

Todo list

Add description of stochastic policies 27

Contents

1	Introduction	1
2	Discrete-Time Stochastic Optimal Control	3
2.1	Sequential Decision Problems	3
2.2	Markov Decision Processes	4
2.3	Policies	5
2.4	Discounted Reward Formulation	6
2.5	Average Reward Formulation	9
2.6	Risk-Sensitive Formulation	10
2.7	Policy Evaluation and Policy Iteration	12
3	Reinforcement Learning	13
3.1	The Reinforcement Learning Problem	14
3.2	Model-Free RL Methods	14
3.2.1	Model Approximation	16
3.2.2	Value Approximation	16
3.2.3	Policy Approximation	17
4	Policy Gradient	19
4.1	Basics of Policy Gradient Methods	20
4.1.1	Risk-Neutral Framework	21
4.1.2	Risk-Sensitive Framework	21
4.2	Finite Differences	23
4.3	Likelihood Ratio Methods	23
4.3.1	Monte Carlo Policy Gradient	24
4.3.2	GPOMDP	26
4.3.3	Risk-Sensitive Monte Carlo Policy Gradient	27
4.3.4	Stochastic Policies	27
4.3.5	Policy Gradient with Parameter Exploration	27
4.4	Policy Gradient Theorem	31
4.4.1	Risk-Neutral Setting	31

4.4.2	Risk-Sensitive Setting	33
4.4.3	GPOMDP from the Policy Gradient Theorem	33
4.4.4	Compatible Function Approximation	33
4.4.5	Actor-Critic Policy Gradient	33
4.4.6	Natural Policy Gradient	34
4.5	Risk-Sensitive Policy Gradient	34
4.5.1	Risk-Sensitive Actor-Critic Algorithm	35
4.6	Parameter-Based Policy Gradient Theorem	36
4.6.1	Risk-Neutral Setting	36
4.6.2	Risk-Sensitive Setting	36
4.6.3	Natural Policy Gradient	36
5	Asset Allocation	37
5.1	Reward Function	37
5.2	States	39
5.3	Actions	40

List of Figures

2.1	Agent-environment interaction in sequential decision problems.	5
3.1	Solution process for the control problem.	15

List of Tables

List of Algorithms

4.1	General setup for a policy gradient algorithm.	20
4.2	Episodic REINFORCE policy gradient estimate	26

Chapter 1

Introduction

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 [11], [62], [10].

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.

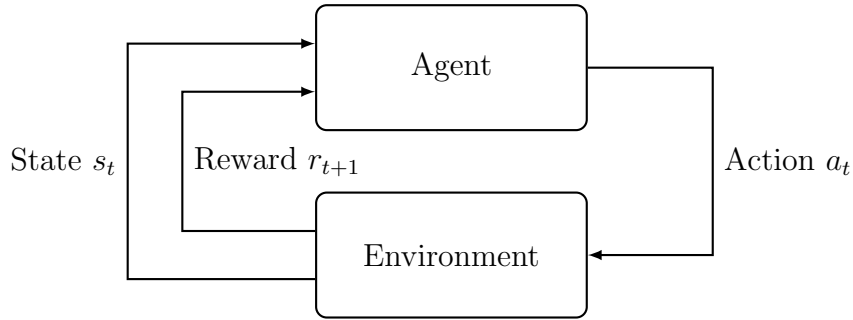


Figure 2.1: Agent-environment interaction in sequential decision problems.

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 π 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 π , 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 $t = 0, 1, \dots$*

$$\begin{cases} S_0 = s_0 \\ A_t \sim \pi(S_t, \cdot) \\ S_{t+1} \sim \mathcal{P}(S_t, A_t, \cdot) \end{cases} \quad (2.3)$$

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

$$\begin{aligned}\mathcal{P}_\pi(s, s') &= \int_{\mathbb{A}} \pi(s, a) \mathcal{P}(s, a, s') da \\ \mathcal{R}_\pi(s) &= \int_{\mathbb{A}} \pi(s, a) \mathcal{R}(s, a) da\end{aligned}\tag{2.4}$$

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 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*

$$G_t = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}\tag{2.5}$$

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 π*

$$V_\pi(s) = \mathbb{E}_\pi[G_t | S_t = s]\tag{2.6}$$

where \mathbb{E}_π indicates that all the actions are selected according to policy π . 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 π*

$$Q_\pi(s, a) = \mathbb{E}_\pi [G_t | S_t = s, A_t = a] \quad (2.7)$$

We have the following relationship between V_π and Q_π

$$V_\pi(s) = \int_{\mathbb{A}} \pi(s, a) Q_\pi(s, a) da \quad (2.8)$$

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

Proposition 2.4.1 (Bellman Expectation Equations).

$$V_\pi(s) = \mathcal{R}_\pi(s) + \gamma T_\pi V_\pi(s) \quad (2.9)$$

$$\begin{aligned} Q_\pi(s, a) &= \mathcal{R}(s, a) + \gamma T_a V_\pi(s) \\ &= \mathcal{R}(s, a) + \gamma \int_{\mathbb{S}} \mathcal{P}(s, a, s') \int_{\mathbb{A}} \pi(s', a') Q_\pi(s', a') da' ds' \end{aligned} \quad (2.10)$$

where we denoted by T_a (resp. T_π) the transition operator for action a (resp. for policy π)

$$\begin{aligned} T_a V(s) &= \int_{\mathbb{S}} \mathcal{P}(s, a, s') V_\pi(s') ds' \\ T_\pi V(s) &= \int_{\mathbb{A}} \pi(s, a) \int_{\mathbb{S}} \mathcal{P}(s, a, s') V_\pi(s') ds' da \end{aligned} \quad (2.11)$$

