and noting that, because of stationarity,

$$\int_{\mathbb{S}} d^\theta(s) \int_{\mathbb{A}} \pi(s, a) \int_{\mathbb{S}} \mathcal{P}(s, a, s') \nabla_\theta(s') ds' dad s = \int_{\mathbb{S}} d^\theta(s) \nabla_\theta V_\theta(s) ds$$

we obtain the result

$$\begin{aligned} \nabla_\theta \rho_\theta &= \int_{\mathbb{S}} d^\theta(s) \int_{\mathbb{A}} \nabla_\theta \pi_\theta(s, a) Q_\theta(s, a) dad s \\ &= \int_{\mathbb{S}} d^\theta(s) \int_{\mathbb{A}} \pi_\theta(s, a) \nabla_\theta \log \pi_\theta(s, a) Q_\theta(s, a) dad s \\ &= \mathbb{E}_{\substack{S \sim d^\theta \\ A \sim \pi_\theta}} [\nabla_\theta \log \pi_\theta(S, A) Q_{\pi_\theta}(S, A)] \end{aligned}$$

Let us now prove the theorem for the start state formulation. The first step is exactly the same. Hence, using the Bellman expectation equation for  $Q_\theta$

$$\begin{aligned} \nabla_\theta Q_\theta(s, a) &= \nabla_\theta \left[ \mathcal{R}(s, a) + \gamma \int_{\mathbb{S}} \mathcal{P}(s, a, s') V_\theta(s') ds' \right] \\ &= \gamma \int_{\mathbb{S}} \mathcal{P}(s, a, s') \nabla_\theta V_\theta(s') ds' \end{aligned}$$

we obtain

$$\begin{aligned} \nabla_\theta V_\theta(s) &= \int_{\mathbb{A}} \left[ \nabla_\theta \pi_\theta(s, a) Q_\theta(s, a) + \gamma \int_{\mathbb{S}} \mathcal{P}(s, a, s') \nabla_\theta V_\theta(s') ds' \right] da \\ &= \int_{\mathbb{S}} \sum_{k=0}^{\infty} \gamma^k \mathcal{P}_\theta^{(k)}(s, x) \int_{\mathbb{A}} \nabla_\theta \pi_\theta(x, a) Q_\theta(x, a) dad x \end{aligned}$$

after unrolling  $\nabla_\theta V_\theta$  infinite times and denoting by  $\mathcal{P}_\theta^{(k)}(s, x)$  the probability of going from state  $s$  to state  $x$  in  $k$  steps under policy  $\pi_\theta$ . Defining the  $\gamma$ -discounted visiting distribution of state  $x$  starting from state  $s$  as

$$d_\gamma^\theta(s, x) = \sum_{k=0}^{\infty} \gamma^k \mathcal{P}_\theta^{(k)}(s, x)$$

we have the result

$$\begin{aligned} \nabla_\theta V_\theta(s) &= \int_{\mathbb{S}} d_\gamma^\theta(s, x) \int_{\mathbb{A}} \nabla_\theta \pi_\theta(x, a) Q_\theta(x, a) dad x \\ &= \mathbb{E}_{\substack{S \sim d_\gamma^\theta(s_0, \cdot) \\ A \sim \pi_\theta}} [\nabla_\theta \log \pi_\theta(S, A) Q_{\pi_\theta}(S, A)] \end{aligned}$$

□



The action-value function is typically unknown and needs to be approximated. As for the REINFORCE algorithm, we can subtract a state-dependent baseline from the action-value function without changing the value of the expectation. Indeed

$$\begin{aligned}\mathbb{E}_{\substack{S \sim d^\theta \\ A \sim \pi_\theta}} [\nabla_\theta \log \pi_\theta(S, A) B_\theta(S)] &= \int_{\mathbb{S}} d^\theta(s) \int_{\mathbb{A}} \pi_\theta(s, a) \nabla_\theta \log \pi_\theta(s, a) B_\theta(s) da ds \\ &= \int_{\mathbb{S}} d^\theta(s) B_\theta(s) \int_{\mathbb{A}} \nabla_\theta \pi_\theta(s, a) da ds \\ &= \int_{\mathbb{S}} d^\theta(s) B_\theta(s) \nabla_\theta \underbrace{\int_{\mathbb{A}} \pi_\theta(s, a) da}_{=1} ds = 0\end{aligned}$$

Hence, we can rewrite the policy gradient for the average reward formulation as

$$\nabla_\theta \rho(\theta) = \mathbb{E}_{\substack{S \sim d^\theta \\ A \sim \pi_\theta}} [\nabla_\theta \log \pi_\theta(S, A) (Q_{\pi_\theta}(S, A) - B_\theta(S))] \quad (4.43)$$

and similarly for the start state formulation. This result can be used as the starting point to derive several policy gradient methods that use different approximation of the action-value function.

#### 4.4.2 GPOMDP

For an episodic MDP, the action-value function can be estimated with the total return obtained on a sample trajectory

$$Q_\theta(s_0, a_0) \approx \sum_{t=0}^{T^{(m)}} \gamma^t r_{t+1}^{(m)}$$

Combining this remark with a Monte Carlo approximation of Eq. (4.43), we obtain the GPOMDP gradient estimate (4.23). Therefore, the GPOMDP algorithm is a simple application of the policy gradient algorithm, which is much more general than the likelihood ratio technique used to derive the algorithm in the previous sections. In the next subsections we present other ways to exploit the policy gradient theorem to design efficient learning algorithms.

#### 4.4.3 Actor-Critic Policy Gradient

A baseline should ideally measure the typical return obtained by an agent in a certain state when following a certain policy. Therefore, it becomes natural

---

**Algorithm 4.4** Generic structure for an online actor-critic algorithm.

---

**Input:** Initial actor  $\theta^0$ , initial critic  $\psi^0$ , learning rate  $\{\alpha_k\}$

**Output:** Approximation of the optimal policy  $\pi_{\theta^*} \approx \pi_*$

- 1: **repeat**
  - 2:   Observe tuple  $\langle s_k, a_k, r_{k+1}, s_{k+1} \rangle$  sampled from the MDP.
  - 3:   Update critic parameters  $\psi^{k+1}$  using a value-based method.
  - 4:   Estimate policy gradient as  $\hat{g}_{AC}^k = \nabla_{\theta} \log \pi_{\theta}(s_k, a_k) \hat{A}_{\psi^{k+1}}(s_k, a_k)$
  - 5:   Update actor by gradient ascent  $\theta^{k+1} = \theta^k + \alpha_k \hat{g}_{AC}^k$ .
  - 6:    $k \leftarrow k + 1$
  - 7: **until** converged
- 

to use the state-value function as a benchmark for the action-value function. Eq. (4.43) thus becomes

$$\nabla_{\theta} \rho(\theta) = \mathbb{E}_{\substack{S \sim d^{\theta} \\ A \sim \pi_{\theta}}} [\nabla_{\theta} \log \pi_{\theta}(S, A) A_{\theta}(S, A)] \quad (4.44)$$

where we introduced the advantage function

$$A_{\theta}(s, a) = Q_{\theta}(s, a) - V_{\theta}(s) \quad (4.45)$$

which measures how good is to take action  $a$  in state  $s$  compared to simply following the policy  $\pi_{\theta}$ . The advantage function is unknown and should be estimated from samples of the MDP. *Actor-Critic* algorithms consist in all those methods that employ an approximation of the advantage function or of the value function, also known as critic, to estimate the policy gradient. These methods thus maintain two sets of parameters: a *critic* that updates the action-value function parameters  $\psi$  and an *actor* that updates the policy parameters  $\theta$  in the direction suggested by the critic. The general structure for an online actor-critic algorithm is reported in Algorithm 4.4.

