

POLITECNICO DI MILANO

Facoltà di Ingegneria

Scuola di Ingegneria Industriale e dell'Informazione

Dipartimento di Elettronica, Informazione e Bioingegneria

Master of Science in

Computer Science and Engineering – Ingegneria Informatica



Adaptive Batch Size for Policy Gradient Methods

Supervisor:

MARCELLO RESTELLI

Assistant Supervisor:

MATTEO PIROTTA

Master Graduation Thesis by:

MATTEO PAPINI

Student Id n. 850421

Academic Year 2016-2017

ACKNOWLEDGMENTS

Todo

CONTENTS

Abstract	vii
1 INTRODUCTION	1
2 PRELIMINARIES ON POLICY GRADIENT	3
2.1 Continuous Markov Decision Processes	3
2.1.1 The model	3
2.1.2 Problem formulation	4
2.1.3 Value functions	5
2.2 The Policy Gradient Approach	5
2.2.1 Definition	5
2.2.2 Stochastic gradient descent	6
2.2.3 Policy Gradient Theorem and eligibility vector	6
2.2.4 Common policy classes	7
2.2.5 A general policy gradient algorithm	7
3 POLICY GRADIENT METHODS: STATE OF THE ART	9
3.1 Why Policy Gradient	9
3.2 Likelihood Ratio Methods	10
3.2.1 From finite differences to likelihood ratio	10
3.2.2 Baselines	11
3.2.3 The REINFORCE algorithm	11
3.2.4 The PGT/G(PO)MDP algorithm	12
3.3 Beyond likelihood ratio	13
3.3.1 Natural gradients	13
3.3.2 Actor-critic methods	13
3.3.3 Expectation-maximization	13
3.4 Safe policy gradient methods	13
3.4.1 The need for safe policy gradients	13
3.4.2 The step size problem	13
3.4.3 State of the art	13
3.4.4 What to do next	13
4 JOINT STEP-SIZE AND BATCH-SIZE OPTIMIZATION	15
5 NUMERICAL SIMULATIONS	17
6 CONCLUSION	19
BIBLIOGRAPHY	21

LIST OF FIGURES

LIST OF TABLES

LISTINGS

ACRONYMS

MDP	Markov Decision Process
RL	Reinforcement Learning

ABSTRACT

Todo

SOMMARIO

Todo

INTRODUCTION

This is an introduction.

PRELIMINARIES ON POLICY GRADIENT

In this chapter we provide basic knowledge on the reinforcement learning framework and on the policy gradient approach. The aim is to introduce concepts that will be used extensively in the following chapters. For a complete treatment of reinforcement learning, refer to [8].

2.1 CONTINUOUS MARKOV DECISION PROCESSES

In this section we provide a brief treatment on the model at the root of most Reinforcement Learning (RL) algorithms, including policy gradient methods: the Markov Decision Process (MDP). Given the scope of our work, we will focus on continuous MDPs.

2.1.1 *The model*

Markov decision processes model the very general situation in which an agent interacts with a stochastic, partially observable environment in order to reach a goal. At each time step t , the agent receives observations from the environment. From the current observations and previous interactions, the agent builds a representation of the current state of the environment, s_t . Then, basing its decision on s_t and any other source of knowledge, it takes action a_t . Finally, it is given a (possibly negative) reward depending on how much its behavior is compliant with the goal, in the form of a scalar signal r_t . The dynamics of the environment, as represented by the agent, must satisfy the Markov property: s_{t+1} must depend only on s_t and a_t . Under this framework, the goal of the agent can be restated as the problem of collecting as much reward as possible. For now, and every time it is not specified otherwise, we will refer to continuing (i.e. never-ending) tasks, in which future rewards are exponentially discounted.

Formally, a discrete-time continuous Markov decision process is defined as a tuple $\langle \mathcal{S}, \mathcal{A}, \mathcal{P}, \mathcal{R}, \gamma, \mu \rangle$, where:

- $\mathcal{S} \subseteq \mathbb{R}^n$ is an n -dimensional continuous state space. Elements of this set are the possible states of the environment as observed by the agent.
- $\mathcal{A} \subseteq \mathbb{R}^m$ is an m -dimensional continuous action space, from which the agent can draw its actions.