If we introduce the Bellman expectation operator B_π , defined as

$$B_\pi V_\pi(s) = \mathcal{R}_\pi(s) + \gamma T_\pi V_\pi(s) \quad (2.12)$$

Then Eq. (2.9) can be written as a fixed-point equation

$$V_\pi(s) = B_\pi V_\pi(s) \quad (2.13)$$

which, under some simple assumptions on the reward functions, admits a unique solution by the contraction mapping theorem. A similar argument applies to Eq. (2.10). The agent's goal is to select a policy π_* 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*

$$V_*(s) = \sup_{\pi} V_{\pi}(s) \quad (2.14)$$

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*

$$Q_*(s, a) = \sup_{\pi} Q_{\pi}(s, a) \quad (2.15)$$

The optimal value functions satisfy the following Bellman equations.

Proposition 2.4.2 (Bellman Optimality Equations).

$$V_*(s) = \sup_a Q_*(s, a) = \sup_a \{ \mathcal{R}(s, a) + \gamma T_a V_*(s) \} \quad (2.16)$$

$$\begin{aligned} Q_*(s, a) &= \mathcal{R}(s, a) + \gamma T_a V_*(s) \\ &= \mathcal{R}(s, a) + \gamma \int_{\mathbb{S}} \mathcal{P}(s, a, s') \sup_a Q_*(s', a') ds' \end{aligned} \quad (2.17)$$

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

$$\pi \succeq \pi' \Leftrightarrow V_{\pi}(s) \geq V_{\pi'}(s) \quad \forall s \in \mathbb{S} \quad (2.18)$$

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 π_* 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

$$a = \arg \sup_a Q_*(s, a) \quad (2.19)$$

We see that this policy is deterministic and only depends on the current state of the system. The Bellman equations provide the basis for the *value iteration* and *policy iteration* algorithms [72], which however require knowledge of the Markov transition kernel.

2.5 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], [47] 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 2.5.1 (Average Reward). *The average reward ρ_π associated to a policy π is defined as*

$$\begin{aligned}\rho_\pi &= \lim_{T \rightarrow \infty} \frac{1}{T} \mathbb{E}_\pi \left[\sum_{t=0}^{T-1} R_{t+1} \right] \\ &= \mathbb{E}_{\substack{S \sim d_\pi \\ A \sim \pi}} [\mathcal{R}(S, A)] \\ &= \int_{\mathbb{S}} d_\pi(s) \int_{\mathbb{A}} \pi(s, a) \mathcal{R}(s, a) da ds\end{aligned}\tag{2.20}$$

where d_π is the stationary distribution of the Markov process induced by π .

The agent aims to find an *average optimal* policy

$$\pi_* = \arg \sup_{\pi} \rho_\pi\tag{2.21}$$

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

Definition 2.5.2 (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 π*

$$V_\pi(s) = \mathbb{E}_\pi \left[\sum_{t=0}^{\infty} (R_{t+1} - \rho_\pi) \mid S_0 = s \right]\tag{2.22}$$

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. ρ_π serves as a baseline that allows to avoid divergence in the value function definition.

Definition 2.5.3 (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 π*

$$Q_\pi(s, a) = \mathbb{E}_\pi \left[\sum_{t=0}^{\infty} (R_{t+1} - \rho_\pi) \middle| S_0 = s, A_0 = a \right] \quad (2.23)$$

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

$$V_\pi(s) = \int_{\mathbb{A}} \pi(s, a) Q_\pi(s, a) \quad (2.24)$$

The value functions satisfy the following Bellman equation

Proposition 2.5.1 (Bellman Expectation Equations).

$$\rho_\pi + V_\pi(s) = \mathcal{R}_\pi(s) + T_\pi V_\pi(s) \quad (2.25)$$

$$\begin{aligned} \rho_\pi + Q_\pi(s, a) &= \mathcal{R}(s, a) + T_a V_\pi(s) \\ &= \mathcal{R}(s, a) + \int_{\mathbb{S}} \mathcal{P}(s, a, s') \int_{\mathbb{A}} \pi(s', a') Q_\pi(s', a') da' ds' \end{aligned} \quad (2.26)$$

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.

2.6 Risk-Sensitive Formulation

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 [61], the authors consider the long-run variance of π

Definition 2.6.1 (Long-Run Variance). *The long-run variance Λ_π under policy π is defined as*

$$\Lambda_\pi = \lim_{T \rightarrow \infty} \frac{1}{T} \mathbb{E}_\pi \left[\sum_{t=0}^{T-1} (R_{t+1} - \rho_\pi)^2 \right] \quad (2.27)$$

The long-run variance can be decomposed as follows

$$\Lambda_\pi = \eta_\pi - \rho_\pi^2 \quad (2.28)$$

where η_π is the average square reward per step

Definition 2.6.2 (Average Square Reward).