There are two possible approaches to estimate the advantage function: the first one is to estimate both the state-value function  $\hat{V}_{\psi} \approx V_{\theta}$  and the action-value function  $\hat{Q}_{\xi} \approx Q_{\theta}$  and derive an estimate for the advantage function  $\hat{A}_{\psi, \xi} = \hat{Q}_{\xi} - \hat{V}_{\psi} \approx A_{\theta}$ . The second approach directly approximates the advantage function  $\hat{A}_{\psi} \approx A_{\theta}$  and is usually preferred, since it allows us to maintain only one critic. Actor-critic algorithms typically use a temporal-difference algorithm to update an approximation of the value function  $\hat{V}_{\psi} \approx V_{\theta}$ , from which we can derive an approximation for the advantage function. Assume that the true state-value function  $V_{\theta}$  is given. Then the TD error

$$\delta_{\theta} = R - \rho_{\theta} + V_{\theta}(S') - V_{\theta}(S) \quad (4.46)$$

is an unbiased estimate of the advantage function. Indeed

$$\begin{aligned}\mathbb{E} [\delta_\theta | S = s, A = a] &= \mathbb{E} [R - \rho_\theta + V_\theta(S') | S = s, A = a] - V_\theta(s) \\ &= Q_\theta(s, a) - V_\theta(s) \\ &= A_\theta(s, a)\end{aligned}$$

By a simple conditioning argument, the policy gradient can be rewritten as

$$\nabla_\theta \rho(\theta) = \mathbb{E}_{\substack{S \sim d^\theta \\ A \sim \pi_\theta}} [\nabla_\theta \log \pi_\theta(S, A) \delta_\theta] \quad (4.47)$$

In practice, we can use an approximate TD error

$$\hat{\delta}_\psi^k = r_{k+1} - \hat{\rho}_k + \hat{V}_\psi(s_{k+1}) - \hat{V}_\psi(s_k) \quad (4.48)$$

where  $\hat{\rho}_k$  is an estimate of the average reward. Hence, we can easily obtain an approximation of the policy gradient by

$$\hat{g}_{\text{TD}(0)}^k = \nabla_\theta \log \pi_\theta(s_k, a_k) \hat{\delta}_\psi^k \quad (4.49)$$

which can be used to update the critic parameter in the gradient ascent direction. On the other hand, the critic parameters can be updated using a TD(0) temporal difference scheme

$$\psi^{k+1} = \psi^k + \alpha_k \hat{\delta}_\psi^k \nabla_\psi \hat{V}_\psi(s_k) \quad (4.50)$$

The discussion can be easily extended to more complex value-based methods for estimating the critic, such as the backward-view TD( $\lambda$ ). This method employs eligibility traces to assign credit for the rewards obtained by the agent to all previous states and actions. More formally, the policy parameters update rule becomes

$$\theta^{k+1} = \theta^k + \alpha_k \hat{\delta}_\psi^k e^k \quad (4.51)$$

where the eligibility trace  $e_k$  is defined by the following recursive equation

$$e^{k+1} = \lambda e^k + \nabla_\theta \log \pi_\theta(s_k, a_k) \quad (4.52)$$

$0 \leq \lambda \leq 1$  is a parameter that manages the amount of bootstrap: for  $\lambda = 0$  the method is the standard temporal difference scheme while for  $\lambda = 1$  the method is equivalent to Monte Carlo. The complete scheme is reported in Algorithm 4.5.

---

**Algorithm 4.5** TD( $\lambda$ ) policy gradient algorithm.

---

**Input:** Initial actor  $\theta^0$ , initial critic  $\psi^0$ , learning rate  $\{\alpha_k\}$

**Output:** Approximation of the optimal policy  $\pi_{\theta^*} \approx \pi_*$

---

- 1: **repeat**
  - 2:   Observe tuple  $\langle s_k, a_k, r_{k+1}, s_{k+1} \rangle$  sampled from the MDP.
  - 3:   Update critic parameters  $\psi^{k+1} = \psi^k + \alpha_k \widehat{\delta}_\psi^k \nabla_\psi \widehat{V}_\psi(s_k)$ .
  - 4:   Update actor parameters  $\theta^{k+1} = \theta^k + \alpha_k \alpha_k \widehat{\delta}_\psi^k e^k$ .
  - 5:   Update eligibility trace  $e^{k+1} = \lambda e^k + \nabla_\theta \log \pi_\theta(s_k, a_k)$
  - 6:    $k \leftarrow k + 1$
  - 7: **until** converged
- 

#### 4.4.4 Compatible Function Approximation

In the previous sections we saw that a critic may reduce the variance of the policy gradient estimate. However, this is achieved at the cost of introducing a bias in the approximation which, in certain cases, may endanger the convergence of the method to a good solution. Luckily, the *compatible function approximation* [76] avoids introducing any bias by a careful choice of the critic. This technique consists in estimating the advantage function with a linear regression on a family of suitable basis functions consisting of the gradient of the likelihood score, i.e.

$$\widehat{A}_\theta(s, a) = \psi^T \nabla_\theta \log \pi_\theta(s, a) \quad (4.53)$$

The power of these basis functions becomes apparent in the following theorem

**Theorem 4.4.2** (Compatible Function Approximation). *If*

- i) the advantage function approximator is compatible to the policy, i.e.*

$$\nabla_\psi A_\psi(s, a) = \nabla_\theta \log \pi_\theta(s, a) \quad (4.54)$$

- ii) the advantage function parameters  $\psi$  minimize the mean square error*

$$\theta = \arg \min_{\psi} \mathbb{E}_{\pi_\theta} \left[ \left( A_{\pi_\theta}(S, A) - \widehat{A}_\psi(S, A) \right)^2 \right] \quad (4.55)$$

*Then the policy gradient is exact*

$$\nabla_\theta J(\theta) = \mathbb{E}_{\pi_\theta} \left[ \nabla_\theta \log \pi_\theta(S, A) \widehat{A}_\psi(S, A) \right] \quad (4.56)$$

Hence, the true advantage function in the policy gradient formula can be replaced by the compatible function approximation without changing the

policy gradient value. A linear regression may seem a too simple representation of the advantage function. However, one must bear in mind that this approximation is only used to reduce the variance of the policy gradient estimate and should not be expected to provide an accurate representation of the true-value function. This technique has proved useful in many application.

### 4.4.5 Natural Policy Gradient

Despite all the advances in the variance reduction techniques, these methods still tend to perform surprisingly poorly. Even when applied to simple examples with rather few states, where the gradient can be determined very accurately, they turn out to be quite inefficient. Typically, one of the reasons of this behavior is the presence of large plateaus in the expected return landscape where the gradients are small and often do not point directly towards the optimal solution. In this context, the steepest ascent with respect to the Fisher information metric, called the *natural policy gradient*, turns out to be significantly more efficient than normal gradients for such plateaus. This technique was first proposed in the reinforcement learning setting in [40] and later applied to actor-critic algorithms in [62]. The following properties of natural policy gradient make it one of the more reliable policy-based methods,

- i) Convergence to a local minimum is guaranteed.
- ii) By choosing a more direct path to the optimal solution in parameter space, the natural gradient typically has faster convergence and avoids premature convergence of “vanilla gradients”.
- iii) The natural policy gradient can be shown to be covariant, i.e. independent of the coordinate frame chosen for expressing the policy parameters.
- iv) As the natural gradient analytically averages out the influence of the stochastic policy (including the baseline of the function approximator), it requires fewer data points for a good gradient estimate than “vanilla gradients”.

#### 4.4.5.1 Formalism of Natural Policy Gradients

Policy gradient methods improve the policy  $\pi_\theta$  by iteratively applying “small” changes  $\Delta\theta$  to the policy parameters  $\theta$ . However, the meaning of “small” is ambiguous. For instance, when working with an Euclidean metric, the size of this update  $\|\Delta\theta\| = \sqrt{\Delta\theta^T \Delta\theta}$  and therefore the update size depends on the parameterization of the policy, which often results in unnaturally slow learning even if higher-order gradient methods were employed. This prob-

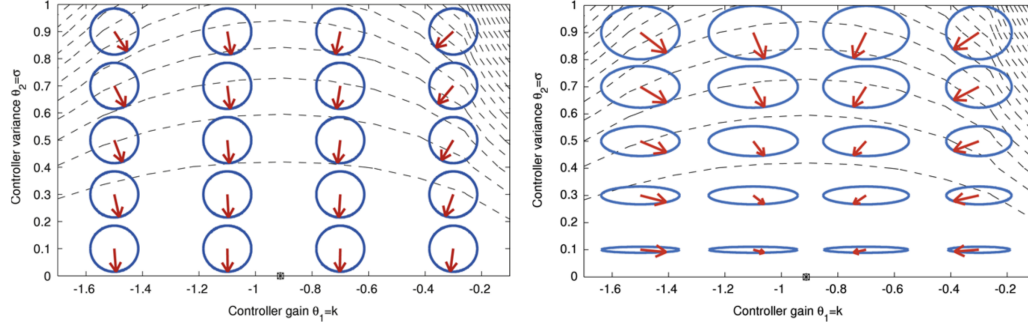


Figure 4.1: “Vanilla” policy gradient (left) vs. natural policy gradient (right). The main difference is how the two approaches punish the change in parameters. This distance is indicated by the blue ellipses in the contour plot while the dashed lines show the expected return. Obtaining a gradient then corresponds to finding a vector pointing from the center of the ellipses to the location with maximum expected return on the ellipse. A vanilla policy gradient (a) considers a change in all parameters as equally distant, thus, it is a search for a maximum on a circle while the natural gradient (b) uses scales determined by the Fisher information which results in a reduction in exploration. The slower reduction in exploration results into a faster convergence to the optimal policy.