- $\mathcal{P}: \mathcal{S} \times \mathcal{A} \mapsto \Delta(\mathcal{S})$ is a Markovian transition model, where $\mathcal{P}(s' | s, a)$ defines the transition density between states s and s' under action a , i.e. the probability distribution over the next state s' when the current state is s and the taken action is a .
- $\mathcal{R}: \mathcal{S} \times \mathcal{A} \rightarrow [-R, R]$ is the reward function, such that $\mathcal{R}(s, a)$ is the expected value of the reward received by taking action a in state s and $R \in \mathbb{R}^+$ is the maximum absolute-value reward.
- $\gamma \in [0, 1)$ is the discount factor for future rewards.
- $\mu \in \Delta(\mathcal{S})$ is the initial state distribution, from which the initial state of the environment is drawn.

The behavior of the agent can be defined as a stationary, possibly stochastic policy:

$$\pi: \mathcal{S} \mapsto \Delta(\mathcal{A}),$$

such that $\pi(s)$ is the probability distribution over the action a to take in state s .

2.1.2 Problem formulation

The goal of [RL](#) is to find an optimal policy π^* by exploiting the information obtained through the interaction with the environment. The objective is to maximize the total discounted reward, called return v :

$$v = \sum_{t=0}^{\infty} \gamma^t r_t$$

To evaluate how good a policy π is, we need to compute the expected value of the return over π and the initial state distribution μ :

$$J_{\mu}^{\pi} = \mathbf{E}_{\mu, \pi}[v].$$

It is convenient to introduce a new distribution, called stationary distribution or future-state distribution d_{μ}^{π} :

$$d_{\mu}^{\pi}(s) = (1 - \gamma) \sum_{t=0}^{\infty} \gamma^t \Pr(s_t = s | \pi, \mu),$$

which allows to define the expected return J as:

$$J_{\mu}^{\pi} = \int_{\mathcal{S}} d_{\mu}^{\pi}(s) \int_{\mathcal{A}} \pi(a | s) \mathcal{R}(s, a) da ds.$$

By introducing the joint distribution ζ between d_{μ}^{π} and π , we can express the expected return even more compactly:

$$J_{\mu}^{\pi} = \mathbf{E}_{(s, a) \sim \zeta} [\mathcal{R}(s, a)]$$

The [RL](#) problem can then be formulated as a stochastic optimization problem:

$$\pi^* \in \arg \max_{\pi} J(\pi).$$

2.1.3 Value functions

Value functions provide a measure of how good a state, or a state-action pair, is under a policy π . They are powerful theoretical tools and are employed in many practical RL algorithms. The state value function $V^\pi: \mathcal{S} \mapsto \mathbb{R}$ assigns to each state the expected return that is obtained by following π starting from state s , and can be defined recursively by the following equation:

$$V^\pi(s) = \int_{\mathcal{A}} \pi(a|s) \left(\mathcal{R}(s, a) + \gamma \int_{\mathcal{S}} \mathcal{P}(s'|s, a) V^\pi(s') ds' \right) da,$$

which is known as Bellman's expectation equation. The expected return $J_\mu(\pi)$ can be expressed in terms of the state-value function as:

$$J_\mu^\pi = \int_{\mathcal{S}} \mu(s) V^\pi(s) ds.$$

For control purposes, it is more useful to define an action-value function $Q^\pi: \mathcal{S} \times \mathcal{A} \mapsto \mathbb{R}$, such that $Q^\pi(s, a)$ is the return obtained starting from state s , taking action a and following policy π from then on. Also the action-value function is defined by a Bellman equation:

$$Q^\pi(s, a) = \mathcal{R}(s, a) + \gamma \int_{\mathcal{S}} \mathcal{P}(s'|s, a) \int_{\mathcal{A}} \pi(a'|s') Q^\pi(s', a') da' ds'.$$

Finally, the difference between the two value functions is known as advantage function:

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

Intuitively, $A(s, a)$ represents the advantage of taking action a in state s instead of drawing an action from the policy.