$$\begin{aligned}\eta_\pi &= \lim_{T \rightarrow \infty} \frac{1}{T} \mathbb{E}_\pi \left[\sum_{t=0}^{T-1} R_{t+1}^2 \right] \\ &= \mathbb{E}_{\substack{S \sim d_\pi \\ A \sim \pi}} [\mathcal{M}(S, A)] \\ &= \int_{\mathbb{S}} d_\pi(s) \int_{\mathbb{A}} \pi(s, a) \mathcal{M}(s, a)\end{aligned}\tag{2.29}$$

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

$$\mathcal{M}(s, a) = \mathbb{E} [R_{t+1}^2 | S_t = s, A_t = a] \tag{2.30}$$

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

Definition 2.6.3 (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 π*

$$U_\pi(s) = \mathbb{E}_\pi \left[\sum_{t=0}^{\infty} (R_{t+1}^2 - \eta_\pi) \mid S_0 = s \right] \tag{2.31}$$

Definition 2.6.4 (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 residual square return that can be obtained starting from a state, taking an action and then following policy π*

$$W_\pi(s, a) = \mathbb{E}_\pi \left[\sum_{t=0}^{\infty} (R_{t+1}^2 - \eta_\pi) \mid S_0 = s, A_0 = a \right] \tag{2.32}$$

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

$$U_\pi(s) = \int_{\mathbb{A}} \pi(s, a) W_\pi(s, a) \tag{2.33}$$

The average adjusted square value functions satisfy the following Bellman equations

Proposition 2.6.1 (Bellman Expectation Equations).

$$\eta_\pi + U_\pi(s) = \mathcal{M}_\pi(s) + T_\pi U_\pi(s) \tag{2.34}$$

$$\begin{aligned}\eta_\pi + W_\pi(s, a) &= \mathcal{M}(s, a) + T_a U_\pi(s) \\ &= \mathcal{M}(s, a) + \int_{\mathbb{S}} \mathcal{P}(s, a, s') \int_{\mathbb{A}} \pi(s', a') W_\pi(s', a') da' ds'\end{aligned}\tag{2.35}$$

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

$$\begin{cases} \max_{\pi} \rho_{\pi} \\ \text{subject to } \Lambda_{\pi} \leq \alpha \end{cases} \quad (2.36)$$

for a given $\alpha > 0$. Using the Lagrangian relaxation procedure, we can recast (2.36) to the following unconstrained problem

$$\max_{\lambda} \min_{\pi} L(\pi, \lambda) = -\rho_{\pi} + \lambda(\Lambda_{\pi} - \alpha) \quad (2.37)$$

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

$$\mathbf{Sh}_{\pi} = \frac{\rho_{\pi}}{\sqrt{\Lambda_{\pi}}} \quad (2.38)$$

The algorithms that will be discussed in the next chapters apply equally to both optimization problems with only minor modifications.

2.7 Policy Evaluation and Policy Iteration

TODO

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 [72], [74] or [82].

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 2.7, 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

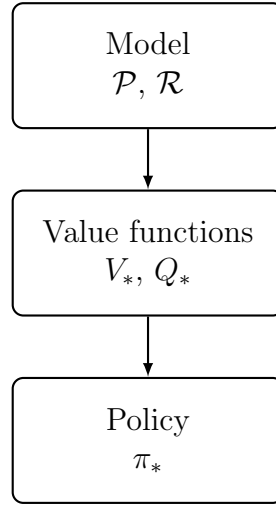


Figure 3.1: Solution process for the control problem.

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 π_* 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 [82].

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 γ 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_π or the action-value function Q_π from samples of the MDP obtained by following the same policy π to be evaluated. Although the optimal policy π_* 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 [73], [41]. 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 [60]. 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 [42] 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 θ_0 , learning rate $\{\alpha_k\}$

Output: Approximation of the optimal policy π_*

- 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 π_θ*

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

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 π_θ*

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

where d^θ is the stationary distribution of the Markov chain induced by π_θ .

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 π_θ*

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

where d^θ is the stationary distribution of the Markov chain induced by π_θ .

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 π_θ . 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 [77]

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 π_θ*

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

In a continuing environment, we might consider the long-run variance [61] defined in Section 2.6

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

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 [48] and the Sharpe ratio [70].

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 ϵ 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_θ and assume that p_θ 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 computational finance when evaluating the greeks of a derivative product [57]. Under some regularity assumptions on the function f and on the probability distribution p_θ , 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

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

which 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_{θ} , 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 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 π_{θ} . Let $G(h)$ denote the expected return obtained on trajectory h

$$G(h) = \mathbb{E}[G_0 | H_t = h] = \sum_{t=1}^{\infty} \gamma^k \mathcal{R}(s_{t-1}, a_{t-1}) \quad (4.16)$$

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

$$J_{\text{avV}}(\theta) = \mathbb{E}_{H \sim p_{\theta}} [G(H)] \quad (4.17)$$

Thanks to the likelihood ratio trick, we obtain

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

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 θ is the policy term,

$$\nabla_{\theta} \log p_{\theta}(H) = \sum_{t=0}^{\infty} \nabla_{\theta} \log \pi_{\theta}(s_t, a_t) \quad (4.19)$$

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

$$\int_{\mathbb{H}} \nabla_{\theta} p_{\theta}(h) dh = 0$$

a constant baseline $b \in \mathbb{R}$ can always be added in the gradient formula

$$\nabla_{\theta} J(\theta) = \mathbb{E} [\nabla_{\theta} \log p_{\theta}(H)(G(H) - b)] \quad (4.20)$$

We will see how this baseline can be chosen in order to minimize the variance of the estimator. In an episodic environment, we can derive an estimate of the objective function gradient by sampling M trajectories $h^{(m)} = \{(s_t^{(m)}, a_t^{(m)})\}_{t=0}^{T^{(m)}}$ from the MDP and by approximating the expected value via Monte Carlo

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

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

$$\hat{g}_{\text{SRF}} = \left[\sum_{i=0}^T \nabla_{\theta} \log \pi_{\theta}(s_i, a_i) \right] \left[\sum_{j=0}^T \gamma^j r_{j+1} - b \right] \quad (4.22)$$

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 benchmark to reduce the estimate variance and will be discussed in detail in the next section.