lem poses the question whether we can achieve a covariant gradient descent, i.e. a gradient descent with respect to an invariant measure of the closeness between the current policy and the updated policy based upon the distribution of the paths generated by each of these. Standard measures of distance between probability distributions are the Kullback-Leibler divergence  $d_{\text{KL}}(p_{\theta}(h) \parallel p_{\theta+\Delta\theta}(h))$  and the Hellinger distance. These two distances can be approximation in first instance by a second-order Taylor expansion

$$d_{\text{KL}}(p_{\theta}(h) \parallel p_{\theta+\Delta\theta}(h)) \approx \frac{1}{2} \Delta\theta^T F_{\theta} \Delta\theta$$

where  $F_{\theta}$  is the Fischer information matrix

$$\begin{aligned} F_{\theta} &= \int_{\mathbb{H}} p_{\theta}(h) \nabla_{\theta} \log p_{\theta}(h) \nabla_{\theta} \log p_{\theta}(h)^T dh \\ &= \mathbb{E}_{H \sim p_{\theta}} [\nabla_{\theta} \log p_{\theta}(H) \nabla_{\theta} \log p_{\theta}(H)^T] \end{aligned} \quad (4.57)$$

The goal is to find the optimal update  $\Delta\theta$  for the policy parameters so as to maximize the objective function, under the constraint that the new policy must be in a radius  $\epsilon$  from the previous policy with respect to the Kullbach-

Leibler divergence, i.e.

$$\begin{cases} \max_{\Delta\theta} J(\theta + \Delta\theta) \approx J(\theta) + \Delta\theta^T \nabla_{\theta} J(\theta) \\ \text{s.t. } d_{\text{KL}}(p_{\theta}(h) \parallel p_{\theta+\Delta\theta}(h)) \approx \frac{1}{2} \Delta\theta^T F_{\theta} \Delta\theta < \epsilon \end{cases}$$

The optimal solution is given by

$$\Delta\theta = \alpha_n F_{\theta}^{-1} \nabla_{\theta} J(\theta)$$

with

$$\alpha_n = \sqrt{\frac{\epsilon}{\nabla_{\theta} J(\theta)^T F_{\theta}^{-1} \nabla_{\theta} J(\theta)}}$$

The direction  $\tilde{\nabla}_{\theta} J(\theta) = \Delta\theta/\alpha_n$  is called the natural gradient and learning algorithms that use this gradient instead of the standard one are called natural policy gradient algorithms. The strongest theoretical advantage of this approach is that its performance no longer depends on the parameterization of the policy and it is therefore safe to use for arbitrary policies. In practice, the learning process converges significantly faster in most practical cases and requires less manual parameter tuning of the learning algorithm. The only remaining point to discuss is how to compute the inverse of the Fischer information matrix. This is not an easy task and the matrix will usually need to be estimated from sample trajectories. However, we will see that this matrix can be computed analytically after reformulating the policy gradient theorem for the parameter-based search methods such as PGPE. The formal derivation will be presented in the last section of this chapter and this result will lead to the *Natural PGPE* (NPGPE) algorithm [51].

## 4.5 Risk-Sensitive Policy Gradient Theorem

In this section the policy gradient theorem is extended to the risk-sensitive framework. In the average reward formulation of the control problem, the result and its derivation are analogous to those presented in the last section for the risk-neutral framework. This will allow to derive in a trivial way the risk-sensitive versions of all the learning algorithms seen above. The risk-sensitive policy gradient theorem for the average-reward formulation was first derived in [63] and the presentation of this section closely follows the original article. Obtaining an equivalent theorem for the discounted reward formulation is more challenging. In the article cited above, the authors prove the theorem under a very strong assumption on the dependence of rewards obtained at different time steps which is not verified in many applications, among which is the asset allocation problem that we will consider in the

next chapter. We will discuss the problems arising in the discounted setting and why it is not easy to derive a policy gradient theorem in this case. In the original article the authors mostly considered the mean-variance criterion since all the algorithms they propose are easily adapted to the Sharpe ratio criterion. Here we take the opposite direction and present the algorithms for the Sharpe ratio, referring to their article for the mean-variance counterparts.

Let us consider a family of parametrized policies  $\pi_\theta$ , with  $\theta \in \Theta \subseteq \mathbb{R}^{D_\theta}$ . The optimization problem then becomes

$$\max_{\lambda} \min_{\theta} L(\theta, \lambda) = -\rho(\theta) + \lambda(\Lambda(\theta) - \alpha) \quad (4.58)$$

Using a policy gradient approach, the policy parameters are updated following the gradient ascent direction

$$\nabla_{\theta} L(\theta, \lambda) = -\nabla_{\theta} \rho(\theta) + \lambda \nabla_{\theta} \Lambda(\theta) \quad (4.59)$$

while the Lagrange multiplier is updated following the gradient descent direction

$$\nabla_{\lambda} L(\theta, \lambda) = \Lambda(\theta) - \alpha \quad (4.60)$$

Since

$$\nabla_{\theta} \Lambda(\theta) = \nabla_{\theta} \eta(\theta) - 2\rho(\theta) \nabla_{\theta} \rho(\theta) \quad (4.61)$$

it is enough to compute  $\nabla_{\theta} \eta(\theta)$  and  $\nabla_{\theta} \rho(\theta)$ . These quantities are provided by the policy gradient theorem

**Theorem 4.5.1** (Policy Gradient).

$$\begin{aligned} \nabla_{\theta} \rho(\theta) &= \mathbb{E}_{\substack{S \sim d_{\pi} \\ A \sim \pi}} [\nabla_{\theta} \log \pi_{\theta}(S, A) Q_{\theta}(S, A)] \\ &= \int_{\mathbb{S}} d_{\pi}(s) \int_{\mathbb{A}} \pi(s, a) \nabla_{\theta} \log \pi_{\theta}(s, a) Q_{\theta}(s, a) da ds \end{aligned} \quad (4.62)$$

$$\begin{aligned} \nabla_{\theta} \eta(\theta) &= \mathbb{E}_{\substack{S \sim d_{\pi} \\ A \sim \pi}} [\nabla_{\theta} \log \pi_{\theta}(S, A) W_{\theta}(S, A)] \\ &= \int_{\mathbb{S}} d_{\pi}(s) \int_{\mathbb{A}} \pi(s, a) \nabla_{\theta} \log \pi_{\theta}(s, a) W_{\theta}(s, a) da ds \end{aligned} \quad (4.63)$$

*Proof.* Eq. (4.62) is the standard policy gradient and its proof can be found in the literature, e.g. [76]. Eq. (4.63) can be shown in a similar fashion. □

As in the standard risk-neutral case, a state-dependent baseline can be introduced in both gradients without changing the result. In particular, by using the average adjusted value functions as baseline, we can replace



the average adjusted action-value functions with the following advantage functions

$$A_\theta(s, a) = Q_\theta(s, a) - V_\theta(s) \quad (4.64)$$

$$B_\theta(s, a) = W_\theta(s, a) - U_\theta(s) \quad (4.65)$$

Intuitively, the advantage functions measure how good is to take action  $a$  in state  $s$  compared to simply following the policy  $\pi_\theta$ . The use of a baseline also serves the purpose of reducing the variance of the gradient estimates. Thus, the gradients can be written as

$$\nabla_\theta \rho(\theta) = \mathbb{E}_{\substack{S \sim d_\pi \\ A \sim \pi}} [\nabla_\theta \log \pi_\theta(S, A) A_\theta(S, A)] \quad (4.66)$$

$$\nabla_\theta \eta(\theta) = \mathbb{E}_{\substack{S \sim d_\pi \\ A \sim \pi}} [\nabla_\theta \log \pi_\theta(S, A) B_\theta(S, A)] \quad (4.67)$$

Results of this type form the basis of many actor-critic algorithms, which employ an approximation of the advantage function to obtain a more accurate estimate of the objective function.

### 4.5.1 Risk-Sensitive Actor-Critic Algorithm

In [63], starting from (4.66) and (4.67), the authors derive a risk-sensitive actor-critic algorithm for the average reward setting. In the algorithm proposed, the advantage functions are approximated using a temporal difference (TD) scheme in which a critic maintains a linear approximation of the value functions. More in detail, let  $\delta_n^A$  and  $\delta_n^B$  be the TD errors for residual value and square value functions

$$\begin{aligned} \delta_t^A &= R_{t+1} - \hat{\rho}_{t+1} + \hat{V}(S_{t+1}) - \hat{V}(S_t) \\ \delta_t^B &= R_{t+1}^2 - \hat{\eta}_{t+1} + \hat{U}(S_{t+1}) - \hat{U}(S_t) \end{aligned} \quad (4.68)$$

where  $\hat{V}$ ,  $\hat{U}$ ,  $\hat{\rho}$  and  $\hat{\eta}$  are unbiased estimate of  $V_\theta$ ,  $U_\theta$ ,  $\rho(\theta)$  and  $\eta(\theta)$  respectively. It is easy to show that  $\delta_t^A$  and  $\delta_t^B$  are unbiased estimates of the advantage functions.