2.2 THE POLICY GRADIENT APPROACH

In this section we describe in general terms a special class of RL algorithms known as policy gradient methods, and we report some theoretical results that are of key importance for any treatment of such methods.

2.2.1 Definition

Direct policy search is an approach to RL consisting in exploring a subset of policy space directly to find an optimal policy. It is opposed to the value function approach, in which the optimal policy is derived from an estimate of the action value function Q^π . The advantages of direct policy search over value function methods for continuous MDPs are discussed in the following chapter. For a recent survey on policy search methods refer to [2].

Policy gradient is a kind of policy search in which the search is restricted to a class of parameterized policies:

$$\Pi_{\theta} = \{\pi_{\theta} : \theta \in \mathbb{R}^m\},$$

where θ is the parameter vector. The parametric policy π_{θ} can have any shape, but is required to be differentiable in θ . Like all gradient descent methods, the parameter vector is updated in the direction of the gradient of a performance measure, which is guaranteed to be the direction of maximum improvement. In our case the performance measure is the expected return $J_{\mu}^{\pi_{\theta}}$, which we rename $J_{\mu}(\theta)$ for readability. The gradient of the expected return w.r.t. the parameter vector, $\nabla_{\theta} J_{\mu}(\theta)$, is called policy gradient. Starting from an arbitrary initial parametrization θ_0 , policy gradient methods performs updates of the kind:

$$\theta \leftarrow \theta + \alpha \nabla_{\theta} J_{\mu}(\theta),$$

where $\alpha \geq 0$ is the learning rate, also called step size. The update is performed at each learning iteration. Convergence to at least a local optimum is guaranteed by the gradient descent update.

2.2.2 Stochastic gradient descent

In many practical tasks it is impossible, or unfeasible, to compute the exact policy gradient $\nabla_{\theta} J_{\mu}(\theta)$, but a gradient estimate $\hat{\nabla}_{\theta} J_{\mu}(\theta)$ can be estimated from sample trajectories. With the term trajectory we mean an episode of the task starting from an initial state s_0 , drawn from μ , and stopping after a fixed number H of time steps. The gradient estimate is often averaged over a number N of such trajectories, also called a batch. The size N of the batch is called batch size. The stochastic gradient descent update is:

$$\theta \leftarrow \theta + \alpha \hat{\nabla}_{\theta} J_{\mu}(\theta),$$

where, at each learning iteration, N trajectories need to be performed. Convergence to a local optimum is still guaranteed, provided that the angular difference between $\nabla_{\theta} J_{\mu}(\theta)$ and $\hat{\nabla}_{\theta} J_{\mu}(\theta)$ is less than $\frac{\pi}{2}$.

2.2.3 Policy Gradient Theorem and eligibility vector

The Policy Gradient Theorem is a result by Sutton [9] that relates the policy gradient to the value function Q^{π} :

Theorem 2.1 (Continuous version of Theorem 1 from [9]). For any MDP

$$\nabla_{\theta} J_{\mu}(\theta) = \int_{\mathcal{S}} d_{\mu}^{\pi}(s) \int_{\mathcal{A}} \nabla_{\theta} \pi_{\theta}(a | s) Q^{\pi}(s, a) da ds.$$

The Policy Gradient Theorem is of key importance for many policy gradient algorithms.