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

$$b_k^* = \arg \min_b \text{Var}(\hat{g}_k) \quad (4.23)$$

Algorithm 4.2 Episodic REINFORCE policy gradient estimate**Input:** Policy parameterization θ , 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 π_{θ}
- 2: For all k , Compute the optimal baseline b_k^* according to Eq. 4.25
- 3: Compute \hat{g}_{RF} according to Eq. (4.21)

It is easy to show that

$$b_k^* = \frac{\mathbb{E} [G(H) (\partial_{\theta_k} \log p_{\theta}(H))^2]}{\mathbb{E} [(\partial_{\theta_k} \log p_{\theta}(H))^2]} \quad (4.24)$$

which can be approximated by

$$\hat{b}_k^* = \frac{\sum_{m=1}^M \left[\sum_{i=0}^{T^{(m)}} \partial_{\theta_k} \log \pi_{\theta} \left(s_i^{(m)}, a_i^{(m)} \right) \right]^2 \sum_{j=0}^{T^{(m)}} \gamma^j r_{j+1}^{(m)}}{\sum_{m=1}^M \left[\sum_{i=0}^{T^{(m)}} \partial_{\theta_k} \log \pi_{\theta} \left(s_i^{(m)}, a_i^{(m)} \right) \right]^2} \quad (4.25)$$

4.3.2 GPOMDP

In the last section, we said that the REINFORCE gradient estimate is characterized by a large variance, which may slow the method's convergence. To improve the estimate, it is sufficient to notice that future actions do not depend on past rewards, unless the policy has been changed. Therefore,

$$\mathbb{E}_{\pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta}(S_t, A_t) \mathcal{R}(S_s, A_s)] = 0 \quad \forall t > s \quad (4.26)$$

From this trivial remark, we can derive two estimates of the objective function gradient

$$g_{PG} = \frac{1}{M} \sum_{m=1}^M \sum_{i=0}^{T^{(m)}} \nabla_{\theta} \log \pi_{\theta}(s_i^{(m)}, a_i^{(m)}) \left(\sum_{j=i}^{T^{(m)}} \gamma^j r_{j+1}^{(m)} - b \right) \quad (4.27)$$

$$g_{GPOMDP} = \frac{1}{M} \sum_{m=1}^M \sum_{j=0}^{T^{(m)}} \left[\sum_{i=j}^{T^{(m)}} \nabla_{\theta} \log \pi_{\theta}(s_i^{(m)}, a_i^{(m)}) \right] \left(\gamma^j r_{j+1}^{(m)} - b \right) \quad (4.28)$$

These two estimates are exactly equivalent. This simple trick greatly reduces the estimate variance and this can speed up convergence.

4.3.3 Risk-Sensitive Monte Carlo Policy Gradient

Cf. Tamar

4.3.4 Stochastic Policies

However, this trick only works if the policy is differentiable with respect to θ and it is stochastic. In most cases this is not a big problem, since stochastic policies are needed anyway to ensure sufficient exploration. TODO

Add description of stochastic policies

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 , this policy selects an action a with probability

$$\pi_{\theta}(s, a) = \frac{e^{\theta^T \Phi(s, a)}}{\sum_{b \in \mathbb{A}} e^{\theta^T \Phi(s, b)}} \quad (4.29)$$

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

$$\nabla_{\theta} \log \pi_{\theta}(s, a) = \Phi(s, a) - \sum_{b \in \mathbb{A}} \pi_{\theta}(s, b) \Phi(s, b) \quad (4.30)$$

4.3.4.2 Gaussian Exploration Policy

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

$$\nabla_{\psi} \log \pi_{\theta}(s, a) = \left(\frac{\partial \mu_{\psi}(s)}{\partial \psi} \right)^T \Sigma^{-1} (a - \mu_{\psi}(s)) \quad (4.31)$$

$$\nabla_{\Sigma} \log \pi_{\theta}(s, a) = \frac{1}{2} \left[\Sigma^{-1} (a - \mu_{\psi}(s)) (a - \mu_{\psi}(s))^T \Sigma^{-1} - \Sigma^{-1} \right] \quad (4.32)$$

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

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 π_{θ} 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 [68] 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 policy $F : \mathbb{S} \times \Theta \rightarrow \mathbb{A}$ 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 θ are random variables distributed according to a probability distribution parametrized by some hyperparameters $\xi \in \Xi \subseteq \mathbb{R}^{D_\xi}$, i.e. $\theta \sim p(\theta; \xi) = p_\xi(\theta)$. Combining these two hypotheses, we obtain a stochastic policy parametrized by the hyperparameters ξ

$$\pi_\xi(s, a) = \pi(s, a; \xi) = \int_{\Theta} p_\xi(\theta) \mathbb{1}_{F_\theta(s)=a} d\theta \quad (4.33)$$