**Proposition 4.5.1** (Temporal Difference Errors).

$$\begin{aligned} \mathbb{E}_\theta [\delta_t^A | S_t = s, A_t = a] &= A_\theta(s, a) \\ \mathbb{E}_\theta [\delta_t^B | S_t = s, A_t = a] &= B_\theta(s, a) \end{aligned} \quad (4.69)$$

*Proof.* TODO

□

Denoting by  $\psi_t = \psi(S_t, A_t) = \nabla_{\theta} \log \pi(S_t, A_t)$  the compatible feature [76], we can easily obtain an unbiased estimate of the gradients

$$\begin{aligned}\nabla_{\theta} \rho(\theta) &\approx \psi_t \delta_t^A \\ \nabla_{\theta} \eta(\theta) &\approx \psi_t \delta_t^B\end{aligned}\tag{4.70}$$

The value functions are linearly approximated using some features vectors  $\Phi_V : \mathbb{S} \rightarrow \mathbb{R}^{D_V}$  and  $\Phi_U : \mathbb{S} \rightarrow \mathbb{R}^{D_U}$  as follows

$$\begin{aligned}\hat{V}(s) &= v^T \Phi_V(s) \\ \hat{U}(s) &= u^T \Phi_U(s)\end{aligned}\tag{4.71}$$

Given all these ingredients, the authors propose a three time-scale stochastic approximation algorithm whose pseudo-code is illustrated in Algorithm

## 4.6 Parameter-Based Policy Gradient Theorem

### 4.6.1 Risk-Neutral Setting

### 4.6.2 Risk-Sensitive Setting

### 4.6.3 Natural Policy Gradient

#### 4.6.3.1 NPGPE

#### 4.6.3.2 Risk-Sensitive NPGPE

# Chapter 5

## Financial Applications of Reinforcement Learning

### 5.1 Bibliographical Survey

### 5.2 Asset Allocation

The asset allocation problem consists of determining how to dynamically invest the available capital in a portfolio of different assets in order to maximize the expected total return or another relevant performance measure. Let us consider a financial market consisting of  $I + 1$  different stocks that are traded only at discrete times  $t \in \{0, 1, 2, \dots\}$  and denote by  $Z_t = (Z_t^0, Z_t^1, \dots, Z_t^I)^T$  their prices at time  $t$ . Typically,  $Z_t^0$  refers to a riskless asset whose dynamic is given by  $Z_t^0 = (1 + X)^t$  where  $X$  is the deterministic risk-free interest rate. The investment process works as follows: at time  $t$ , the investor observes the state of the market  $S_t$ , consisting for example of the past asset prices and other relevant economic variables, and subsequently chooses how to rebalance his portfolio, by specifying the units of each stock  $n_t = (n_t^0, n_t^1, \dots, n_t^I)^T$  to be held between  $t$  and  $t + 1$ . In doing so, he needs to take into account the transaction costs that he has to pay to the broker to change his position. At time  $t + 1$ , the investor realizes a profit or a loss from his investment due to the stochastic variation of the stock values. The investor's goal is to maximize a given performance measure.

#### 5.2.1 Reward Function

Let  $W_t$  denote the wealth of the investor at time  $t$ . The profit realized between  $t$  and  $t + 1$  is simply given by the difference between the trading

results and the transaction costs paid to the broker. More formally

$$\Delta W_{t+1} = W_{t+1} - W_t = \text{PNL}_{t+1} - \text{TC}_t$$

where  $\text{PNL}_{t+1}$  denotes the profit due to the variation of the portfolio asset prices between  $t$  and  $t + 1$

$$\text{PNL}_{t+1} = n_t \cdot \Delta Z_{t+1} = \sum_{i=0}^I n_t^i (Z_{t+1}^i - Z_t^i)$$

and  $\text{TC}_t$  denotes the fees paid to the broker to change the portfolio allocation and on the short positions

$$\text{TC}_t = \sum_{i=0}^I \delta_p^i |n_t^i - n_{t-1}^i| Z_t^i - \delta_f W_t \mathbb{1}_{n_t \neq n_{t-1}} - \sum_{i=0}^I \delta_s^i (n_t^i)^- Z_t^i$$

The transaction costs consist of three different components. The first term represent a transaction cost that is proportional to the change in value of the position in each asset. The second term is a fixed fraction of the total value of the portfolio which is payed only if the allocation is changed. The last term represents the fees payed to the broker for the shares borrowed to build a short position. The portfolio return between  $t$  and  $t + 1$  is thus given by

$$X_{t+1} = \frac{\Delta W_{t+1}}{W_t} = \sum_{i=0}^I \left[ a_t^i X_{t+1}^i - \delta_i |a_t^i - \tilde{a}_t^i| - \delta_s (a_t^i)^- \right] - \delta_f \mathbb{1}_{a_t \neq \tilde{a}_{t-1}} \quad (5.1)$$

where

$$X_{t+1}^i = \frac{\Delta Z_{t+1}^i}{Z_t^i}$$

is the return of the  $i$ -th stock between  $t$  and  $t + 1$ ,

$$a_t^i = \frac{n_t^i Z_t^i}{W_t}$$

is the fraction of wealth invested in the  $i$ -th stock between time  $t$  and  $t + 1$  and finally

$$\tilde{a}_t^i = \frac{n_{t-1}^i Z_t^i}{W_t} = \frac{a_{t-1}^i (1 + X_t^i)}{1 + X_t}$$

is the fraction of wealth invested in the  $i$ -th stock just before the reallocation. We assume that the agent invests all his wealth at each step, so that  $W_t$  can

be also interpreted as the value of his portfolio. This assumption leads to the following constraint on the portfolio weights

$$\sum_{i=0}^I a_t^i = 1 \quad \forall t \in \{0, 1, 2, \dots\} \quad (5.2)$$

We notice that we are neglecting the typical margin requirements on the short positions, which would reduce the available capital at time  $t$ . Considering margin requirements would lead to a more complex constraint on the portfolio weights which would be difficult to treat in the reinforcement learning framework. Plugging this constraint into Eq. (5.1), we obtain

$$X_{t+1} = X + \sum_{i=1}^I a_t^i (X_{t+1}^i - X) - \sum_{i=0}^I \left[ \delta_i |a_t^i - \tilde{a}_t^i| - \delta_s^i (a_t^i)^- \right] - \delta_f \mathbb{1}_{a_t \neq \tilde{a}_{t-1}} \quad (5.3)$$

which highlights the role of the risk-free asset as a benchmark for the portfolio returns. The total profit realized by the investor between  $t = 0$  and  $T$  is

$$\Pi_T = W_T - W_0 = \sum_{t=1}^T \Delta W_t = \sum_{t=1}^T W_t X_t$$

The portfolio return between  $t = 0$  and  $T$  is given by

$$X_{0,T} = \frac{W_T}{W_0} - 1 = \prod_{t=1}^T (1 + X_t) - 1$$

In order to cast the asset allocation problem in the reinforcement learning framework, we consider the log-return of the portfolio between  $t = 0$  and  $T$

$$R_{0,T} = \log \frac{W_T}{W_0} = \sum_{t=1}^T \log(1 + X_t) = \sum_{t=1}^T R_t \quad (5.4)$$

where  $R_{t+1}$  is the log-return of the portfolio between  $t$  and  $t + 1$

$$R_{t+1} = \log \left\{ 1 + \sum_{i=0}^I \left[ a_t^i X_{t+1}^i - \delta_i |a_t^i - \tilde{a}_t^i| - \delta_s^i (a_t^i)^- \right] - \delta_f \mathbb{1}_{a_t \neq \tilde{a}_{t-1}} \right\} \quad (5.5)$$

The portfolio return and log-return can be used as the reward function of a RL algorithm, either in a offline or in an online approach.