2.2.4 Common policy classes

Although there are no restriction on parameterized policies, some specific classes are used extensively in applications and have well studied properties. One of the most popular is the Gaussian policy, which in its most general form is defined as:

$$\pi_{\theta}(a | s) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left(-\frac{1}{2} \left(\frac{a - \mu_{\theta}(s)}{\sigma_{\theta}(s)} \right)^2 \right).$$

A common form is the Gaussian policy with fixed standard deviation σ and mean linear in a feature vector $\Phi(s)$:

$$\pi_{\theta}(a | s) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left(-\frac{1}{2} \left(\frac{a - \theta^T \Phi(s)}{\sigma} \right)^2 \right),$$

having eligibility vector:

$$\nabla_{\theta} \log \pi_{\theta}(a | s) = \frac{(a - \theta^T \Phi(s)) \Phi(s)}{\sigma^2}.$$

When the action space \mathcal{A} is discrete, a common choice is the softmax policy:

$$\pi_{\theta}(a | s) = \frac{e^{\theta^T \Phi(s, a)}}{\sum_{a' \in \mathcal{A}} e^{\theta^T \Phi(s, a')}},$$

having eligibility vector:

$$\nabla_{\theta} \log \pi_{\theta}(a | s) = \Phi(s, a) - \sum_{a' \in \mathcal{A}} \pi_{\theta}(s | a') \Phi(s, a')$$

2.2.5 A general policy gradient algorithm

TODO

POLICY GRADIENT METHODS: STATE OF THE ART

In this chapter we provide an overview of the most commonly used policy gradient methods. The aim is not to provide an exhaustive survey of policy gradient algorithms, but to give a proper context to the safe policy gradient approach, which represent the main focus of our work.

3.1 WHY POLICY GRADIENT

Traditional RL methods rely on the computation (or the estimation) of the action-value function Q^π . These so-called value function methods were first proposed for solving simple MDPs characterized by small, discrete state and action spaces. They are also known as tabular methods, since the value function can be stored in a finite table. In this limited scenario, tabular algorithms have strong convergence guarantees and perform well in practice. However, they become unfeasible when the number of possible states and actions is too big or even infinite, such as in continuous MDPs. The natural extension of tabular methods is to employ the function approximation tools known from supervised learning, such as neural networks, to represent Q^π . However, this approach have many problems, especially when applied to control tasks:

- Value function approximation methods have no guarantee of convergence, not even to locally optimal policies. In some applications, divergence may even damage the system.
- Greedy policies based on approximated value functions can be highly sensitive to small perturbations in the state, being unfeasible for tasks characterized by noise, e.g. control tasks with noisy sensors.
- Given the complexity of the value function representation (think of a multi-layer perceptron), it is difficult to isolate potentially dangerous behaviors.

Policy gradient methods bypass these problems by searching directly for the optimal policy. Convergence to a local optimum is guaranteed (also for stochastic methods), uncertainty in the state have controllable effects and policies can be made safe by design. Moreover, prior domain knowledge can be exploited to design policies suited for the specific tasks.

Indeed, policy gradient methods have been successfully applied to complex control tasks. In [7] the authors were able to solve a robot standing task, where a (simulated) biped robot must keep an erect position while perturbed by external forces. In [3] Kober and Peters used policy search in combination with imitation learning to teach the Ball-in-a-Cup game to a real robotic arm. A similar approach was used in [5] for a baseball swing task. In [4] policy search was used to play simulated robot table tennis.

3.2 LIKELIHOOD RATIO METHODS

Likelihood ratio methods include a large number of widely used policy gradient algorithms.

3.2.1 From finite differences to likelihood ratio

Among the oldest policy gradient methods we find finite-difference methods, which arose in the stochastic simulation literature. The idea is to perturb the policy parameter θ many times, obtaining $\theta + \Delta\theta_1 \dots \theta + \Delta\theta_K$, and estimate for each perturbation the expected return difference $\Delta\hat{J}_i \simeq J(\theta + \Delta\theta_i) - J(\theta)$, by performing a number of sample trajectories. After collecting all these data, the policy gradient can be estimated by regression as:

$$\hat{\nabla}_{\theta} J_{\mu}(\theta) = (\Delta\Theta^T \Delta\Theta)^{-1} \Delta\Theta^T \Delta\hat{J}$$

where $\Delta\Theta^T = [\Delta\theta_1, \dots, \Delta\theta_K]^T$ and $\Delta\hat{J} = [\Delta\hat{J}_1, \dots, \Delta\hat{J}_K]^T$. This method is easy to implement and very efficient when applied to deterministic tasks or pseudo-random number simulations. In real control tasks though, perturbing the policy parametrization without incurring in instability may not be trivial and noise can have a disastrous impact on performance. Moreover, gradient estimation requires a large number of trajectories.