The advantage of this approach is that the policy 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 θ once at the beginning of the episode and then generate an entire trajectory following the deterministic policy F_θ . 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 ξ will be updated following the gradient of the expected reward, which can be rewritten as

$$J(\xi) = \int_{\Theta} \int_{\mathbb{H}} p_\xi(\theta, h) G(h) dh d\theta \quad (4.34)$$

Hence, using the fact that h is conditionally independent from ξ given θ , so that $p_\xi(\theta, h) = p_\xi(\theta)p_\theta(h)$, and the likelihood trick

$$\begin{aligned} \nabla_\xi J(\xi) &= \int_{\Theta} \int_{\mathbb{H}} \nabla_\xi p_\xi(\theta) p_\theta(h) G(h) dh d\theta \\ &= \int_{\Theta} \int_{\mathbb{H}} p_\xi(\theta) p_\theta(h) \nabla_\xi \log p_\xi(\theta) G(h) dh d\theta \\ &= \mathbb{E} [\nabla_\xi \log p_\xi(\theta) G(H)] \end{aligned} \quad (4.35)$$

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

$$\nabla_{\xi} J(\xi) = \mathbb{E} [\nabla_{\xi} \log p_{\xi}(\theta) (G(H) - b)] \quad (4.36)$$

the hyperparameters can then be updated by gradient ascent

$$\xi_{k+1} = \xi_k + \alpha_k \nabla_{\xi} J(\xi) \quad (4.37)$$

The gradient can be approximated via Monte Carlo by sampling $\theta \sim p_{\xi}$ and then generating M trajectories $h^{(m)} = \{(s_t^{(m)}, a_t^{(m)})\}_{t \geq 0}$ following the deterministic policy F_{θ}

$$g_{\text{PGPE}} = \frac{1}{M} \sum_{m=1}^M \nabla_{\xi} \log p_{\xi}(\theta) [G(h^{(m)}) - b] \quad (4.38)$$

Alternatively, we can use a fully stochastic approximation by sampling a single trajectory and approximating the gradient as

$$g_{\text{PGPE}} = \nabla_{\xi} \log p_{\xi}(\theta) [G(h) - b] \quad (4.39)$$

In the following paragraphs, we discuss some possible choices for the policy parameters distribution p_{ξ} , which is the last component of the algorithm.

Gaussian Parameter Distribution A simple approach is to assume that all the components of the parameter vector θ are independent and normally distributed with mean μ_i and variance σ_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

$$\begin{aligned} \frac{\partial \log p_{\xi}(\theta)}{\partial \mu_i} &= \frac{\theta_i - \mu_i}{\sigma_i^2} \\ \frac{\partial \log p_{\xi}(\theta)}{\partial \sigma_i} &= \frac{(\theta_i - \mu_i)^2 - \sigma_i^2}{\sigma_i^3} \end{aligned} \quad (4.40)$$

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

$$\begin{aligned} \mu_i^{k+1} &= \mu_i^k + \alpha [G(h) - b] (\theta_i - \mu_i) \\ \sigma_i^{k+1} &= \sigma_i^k + \alpha [G(h) - b] \frac{(\theta_i - \mu_i)^2}{\sigma_i} \end{aligned} \quad (4.41)$$

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

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 θ^+ (resp. θ^-), the objective function gradient can be approximated with

$$\begin{aligned}\nabla_{\mu_i} J(\xi) &\approx \frac{\epsilon_i(G^+ - G^-)}{2\sigma_i^2} \\ \nabla_{\sigma_i} J(\xi) &\approx \frac{\epsilon_i^2 - \sigma_i^2}{\sigma_i^3} \left(\frac{G^+ + G^-}{2} - b \right)\end{aligned}\tag{4.42}$$

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

$$\begin{aligned}\mu_i^{k+1} &= \mu_i^k + \alpha\epsilon_i(G^+ - G^-) \\ \sigma_i^{k+1} &= \sigma_i^k + \alpha \frac{\epsilon_i^2 - \sigma_i^2}{\sigma_i} (G^+ + G^- - 2b)\end{aligned}\tag{4.43}$$

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

$$\begin{aligned}\mu_i^{k+1} &= \mu_i^k + \alpha\epsilon_i \frac{(G^+ - G^-)}{2m - G^+ - G^-} \\ \sigma_i^{k+1} &= \sigma_i^k + \alpha \frac{\epsilon_i^2 - \sigma_i^2}{\sigma_i} \frac{(G^+ + G^- - 2b)}{m - b}\end{aligned}\tag{4.44}$$

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

In this section we discuss the extension of the PGPE method to the infinite horizon case [67]. While in the episodic PGPE the parameters θ are sampled only at the beginning of each episode, in *Infinite Horizon PGPE* (IHPGPE) 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 θ_t can be sampled consecutively as follows

$$p_\xi(\theta_{i,t+1}) = \varepsilon \mathcal{N}(\mu_{i,t}, \sigma_{i,t}^2) + (1 - \varepsilon) \delta_{\theta_{i,t}} \quad (4.45)$$

In practice, ε 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

$$p_\xi(\theta_{t+1}) = \varepsilon \mathcal{N}(\mu_t, \Sigma_t) + (1 - \varepsilon) \delta_{\theta_t} \quad (4.46)$$