### 5.2.2 States

At each time step, the agent observes the state of the system to make his decisions. First, we consider the  $P + 1$  past returns of all risky assets, i.e.  $\{X_t, X_{t-1}, \dots, X_{t-P}\}$ . These input variables may be used to construct more complex features for example using some deep learning techniques, such as a deep auto-encoder. In order to properly incorporate the effects of transaction costs into his decision process, the agent must keep track of its current position  $\tilde{a}_t$ . Finally, we might consider some external variables  $Y_t$  that may be relevant to the trader, such as the common technical indicator used in practice. Summing up, we assume that the state of the system is composed in the following way

$$S_t = \{X_t, X_{t-1}, \dots, X_{t-P}, \tilde{a}_t, Y_t, Y_{t-1}, \dots, Y_{t-P}\} \quad (5.6)$$

### 5.2.3 Actions

The agent, or trading system, only specifies the portfolio weights  $a_t = (a_t^0, \dots, a_t^I)^T$  and therefore determines the allocation for the time interval  $[t, t + 1)$ . If we assume that the agent invests all of his capital at each time step and that short-selling is not allowed, then the portfolio weights should satisfy, for every  $t \in \{0, 1, \dots\}$ , the following constraints

$$\begin{cases} a_t^i \geq 0, \forall i \in \{0, \dots, I\} \\ \sum_{i=0}^I a_t^i = 1 \end{cases} \quad (5.7)$$

This constraint can be easily enforced by considering a parametric softmax policy. If short-selling is allowed, weights might also be negative. A simple approach is to assume that, for every  $t \in \{0, 1, \dots\}$ , the weights satisfy

$$\begin{cases} a_t^i \in \mathbb{R}, \forall i \in \{1, \dots, I\} \\ a_t^0 = 1 - \sum_{i=1}^I a_t^i \end{cases} \quad (5.8)$$

Since  $a_t^0$  is uniquely determined by the other weights, it is enough to define a policy that specifies the allocation in the risky assets, e.g. a Gaussian policy in  $\mathbb{R}^I$ . The obvious shortcoming of this approach is that the agent might enter in huge short positions, which is not realistic. A first observation is that, if the stochastic policy is concentrated around the origin, huge long or short positions would have very small probabilities of being selected. Moreover, we notice that the agent would pay large fees for entering into large short positions, independently of the trading profits. Therefore, we expect the agent to learn that short positions are very expensive and would therefore try to

avoid them.

Working in a continuous action space is computationally difficult and only few reinforcement learning algorithms are well-suited to this setting, e.g. policy gradient methods. A simpler approach is to reduce the action space to a discrete space. For instance, in the two assets scenario we might assume that  $a_t \in \{-1, 0, +1\}$ . Thus the agent may be long (+1), neutral (0) or short (-1) on the risky-asset. Working in a discrete action space is more simple, and standard value-based approaches might also be employed.





# Chapter 6

## Numerical Results for the Asset Allocation Problem

In this chapter we present the numerical results of some of the policy gradient algorithms discussed in Chapter 4 for the asset allocation problem. Two different type of markets are analyzed: a market with only one risky asset and a market where multiple risky assets are available, for which finding a trading strategy is more difficult since the state and action spaces are much larger. The learning algorithms are first applied both in their risk-neutral version and risk-sensitive formulation to synthetically generated data, which present profitably tradable features. Once the behavior of these algorithms is validated in this controlled environment, the various methods are also applied to historical price series.

### 6.1 Synthetic Risky Asset

To test the different reinforcement learning methods in a controlled environment, we generated log-price series for the risky asset as random walks with autoregressive trend processes. The two-parameter model is thus given by

$$\begin{aligned} z_t &= z_{t-1} + \beta_{t-1} + \kappa \epsilon_t \\ \beta_t &= \alpha \beta_{t-1} + \nu_t \end{aligned}$$

We then define the synthetic price series as

$$Z_t = \exp \left( \frac{z_t}{\max_t z_t - \min_t z_t} \right)$$

This model is often taken as a benchmark test in the automated trading literature [55], because the price series generated in this way presents some

patterns that can be profitably exploited. Moreover the model is stationary and therefore the policy learned on the training set should generalize well on the test set, also known as backtest in the financial jargon. Thus we would expect our learning algorithms to perform well on this test case. If this wasn't the case, we should go back and improve the learning algorithms. In this setting, we compare the results of three long-short strategies obtained with ARAC, PGPE and NPGPE in both the risk-neutral and risk-sensitive framework. This means that the agent can either go long on the risky asset (i.e.  $a_t^1 = 1$ ) or short the security (i.e.  $a_t^1 = -1$ ) and invest the proceedings in the risk-less asset. Given the current conditions of the financial markets, we always assume a risk-free rate  $X = 0$ . Let us describe in more detail the choice we made for each of the algorithms.

### 6.1.1 Learning Algorithm Specifications

#### 6.1.1.1 ARAC

For the ARAC algorithm we considered a Boltzmann exploration policy on the two actions  $a_t^1 \in \{-1, 1\}$  and a linear critic in which the features coincide with the agent's observation of the system state. This critic is extremely simple and there is surely some work to be done to improve it.

#### 6.1.1.2 PGPE

For the PGPE algorithm we considered a binary deterministic controller

$$F_\theta(s) = \text{sign}(\theta \cdot s)$$

where the parameters and the state also include a bias term. The controller parameters are sampled from a multi-variate Gaussian distribution

$$\theta \sim \mathcal{N}(\mu, \text{diag}(\sigma))$$

#### 6.1.1.3 NPGPE

In NPGPE, we used the same controller as for PGPE but we assumed that the controller parameters are sampled from a Gaussian distribution parameterized by its mean and Cholesky factor

$$\theta \sim \mathcal{N}(\mu, C^T C)$$

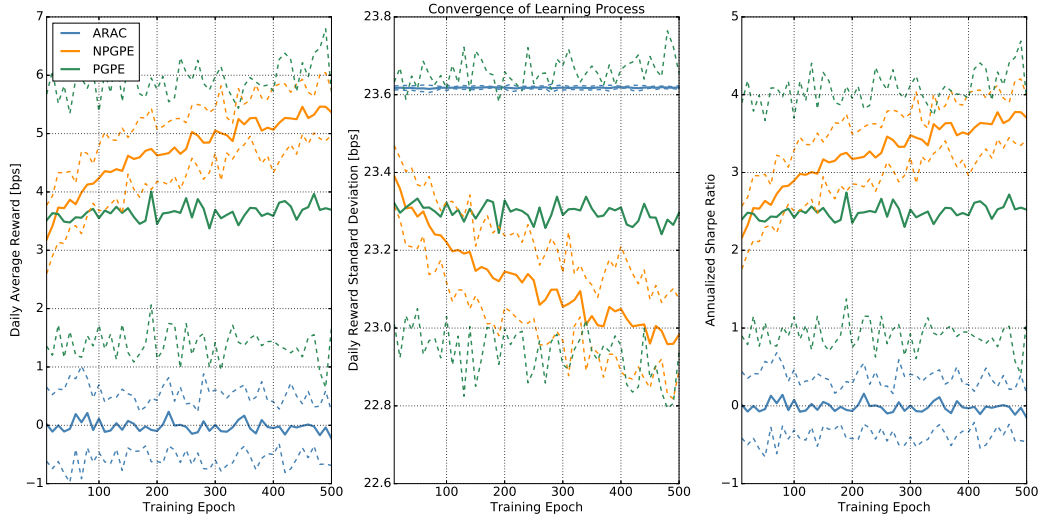


Figure 6.1: Risk-neutral learning process for the asset allocation problem with one synthetic risky asset.

### 6.1.2 Experimental Setup

All the algorithms were tested on the same price series of size 9000, generated from the process described above using  $\alpha = 0.9$  and  $\kappa = 3$ . The learning process consisted of 500 training epochs on the first 7000 days of the series with a learning rate that decreased at each epoch according to a polynomial schedule. The trained agents were subsequently backtested on the final 2000 days, during which the agents kept learning online in order to try to adapt to the changing environment. The results that we present are the average of 10 independent experiments that used slightly different random initialization of the policy parameters.