The likelihood ratio approach allows to estimate the gradient even from a single trajectory, without the need of perturbing the parametrization. This method is due to Williams [11], but is better understood from the perspective of the Policy Gradient Theorem [9]. By applying trivial differentiation rules to the expression of the policy gradient (see Theorem 2.1), we have:

$$\nabla_{\theta} J_{\mu}(\theta) = \int_S d_{\mu}^{\pi}(s) \int_{\mathcal{A}} \pi_{\theta}(a | s) \nabla_{\theta} \log \pi_{\theta}(a | s) Q^{\pi}(s, a) da ds.$$

This is known as the REINFORCE trick. The term

$$\nabla_{\theta} \log \pi_{\theta}(a | s) = \frac{\nabla_{\theta} \pi_{\theta}(s, a)}{\pi_{\theta}(s, a)}$$

has many names: eligibility vector, characteristic eligibility, policy score and, of course, likelihood ratio. In terms of the joint distribution ζ , the policy gradient is just:

$$\nabla_{\theta} J_{\mu}(\theta) = \mathbf{E}_{(s,a) \sim \zeta} [\nabla_{\theta} \log \pi_{\theta}(a | s) Q^{\pi}(s, a)].$$

Since sampling state-action pairs from ζ is equivalent to following policy π_{θ} , the gradient can be estimated from a single trajectory as:

$$\tilde{\nabla}_{\theta} J_{\mu}(\theta) = \langle \nabla_{\theta} \log \pi_{\theta}(a | s) Q^{\pi}(s, a) \rangle_H,$$

where $\langle \cdot \rangle^H$ denotes sample average over H time steps. Of course a more stable estimate can be obtained by averaging again over a batch of N trajectories:

$$\hat{\nabla}_{\theta} J_{\mu}(\theta) = \langle \langle \nabla_{\theta} \log \pi_{\theta}(a | s) Q^{\pi}(s, a) \rangle_H \rangle_N.$$

3.2.2 Baselines

A problem of the REINFORCE approach is the large variance of the estimate, which results in slower convergence. The Policy Gradient Theorem still holds when an arbitrary baseline $b(s)$ is subtracted from the action-value function, as shown in [11]:

$$\nabla_{\theta} J_{\mu}(\theta) = \mathbf{E}_{(s,a) \sim \zeta} [\nabla_{\theta} \log \pi_{\theta}(a | s) (Q^{\pi}(s, a) - b(s))].$$

The baseline can be chosen to reduce the variance of the gradient estimate without introducing any bias, speeding up convergence. A natural choice of baseline is the state-value function V^{π} :

$$\begin{aligned} \nabla_{\theta} J_{\mu}(\theta) &= \mathbf{E}_{(s,a) \sim \zeta} [\nabla_{\theta} \log \pi_{\theta}(a | s) (Q^{\pi}(s, a) - V^{\pi}(s))] \\ &= \mathbf{E}_{(s,a) \sim \zeta} [\nabla_{\theta} \log \pi_{\theta}(a | s) (A^{\pi}(s, a))]. \end{aligned}$$

The usage of the advantage function with the eligibility vector has an intuitive justification: to follow the policy gradient direction means to assign to the most advantageous actions the highest probability of being taken.

3.2.3 The REINFORCE algorithm

In practical tasks, the exact value function Q^{π} is not available and must be estimated. Episodic likelihood ratio algorithms estimate it from the total return at the end of the episode, assigning to each visited state-action pair a value in retrospective. The REINFORCE gradient estimator [11] uses the total return directly:

$$\hat{\nabla}_{\theta} J_{\mu}^{\text{RF}}(\theta) = \left\langle \sum_{k=1}^H \nabla_{\theta} \log \pi_{\theta}(a_k^n | s_k^n) \left(\sum_{l=1}^H \gamma^{l-1} r_l^n - b \right) \right\rangle_N.$$