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

$$\begin{aligned} \mu_{i,t+1} &= \mu_{i,t} + \alpha [G_t(h) - b] (\theta_{i,t} - \mu_{i,t}) \\ \sigma_{i,t+1} &= \sigma_{i,t} + \alpha [G_t(h) - b] \frac{(\theta_{i,t} - \mu_{i,t})^2}{\sigma_{i,t}} \end{aligned} \quad (4.47)$$

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

4.4 Policy Gradient Theorem

4.4.1 Risk-Neutral Setting

These algorithms can be derived in an alternative way from a more general result: the policy gradient theorem [73]

Theorem 4.4.1 (Policy Gradient). *For any differentiable policy π_θ , for any of the policy objective functions $J = J_{start}, J_{avV}, J_{avR}, \frac{1}{1-\gamma} J_{avV}$, the policy gradient is*

$$\nabla_\theta J(\theta) = \mathbb{E}_{\substack{S \sim d^\theta \\ A \sim \pi_\theta}} [\nabla_\theta \log \pi_\theta(S, A) Q_{\pi_\theta}(S, A)] \quad (4.48)$$

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_θ

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

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^θ and noting that, because of stationarity,

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

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. Using the Bellman expectation equation for Q_θ ,

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

Hence,

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

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 π_θ . Defining the γ -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

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

□

The action-value function is typically unknown and needs to be approximated. For instance, REINFORCE and GPOMDP replace it with the realized return achieved on a sample trajectory, that are an unbiased estimate of the action-value function. However, the theorem can be used as the starting point to derive many other policy gradient methods that use different approximation of the action-value function.

4.4.2 Risk-Sensitive Setting

4.4.3 GPOMDP from the Policy Gradient Theorem

4.4.4 Compatible Function Approximation

4.4.5 Actor-Critic Policy Gradient

Actor-Critic Policy Gradient employs a *critic*, that is a parametric approximation $\hat{Q} : \mathbb{S} \times \mathbb{A} \times \Psi \rightarrow \mathbb{R}$, where $\Psi \in \mathbb{R}^{D_{\psi}}$ such that $\hat{Q}_{\psi}(s, a) = \hat{Q}(s, a; \psi) \approx Q_{\pi_{\theta}}(s, a)$. Therefore these methods maintain two sets of parameters: a *critic* that updates the action-value function parameters ψ and an *actor* that updates the policy parameters θ in the direction suggested by the critic. More formally, given \hat{Q}_{ψ} , the current state of the system s_t and the action a_t selected using policy π_{θ} , the policy parameters are updated in the approximated gradient direction

$$\theta_{k+1} = \theta_k + \alpha_k \nabla_{\theta} \pi_{\theta}(s_t, a_t) \hat{Q}_{\psi}(s_t, a_t) \quad (4.49)$$

On the other hand, the action-value function parameters ψ can be updated using any value-based approach, for instance TD(0). This leads to the QAC algorithm, for which the pseudo code is reported in ??.

[TODO: TD policy gradient, compatible function approximation, advantage function and actor-critic policy gradient]

4.4.5.1 QAC

4.4.5.2 ARRSAC

4.4.6 Natural Policy Gradient

4.5 Risk-Sensitive Policy Gradient

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

$$\max_{\lambda} \min_{\theta} L(\theta, \lambda) = -\rho(\theta) + \lambda(\Lambda(\theta) - \alpha) \quad (4.50)$$

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

$$\nabla_{\theta} L(\theta, \lambda) = -\nabla_{\theta} \rho(\theta) + \lambda \nabla_{\theta} \Lambda(\theta) \quad (4.51)$$

while the Lagrange multiplier is updated following the gradient descent direction

$$\nabla_{\lambda} L(\theta, \lambda) = \Lambda(\theta) - \alpha \quad (4.52)$$

Since

$$\nabla_{\theta} \Lambda(\theta) = \nabla_{\theta} \eta(\theta) - 2\rho(\theta) \nabla_{\theta} \rho(\theta) \quad (4.53)$$

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).

$$\begin{aligned} \nabla_{\theta} \rho(\theta) &= \mathbb{E}_{\substack{S \sim d_{\pi} \\ A \sim \pi}} [\nabla_{\theta} \log \pi_{\theta}(S, A) Q_{\theta}(S, A)] \\ &= \int_{\mathbb{S}} d_{\pi}(s) \int_{\mathbb{A}} \pi(s, a) \nabla_{\theta} \log \pi_{\theta}(s, a) Q_{\theta}(s, a) da ds \end{aligned} \quad (4.54)$$

$$\begin{aligned} \nabla_{\theta} \eta(\theta) &= \mathbb{E}_{\substack{S \sim d_{\pi} \\ A \sim \pi}} [\nabla_{\theta} \log \pi_{\theta}(S, A) W_{\theta}(S, A)] \\ &= \int_{\mathbb{S}} d_{\pi}(s) \int_{\mathbb{A}} \pi(s, a) \nabla_{\theta} \log \pi_{\theta}(s, a) W_{\theta}(s, a) da ds \end{aligned} \quad (4.55)$$

Proof. Eq. (4.54) is the standard policy gradient and its proof can be found in the literature, e.g. [73]. Eq. (4.55) 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

$$A_\theta(s, a) = Q_\theta(s, a) - V_\theta(s) \quad (4.56)$$

$$B_\theta(s, a) = W_\theta(s, a) - U_\theta(s) \quad (4.57)$$

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

$$\nabla_\theta \rho(\theta) = \mathbb{E}_{\substack{S \sim d_\pi \\ A \sim \pi}} [\nabla_\theta \log \pi_\theta(S, A) A_\theta(S, A)] \quad (4.58)$$

$$\nabla_\theta \eta(\theta) = \mathbb{E}_{\substack{S \sim d_\pi \\ A \sim \pi}} [\nabla_\theta \log \pi_\theta(S, A) B_\theta(S, A)] \quad (4.59)$$