### 6.1.3 Risk-Neutral Framework

Let us first discuss the case with no transaction costs. Figure 6.1 shows the learning curves for the three risk-neutral algorithms in terms of average daily reward, which is the quantity being maximized by the algorithms, the daily reward standard deviation and the annualized Sharpe ratio. The first thing we observe is the ARAC algorithm seems not to be improving the trading strategy as the training epochs go by. The average reward obtained is close to zero and will be surely be negative once transaction costs are introduced. On the other hand, PGPE converges very quickly to a profitable strategy characterized by a large variance in the rewards. The learned policy is however suboptimal compared to the one found by NPGPE, which is better

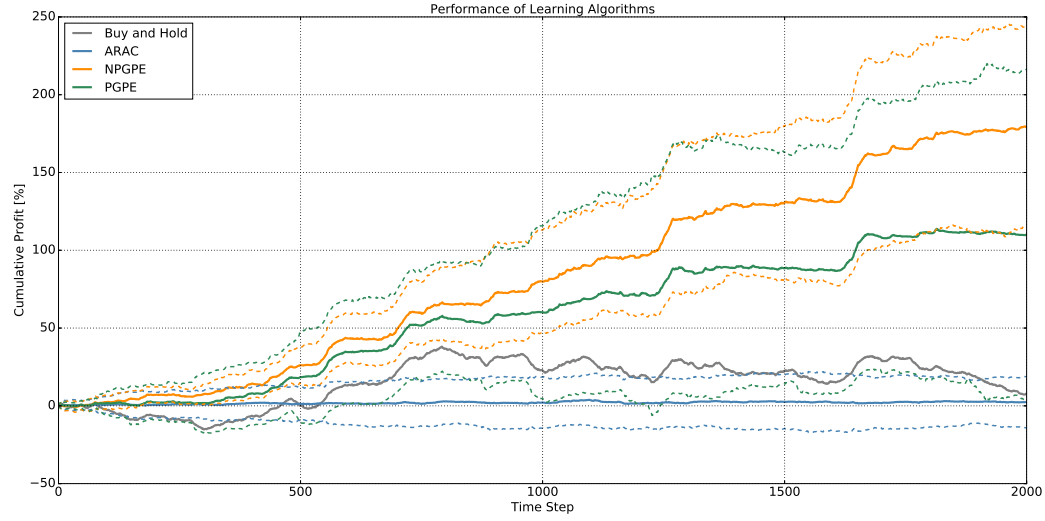


Figure 6.2: Backtest performance of trained trading systems for the asset allocation problem with one synthetic risky asset.

in all three measures considered. This confirms the fact that natural version of a policy gradient algorithm often outperforms its “vanilla” counterpart. It is interesting to notice that NPGPE yields a learning curve for the Sharpe ratio very similar to the one for the average reward. Even if the algorithm is risk-neutral, it manages to improve a risk-sensitive measure at the same time of the average reward. This might be simply a peculiarity of the very simple model assumed for the risky asset.

Figure 6.2 compares the backtest performances of the three learned policies and a Buy and Hold strategy, which simply consists in investing all the available capital in the risky asset. Let us repeat that the solid lines are the averages of 10 independent experiments, which allows us to determine the 95% confidence intervals represented with the dashed lines. We clearly see that PGPE and NPGPE easily beat the market, realizing a total profit of 101.0% and 179.0% respectively against the 7.8% profit of the Buy and Hold strategy over the same period. More statistics of the trading strategies are reported in Table 6.1

#### 6.1.4 Risk-Sensitive Framework

In this section we present the results in the risk-sensitive framework, in which all the algorithms optimize the Sharpe ratio of the policy. Figure 6.3 shows the learning curves for the three risk-sensitive algorithms RSARAC, RSPGPE and RSNPGPE, which show similar behaviors to their risk-neutral counterparts. However, it is surprising that the strategies learned by RSPGPE

	Buy and Hold	ARAC	NPGPE	PGPE
Total Return	7.8%	2.4%	179.0%	101.0%
Daily Sharpe	0.27	0.09	3.49	2.43
Monthly Sharpe	0.19	0.08	2.32	1.50
Yearly Sharpe	0.23	0.08	1.36	0.87
Max Drawdown	-22.4%	-11.5%	-4.6%	-1.0%
Avg Drawdown	-1.8%	-1.5%	-0.5%	-0.9%
Avg Up Month	2.87%	1.1%	2.5%	2.6%
Avg Down Month	-2.6%	-1.0%	-0.9%	-1.5%
Win Year [	40.0%	48.0%	92.0%	84.0%
Win 12m	56.4%	51.3%	99.6%	92.9%

Table 6.1: Backtest statistics of the risk-neutral trading strategies for the asset allocation problem with one synthetic risky asset.

and RSNPGPE have smaller Sharpe ratio than in the risk-neutral version of these algorithms, which optimize a different quantity. Figure 6.4 shows the backtest performances for the three risk-sensitive trading strategies. Again, RSPGPE and NPGPE beat the market even if in smaller measure than in the risk-neutral setting.

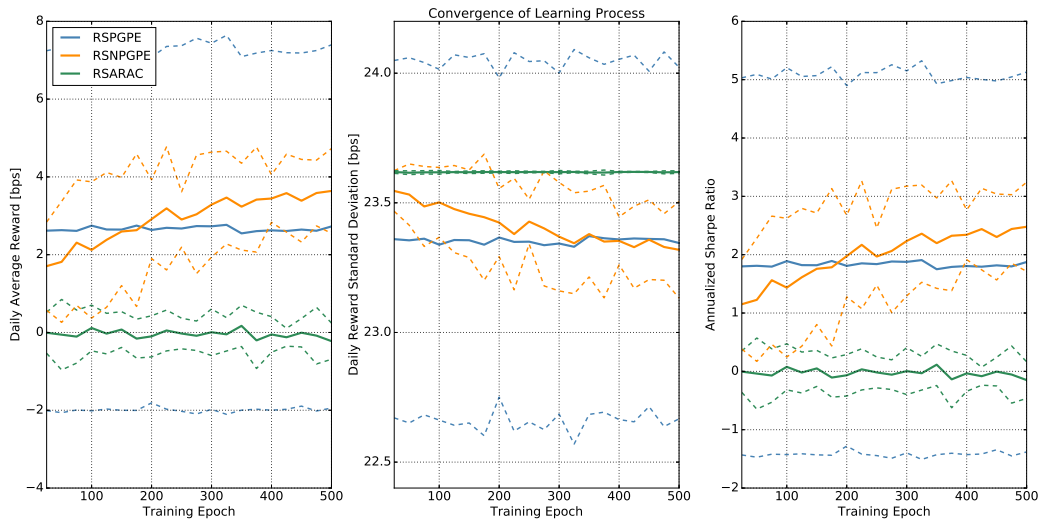


Figure 6.3: Risk-sensitive learning process for the asset allocation problem with one synthetic risky asset.

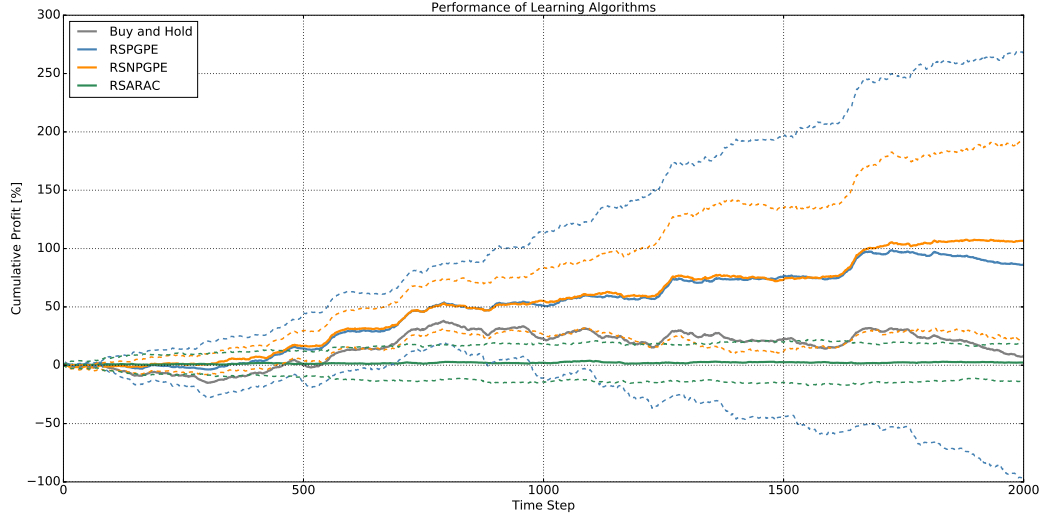


Figure 6.4: Backtest performance of the risk-sensitive trading strategies for the asset allocation problem with one synthetic risky asset.

	Buy and Hold	RSPGPE	RSNPGPE	RSARAC
Total Return	7.8%	86.0%	106.5%	2.3%
Daily Sharpe	0.27	1.77	2.40	0.09
Monthly Sharpe	0.19	1.18	1.65	0.08
Yearly Sharpe	0.23	0.70	0.95	0.07
Max Drawdown	-22.4%	-13.8%	-6.5%	-11.5%
Avg Drawdown	-1.8%	-1.2%	-0.7%	-1.5%
Avg Up Month	2.9%	2.8%	2.2%	1.1%
Avg Down Month	-2.6%	-1.6%	-1.1%	-1.0%
Win Year %	40.0%	70.0%	86.0%	48.0%
Win 12m %	56.4%	78.2%	96.9%	50.4%

Table 6.2: Backtest statistics of the risk-sensitive trading strategies for the asset allocation problem with one synthetic risky asset.

## **6.2 Historic Risky Asset**

## **6.3 Multiple Synthetic Risky Assets**

## **6.4 Historic Multiple Risky Assets**





# Chapter 7

## Conclusions

### 7.1 Summary

### 7.2 Further Developments



# Bibliography

- [1] A. AGARWAL, P. BARTLETT, AND M. DAMA, *Optimal allocation strategies for the dark pool problem*, arXiv preprint arXiv:1003.2245, (2010).
- [2] R. ALMGREN AND N. CHRISS, *Optimal execution of portfolio transactions*, Journal of Risk, 3 (2001), pp. 5–40.
- [3] A. ARAPOSTATHIS, V. S. BORKAR, E. FERNÁNDEZ-GAUCHERAND, M. K. GHOSH, AND S. I. MARCUS, *Discrete-time controlled markov processes with average cost criterion: a survey*, SIAM Journal on Control and Optimization, 31 (1993), pp. 282–344.
- [4] A. BASU, T. BHATTACHARYYA, AND V. S. BORKAR, *A learning algorithm for risk-sensitive cost*, Mathematics of Operations Research, 33 (2008), pp. 880–898.
- [5] N. BÄUERLE AND U. RIEDER, *Markov decision processes with applications to finance*, Springer Science & Business Media, 2011.
- [6] J. BAXTER AND P. L. BARTLETT, *Infinite-horizon policy-gradient estimation*, Journal of Artificial Intelligence Research, 15 (2001), pp. 319–350.
- [7] S. D. BEKIROU, *Heterogeneous trading strategies with adaptive fuzzy actor–critic reinforcement learning: A behavioral approach*, Journal of Economic Dynamics and Control, 34 (2010), pp. 1153–1170.
- [8] F. BERTOLUZZO AND M. CORAZZA, *Testing different reinforcement learning configurations for financial trading: Introduction and applications*, Procedia Economics and Finance, 3 (2012), pp. 68–77.
- [9] ———, *Q-learning-based financial trading systems with applications*, University Ca’Foscari of Venice, Dept. of Economics Working Paper Series No, 15 (2014).

- [10] —, *Reinforcement learning for automated financial trading: Basics and applications*, in Recent Advances of Neural Network Models and Applications, Springer, 2014, pp. 197–213.
- [11] D. P. BERTSEKAS, *Dynamic programming and optimal control*, vol. 1, Athena Scientific, Belmont, 1995.
- [12] D. P. BERTSEKAS AND S. E. SHREVE, *Stochastic optimal control: the discrete-time case*, vol. 23, Academic Press New York, 1978.
- [13] D. P. BERTSEKAS AND J. N. TSITSIKLIS, *Neuro-Dynamic Programming*, Optimization and neural computation series, Athena Scientific, Belmont, 1 ed., 1996.
- [14] S. BHATNAGAR, R. S. SUTTON, M. GHAVAMZADEH, AND M. LEE, *Natural actor-critic algorithms*, Automatica, 45 (2009), pp. 2471–2482.
- [15] C. M. BISHOP, *Pattern Recognition and Machine Learning*, Springer, 2006.
- [16] V. S. BORKAR, *A sensitivity formula for risk-sensitive cost and the actor-critic algorithm*, Systems & Control Letters, 44 (2001), pp. 339–346.
- [17] —, *Q-learning for risk-sensitive control*, Mathematics of operations research, 27 (2002), pp. 294–311.
- [18] V. S. BORKAR AND S. P. MEYN, *Risk-sensitive optimal control for markov decision processes with monotone cost*, Mathematics of Operations Research, 27 (2002), pp. 192–209.
- [19] L. BUSONIU, R. BABUSKA, B. DE SCHUTTER, AND D. ERNST, *Reinforcement learning and dynamic programming using function approximators*, vol. 39, CRC press, 2010.
- [20] P. X. CASQUEIRO AND A. J. RODRIGUES, *Neuro-dynamic trading methods*, European Journal of Operational Research, 175 (2006), pp. 1400–1412.
- [21] N. CHAPADOS AND Y. BENGIO, *Cost functions and model combination for var-based asset allocation using neural networks*, IEEE Transactions on Neural Networks, 12 (2001), pp. 890–906.
- [22] M. CHOEY AND A. S. WEIGEND, *Nonlinear trading models through sharpe ratio maximization*, International Journal of Neural Systems, 8 (1997), pp. 417–431.

- [23] Y. CHOW, A. TAMAR, S. MANNOR, AND M. PAVONE, *Risk-sensitive and robust decision-making: a cvar optimization approach*, in Advances in Neural Information Processing Systems, 2015, pp. 1522–1530.
- [24] M. CORAZZA AND A. SANGALLI, *Q-learning vs. sarsa: comparing two intelligent stochastic control approaches for financial trading*.
- [25] J. CUMMING, D. ALRAJEH, AND L. DICKENS, *An investigation into the use of reinforcement learning techniques within the algorithmic trading domain*, (2015).
- [26] M. A. DEMPSTER AND V. LEEMANS, *An automated fx trading system using adaptive reinforcement learning*, Expert Systems with Applications, 30 (2006), pp. 543–552.
- [27] M. A. H. DEMPSTER AND Y. S. ROMAHI, *Intraday FX trading: An evolutionary reinforcement learning approach*, Springer, 2002.
- [28] Y. DENG, F. BAO, Y. KONG, Z. REN, AND Q. DAI, *Deep direct reinforcement learning for financial signal representation and trading*, IEEE Transactions on Neural Networks and Learning Systems, (2016).
- [29] Y. DENG, Y. KONG, F. BAO, AND Q. DAI, *Sparse coding inspired optimal trading system for hft industry*, IEEE Transactions on Industrial Informatics, 11 (2015), pp. 467–475.
- [30] X. DU, J. ZHAI, AND K. LV, *Algorithm trading using q-learning and recurrent reinforcement learning*, positions, 1, p. 1.
- [31] T. ELDER, *Creating algorithmic traders with hierarchical reinforcement learning*, (2008).
- [32] L. A. FELDKAMP, D. V. PROKHOROV, C. F. EAGEN, AND F. YUAN, *Enhanced multi-stream kalman filter training for recurrent networks*, in Nonlinear Modeling, Springer, 1998, pp. 29–53.
- [33] K. GANCHEV, Y. NEVMYVAKA, M. KEARNS, AND J. W. VAUGHAN, *Censored exploration and the dark pool problem*, Communications of the ACM, 53 (2010), pp. 99–107.
- [34] C. GOLD, *Fx trading via recurrent reinforcement learning*, in IEEE 2003 IEEE International Conference on Computational Intelligence for Financial Engineering. Proceedings, 2003.

- [35] I. GOODFELLOW, Y. BENGIO, AND A. COURVILLE, *Deep learning*. Book in preparation for MIT Press, 2016.
- [36] T. HASTIE, R. TIBSHIRANI, AND J. FRIEDMAN, *The elements of statistical learning: data mining, inference, and prediction*, Springer, 2009.
- [37] D. HENDRICKS AND D. WILCOX, *A reinforcement learning extension to the almgren-chriss framework for optimal trade execution*, in IEEE Conference on Computational Intelligence for Financial Engineering & Economics (CIFEr), 2014, pp. 457–464.
- [38] H. JAEGER, *Tutorial on training recurrent neural networks, covering BPPT, RTRL, EKF and the echo state network approach*, GMD-Forschungszentrum Informationstechnik, 2002.
- [39] M. S. JOSHI, *C++ design patterns and derivatives pricing*, vol. 2, Cambridge University Press, 2008.
- [40] S. KAKADE, *A natural policy gradient.*, in NIPS, vol. 14, 2001, pp. 1531–1538.
- [41] K. KAMIJO AND T. TANIGAWA, *Stock price pattern recognition-a recurrent neural network approach*, in IEEE 1990 IJCNN International Joint Conference on Neural Networks, 1990.
- [42] M. KEARNS AND Y. NEVMYVAKA, *Machine learning for market microstructure and high frequency trading*, High-Frequency Trading—New Realities for Traders, Markets and Regulators, (2013), pp. 91–124.
- [43] V. R. KONDA AND J. N. TSITSIKLIS, *Actor-critic algorithms.*, in NIPS, vol. 13, 1999, pp. 1008–1014.
- [44] H. KUSHNER AND G. G. YIN, *Stochastic approximation and recursive algorithms and applications*, vol. 35, Springer Science & Business Media, 2003.
- [45] P. L.A. AND M. GHAVAMZADEH, *Actor-critic algorithms for risk-sensitive mdps*, in Advances in Neural Information Processing Systems 26, C. J. C. Burges, L. Bottou, M. Welling, Z. Ghahramani, and K. Q. Weinberger, eds., Curran Associates, Inc., 2013, pp. 252–260.
- [46] S. LARUELLE, C.-A. LEHALLE, AND G. PAGES, *Optimal split of orders across liquidity pools: a stochastic algorithm approach*, SIAM Journal on Financial Mathematics, 2 (2011), pp. 1042–1076.