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 [61], starting from (4.58) and (4.59), 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 δ_n^A and δ_n^B be the TD errors for residual value and square value functions

$$\begin{aligned} \delta_t^A &= R_{t+1} - \hat{\rho}_{t+1} + \hat{V}(S_{t+1}) - \hat{V}(S_t) \\ \delta_t^B &= R_{t+1}^2 - \hat{\eta}_{t+1} + \hat{U}(S_{t+1}) - \hat{U}(S_t) \end{aligned} \quad (4.60)$$

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

Proposition 4.5.1 (Temporal Difference Errors).

$$\begin{aligned} \mathbb{E}_\theta [\delta_t^A | S_t = s, A_t = a] &= A_\theta(s, a) \\ \mathbb{E}_\theta [\delta_t^B | S_t = s, A_t = a] &= B_\theta(s, a) \end{aligned} \quad (4.61)$$

Proof. TODO □

Denoting by $\psi_t = \psi(S_t, A_t) = \nabla_{\theta} \log \pi(S_t, A_t)$ the compatible feature [73], we can easily obtain an unbiased estimate of the gradients

$$\begin{aligned}\nabla_{\theta} \rho(\theta) &\approx \psi_t \delta_t^A \\ \nabla_{\theta} \eta(\theta) &\approx \psi_t \delta_t^B\end{aligned}\tag{4.62}$$

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

$$\begin{aligned}\widehat{V}(s) &= v^T \Phi_V(s) \\ \widehat{U}(s) &= u^T \Phi_U(s)\end{aligned}\tag{4.63}$$

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

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.

5.1 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 P&L and the transaction costs paid to the broker. More formally

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

where PL_{t+1} denotes the profit due to the variation of the portfolio asset prices between t and $t + 1$

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

and TC_t denotes the fees paid to the broker to change the portfolio allocation and on the short positions

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

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 States

In this section, we specify the state of the system that the agent observes in order 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. Summing up, we assume that the state of the system is composed in the following way

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

5.3 Actions

In this section, we specify the actions that the trading system may take. We assume that the agent 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, which has been introduced in Section ?? . 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 be employed.

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] 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.
- [7] 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.
- [8] —, *Q-learning-based financial trading systems with applications*, University Ca’Foscari of Venice, Dept. of Economics Working Paper Series No, 15 (2014).
- [9] —, *Reinforcement learning for automated financial trading: Basics and applications*, in Recent Advances of Neural Network Models and Applications, Springer, 2014, pp. 197–213.

-
- [10] D. P. BERTSEKAS, *Dynamic programming and optimal control*, vol. 1, Athena Scientific, Belmont, 1995.
 - [11] D. P. BERTSEKAS AND S. E. SHREVE, *Stochastic optimal control: the discrete-time case*, vol. 23, Academic Press New York, 1978.
 - [12] D. P. BERTSEKAS AND J. N. TSITSIKLIS, *Neuro-Dynamic Programming*, Optimization and neural computation series, Athena Scientific, Belmont, 1 ed., 1996.
 - [13] S. BHATNAGAR, R. S. SUTTON, M. GHAVAMZADEH, AND M. LEE, *Natural actor-critic algorithms*, Automatica, 45 (2009), pp. 2471–2482.
 - [14] C. M. BISHOP, *Pattern Recognition and Machine Learning*, Springer, 2006.
 - [15] V. S. BORKAR, *A sensitivity formula for risk-sensitive cost and the actor-critic algorithm*, Systems & Control Letters, 44 (2001), pp. 339–346.
 - [16] —, *Q-learning for risk-sensitive control*, Mathematics of operations research, 27 (2002), pp. 294–311.
 - [17] 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.
 - [18] L. BUSONI, R. BABUSKA, B. DE SCHUTTER, AND D. ERNST, *Reinforcement learning and dynamic programming using function approximators*, vol. 39, CRC press, 2010.
 - [19] P. X. CASQUEIRO AND A. J. RODRIGUES, *Neuro-dynamic trading methods*, European Journal of Operational Research, 175 (2006), pp. 1400–1412.
 - [20] 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.
 - [21] M. CHOEY AND A. S. WEIGEND, *Nonlinear trading models through sharpe ratio maximization*, International Journal of Neural Systems, 8 (1997), pp. 417–431.
 - [22] 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.

- [23] M. CORAZZA AND A. SANGALLI, *Q-learning vs. sarsa: comparing two intelligent stochastic control approaches for financial trading*.
- [24] J. CUMMING, D. ALRAJEH, AND L. DICKENS, *An investigation into the use of reinforcement learning techniques within the algorithmic trading domain*, (2015).
- [25] M. A. DEMPSTER AND V. LEEMANS, *An automated fx trading system using adaptive reinforcement learning*, Expert Systems with Applications, 30 (2006), pp. 543–552.
- [26] M. A. H. DEMPSTER AND Y. S. ROMAHI, *Intraday FX trading: An evolutionary reinforcement learning approach*, Springer, 2002.
- [27] 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).
- [28] 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.
- [29] X. DU, J. ZHAI, AND K. LV, *Algorithm trading using q-learning and recurrent reinforcement learning*, positions, 1, p. 1.
- [30] T. ELDER, *Creating algorithmic traders with hierarchical reinforcement learning*, (2008).
- [31] 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.
- [32] 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.
- [33] C. GOLD, *Fx trading via recurrent reinforcement learning*, in IEEE 2003 IEEE International Conference on Computational Intelligence for Financial Engineering. Proceedings, 2003.
- [34] I. GOODFELLOW, Y. BENGIO, AND A. COURVILLE, *Deep learning*. Book in preparation for MIT Press, 2016.
- [35] T. HASTIE, R. TIBSHIRANI, AND J. FRIEDMAN, *The elements of statistical learning: data mining, inference, and prediction*, Springer, 2009.

- [36] 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.
- [37] H. JAEGER, *Tutorial on training recurrent neural networks, covering BPPT, RTRL, EKF and the echo state network approach*, GMD-Forschungszentrum Informationstechnik, 2002.
- [38] M. S. JOSHI, *C++ design patterns and derivatives pricing*, vol. 2, Cambridge University Press, 2008.
- [39] 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.
- [40] 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.
- [41] V. R. KONDA AND J. N. TSITSIKLIS, *Actor-critic algorithms.*, in NIPS, vol. 13, 1999, pp. 1008–1014.
- [42] H. KUSHNER AND G. G. YIN, *Stochastic approximation and recursive algorithms and applications*, vol. 35, Springer Science & Business Media, 2003.
- [43] 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.
- [44] 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.
- [45] —, *Optimal posting price of limit orders: learning by trading*, Mathematics and Financial Economics, 7 (2013), pp. 359–403.
- [46] 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.

- [47] S. MAHADEVAN, *Average reward reinforcement learning: Foundations, algorithms, and empirical results*, Machine learning, 22 (1996), pp. 159–195.
- [48] H. MARKOWITZ, *Portfolio selection*, The journal of finance, 7 (1952), pp. 77–91.
- [49] 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.
- [50] J. MOODY AND M. SAFFELL, *Learning to trade via direct reinforcement*, Neural Networks, IEEE Transactions on, 12 (2001), pp. 875–889.
- [51] 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.
- [52] 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.
- [53] 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.
- [54] 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.
- [55] J. NOCEDAL AND S. WRIGHT, *Numerical optimization*, Springer Science & Business Media, 2006.
- [56] 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.
- [57] G. PAGES, *Introduction to Numerical Probability for Finance*, LPMA-Université Pierre et Marie Curie, 2016.
- [58] J. PETERS, K. MÜLLING, AND Y. ALTUN, *Relative entropy policy search.*, in AAAI, Atlanta, 2010.

- [59] 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.
- [60] —, *Reinforcement learning of motor skills with policy gradients*, Neural networks, 21 (2008), pp. 682–697.
- [61] L. PRASHANTH AND M. GHAVAMZADEH, *Actor-critic algorithms for risk-sensitive reinforcement learning.*, arXiv preprint arXiv:1403.6530, (2014).
- [62] M. L. PUTERMAN, *Markov decision processes: discrete stochastic dynamic programming*, John Wiley & Sons, 1994.
- [63] 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.
- [64] C. SANDERSON, *Armadillo: An open source c++ linear algebra library for fast prototyping and computationally intensive experiments*, (2010).
- [65] 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.
- [66] —, *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.
- [67] F. SEHNKE ET AL., *Parameter exploring policy gradients and their implications*, PhD thesis, Technische Universität München, 2012.
- [68] 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.
- [69] —, *Parameter-exploring policy gradients*, Neural Networks, 23 (2010), pp. 551–559.
- [70] W. F. SHARPE, *The sharpe ratio*, The journal of portfolio management, 21 (1994), pp. 49–58.

- [71] D. SILVER, G. LEVER, N. HEESS, T. DEGRIS, D. WIERSTRA, AND M. RIEDMILLER, *Deterministic policy gradient algorithms*, in ICML, 2014.
- [72] R. S. SUTTON AND A. G. BARTO, *Introduction to reinforcement learning*, vol. 135, MIT Press Cambridge, 1998.
- [73] 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.
- [74] C. SZEPESVÁRI, *Algorithms for reinforcement learning*, Synthesis lectures on artificial intelligence and machine learning, 4 (2010), pp. 1–103.
- [75] 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.
- [76] 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.
- [77] 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.
- [78] A. TAMAR AND S. MANNOR, *Variance adjusted actor critic algorithms*, arXiv preprint arXiv:1310.3697, (2013).
- [79] 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.
- [80] R. S. TSAY, *Analysis of financial time series*, vol. 543, John Wiley & Sons, 2005.
- [81] P. J. WERBOS, *Backpropagation through time: what it does and how to do it*, Proceedings of the IEEE, 78 (1990), pp. 1550–1560.
- [82] M. WIERING AND M. VAN OTTERLO, *Reinforcement Learning: State-of-the-Art*, vol. 12 of Adaptation, Learning, and Optimization, Springer, 1 ed., 2012.

-
- [83] R. J. WILLIAMS AND D. ZIPSER, *A learning algorithm for continually running fully recurrent neural networks*, Neural computation, 1 (1989), pp. 270–280.
 - [84] 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.
 - [85] 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.
 - [86] 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.