- [47] ———, *Optimal posting price of limit orders: learning by trading*, Mathematics and Financial Economics, 7 (2013), pp. 359–403.
- [48] H. LI, C. H. DAGLI, AND D. ENKE, *Short-term stock market timing prediction under reinforcement learning schemes*, in 2007 IEEE International Symposium on Approximate Dynamic Programming and Reinforcement Learning, 2007, pp. 233–240.
- [49] S. MAHADEVAN, *Average reward reinforcement learning: Foundations, algorithms, and empirical results*, Machine learning, 22 (1996), pp. 159–195.
- [50] H. MARKOWITZ, *Portfolio selection*, The journal of finance, 7 (1952), pp. 77–91.
- [51] A. MIYAMAE, Y. NAGATA, I. ONO, AND S. KOBAYASHI, *Natural policy gradient methods with parameter-based exploration for control tasks*, in Advances in neural information processing systems, 2010, pp. 1660–1668.
- [52] J. MOODY AND M. SAFFELL, *Learning to trade via direct reinforcement*, Neural Networks, IEEE Transactions on, 12 (2001), pp. 875–889.
- [53] J. MOODY, M. SAFFELL, Y. LIAO, AND L. WU, *Reinforcement learning for trading systems*, in Decision Technologies for Computational Finance: Proceedings of the fifth International Conference Computational Finance, vol. 2, Springer Science & Business Media, 2013, p. 129.
- [54] J. MOODY AND L. WU, *Optimization of trading systems and portfolios*, in Proceedings of the IEEE/IAFE 1997 Computational Intelligence for Financial Engineering (CIFEr), 1997, pp. 300–307.
- [55] J. MOODY, L. WU, Y. LIAO, AND M. SAFFELL, *Performance functions and reinforcement learning for trading systems and portfolios*, Journal of Forecasting, 17 (1998), pp. 441–470.
- [56] Y. NEVMYVAKA, Y. FENG, AND M. KEARNS, *Reinforcement learning for optimized trade execution*, in Proceedings of the 23rd international conference on Machine learning, ACM, 2006, pp. 673–680.
- [57] J. NOCEDAL AND S. WRIGHT, *Numerical optimization*, Springer Science & Business Media, 2006.

- [58] J. O, J. LEE, J. W. LEE, AND B.-T. ZHANG, *Adaptive stock trading with dynamic asset allocation using reinforcement learning*, Information Sciences, 176 (2006), pp. 2121–2147.
- [59] G. PAGES, *Introduction to Numerical Probability for Finance*, LPMA-Université Pierre et Marie Curie, 2016.
- [60] J. PETERS, K. MÜLLING, AND Y. ALTUN, *Relative entropy policy search.*, in AAAI, Atlanta, 2010.
- [61] J. PETERS AND S. SCHAAL, *Policy gradient methods for robotics*, in Intelligent Robots and Systems, 2006 IEEE/RSJ International Conference on, IEEE, 2006, pp. 2219–2225.
- [62] —, *Reinforcement learning of motor skills with policy gradients*, Neural networks, 21 (2008), pp. 682–697.
- [63] L. PRASHANTH AND M. GHAVAMZADEH, *Actor-critic algorithms for risk-sensitive reinforcement learning.*, arXiv preprint arXiv:1403.6530, (2014).
- [64] M. L. PUTERMAN, *Markov decision processes: discrete stochastic dynamic programming*, John Wiley & Sons, 1994.
- [65] E. W. SAAD, D. V. PROKHOROV, AND D. C. WUNSCH, *Comparative study of stock trend prediction using time delay, recurrent and probabilistic neural networks*, IEEE Transactions on Neural Networks, 9 (1998), pp. 1456–1470.
- [66] C. SANDERSON, *Armadillo: An open source c++ linear algebra library for fast prototyping and computationally intensive experiments*, (2010).
- [67] M. SATO AND S. KOBAYASHI, *Variance-penalized reinforcement learning for risk-averse asset allocation*, in Intelligent Data Engineering and Automated Learning—IDEAL 2000. Data Mining, Financial Engineering, and Intelligent Agents, Springer, 2000, pp. 244–249.
- [68] —, *Average-reward reinforcement learning for variance penalized markov decision problems*, in Proceedings of the Eighteenth International Conference on Machine Learning, ICML '01, San Francisco, CA, USA, 2001, Morgan Kaufmann Publishers Inc., pp. 473–480.
- [69] F. SEHNKE ET AL., *Parameter exploring policy gradients and their implications*, PhD thesis, Technische Universität München, 2012.



- [70] F. SEHNKE, C. OSENDORFER, T. RÜCKSTIESS, A. GRAVES, J. PETERS, AND J. SCHMIDHUBER, *Policy gradients with parameter-based exploration for control*, in Artificial Neural Networks-ICANN 2008, Springer, 2008, pp. 387–396.
- [71] —, *Parameter-exploring policy gradients*, Neural Networks, 23 (2010), pp. 551–559.
- [72] W. F. SHARPE, *The sharpe ratio*, The journal of portfolio management, 21 (1994), pp. 49–58.
- [73] D. SILVER, G. LEVER, N. HEESS, T. DEGRIS, D. WIERSTRA, AND M. RIEDMILLER, *Deterministic policy gradient algorithms*, in ICML, 2014.
- [74] M. J. SOBEL, *The variance of discounted markov decision processes*, Journal of Applied Probability, (1982), pp. 794–802.
- [75] R. S. SUTTON AND A. G. BARTO, *Introduction to reinforcement learning*, vol. 135, MIT Press Cambridge, 1998.
- [76] R. S. SUTTON, D. A. MCALLESTER, S. P. SINGH, Y. MANSOUR, ET AL., *Policy gradient methods for reinforcement learning with function approximation.*, in NIPS, vol. 99, 1999, pp. 1057–1063.
- [77] C. SZEPESVÁRI, *Algorithms for reinforcement learning*, Synthesis lectures on artificial intelligence and machine learning, 4 (2010), pp. 1–103.
- [78] A. TAMAR, D. D. CASTRO, AND S. MANNOR, *Temporal difference methods for the variance of the reward to go*, in Proceedings of the 30th International Conference on Machine Learning (ICML-13), 2013, pp. 495–503.
- [79] A. TAMAR, Y. CHOW, M. GHAVAMZADEH, AND S. MANNOR, *Policy gradient for coherent risk measures*, in Advances in Neural Information Processing Systems, 2015, pp. 1468–1476.
- [80] A. TAMAR, D. DI CASTRO, AND S. MANNOR, *Policy gradients with variance related risk criteria*, in Proceedings of the 29th International Conference on Machine Learning, 2012, pp. 387–396.
- [81] A. TAMAR AND S. MANNOR, *Variance adjusted actor critic algorithms*, arXiv preprint arXiv:1310.3697, (2013).

- [82] Z. TAN, C. QUEK, AND P. Y. CHENG, *Stock trading with cycles: A financial application of anfis and reinforcement learning*, Expert Systems with Applications, 38 (2011), pp. 4741–4755.
- [83] R. S. TSAY, *Analysis of financial time series*, vol. 543, John Wiley & Sons, 2005.
- [84] P. J. WERBOS, *Backpropagation through time: what it does and how to do it*, Proceedings of the IEEE, 78 (1990), pp. 1550–1560.
- [85] M. WIERING AND M. VAN OTTERLO, *Reinforcement Learning: State-of-the-Art*, vol. 12 of Adaptation, Learning, and Optimization, Springer, 1 ed., 2012.
- [86] R. J. WILLIAMS AND D. ZIPSER, *A learning algorithm for continually running fully recurrent neural networks*, Neural computation, 1 (1989), pp. 270–280.
- [87] S. YANG, M. PADDRIK, R. HAYES, A. TODD, A. KIRILENKO, P. BELING, AND W. SCHERER, *Behavior based learning in identifying high frequency trading strategies*, in IEEE Conference on Computational Intelligence for Financial Engineering & Economics (CIFEr), 2012.
- [88] T. ZHAO, H. HACHIYA, G. NIU, AND M. SUGIYAMA, *Analysis and improvement of policy gradient estimation*, in Advances in Neural Information Processing Systems, 2011, pp. 262–270.
- [89] T. ZHAO, G. NIU, N. XIE, J. YANG, AND M. SUGIYAMA, *Regularized policy gradients: Direct variance reduction in policy gradient estimation*, in Proceedings of The 7th Asian Conference on Machine Learning, 2015, pp. 333–348.