A problem of the REINFORCE algorithm is the high variance of the gradient estimate, which can be reduced by using a proper baseline. A variance-minimizing baseline can be found in [6]:

$$b = \frac{\left\langle \left(\sum_{k=1}^H \nabla_{\theta} \log \pi_{\theta}(a_k^n | s_k^n) \right)^2 \sum_{l=1}^H \gamma^{l-1} r_l^n \right\rangle_N}{\left\langle \left(\sum_{k=1}^H \nabla_{\theta} \log \pi_{\theta}(a_k^n | s_k^n) \right)^2 \right\rangle_N}.$$

Note that, in the case $|\theta| > 1$, all operations among vectors are to be intended per-component or, equivalently, each gradient component $\hat{\nabla}_{\theta_i} J_{\mu}(\theta)$ can be computed separately.

3.2.4 The PGT/G(PO)MDP algorithm

A problem of the REINFORCE gradient estimator is that the value estimate for (s_k^n, a_k^n) depends on past rewards. Two refinements are the PGT gradient estimator [10]:

$$\hat{\nabla}_{\theta} J_{\mu}^{\text{PGT}}(\theta) = \left\langle \sum_{k=1}^H \nabla_{\theta} \log \pi_{\theta}(a_k^n | s_k^n) \left(\sum_{l=k}^H \gamma^{l-1} r_l^n - b \right) \right\rangle_N,$$

and the G(PO)MDP gradient estimator [1]:

$$\hat{\nabla}_{\theta} J_{\mu}^{\text{G(PO)MDP}}(\theta) = \left\langle \sum_{l=1}^H \left(\sum_{k=1}^l \nabla_{\theta} \log \pi_{\theta}(a_k^n | s_k^n) \right) (\gamma^{l-1} r_l^n - b) \right\rangle_N.$$

Although the algorithms are different, the gradient estimator that is computed is exactly the same, i.e. $\hat{\nabla}_{\theta} J_{\mu}^{\text{PGT}}(\theta) = \hat{\nabla}_{\theta} J_{\mu}^{\text{G(PO)MDP}}(\theta)$, as shown in [6]. From now on we will refer to the PGT form. In the case $b = 0$, the PGT gradient estimation has been proven to suffer from less variance than REINFORCE [12] under mild assumptions. A variance-minimizing baseline for PGT is provided in [6]. Differently from the REINFORCE case, a separate baseline b_l is used for each time step l :

$$b_l = \frac{\left\langle \left(\sum_{k=1}^l \nabla_{\theta} \log \pi_{\theta}(a_k^n | s_k^n) \right)^2 \gamma^{l-1} r_l^n \right\rangle_N}{\left\langle \left(\sum_{k=1}^l \nabla_{\theta} \log \pi_{\theta}(a_k^n | s_k^n) \right)^2 \right\rangle_N}$$

$$\hat{\nabla}_{\theta} J_{\mu}^{\text{PGT}}(\theta) = \left\langle \sum_{k=1}^H \nabla_{\theta} \log \pi_{\theta}(a_k^n | s_k^n) \left(\sum_{l=k}^H \gamma^{l-1} r_l^n - b_l \right) \right\rangle_N$$

3.3 BEYOND LIKELIHOOD RATIO

3.3.1 *Natural gradients*

3.3.2 *Actor-critic methods*

3.3.3 *Expectation-maximization*

3.4 SAFE POLICY GRADIENT METHODS

3.4.1 *The need for safe policy gradients*

3.4.2 *The step size problem*

3.4.3 *State of the art*

3.4.4 *What to do next*

JOINT STEP-SIZE AND BATCH-SIZE OPTIMIZATION

CONCLUSION

This is a conclusion.

BIBLIOGRAPHY

- [1] Jonathan Baxter and Peter L. Bartlett. “Infinite-Horizon Policy-Gradient Estimation.” In: *Journal of Artificial Intelligence Research* 15 (2001), pp. 319–350 (cit. on p. 12).
- [2] Marc Peter Deisenroth, Gerhard Neumann, Jan Peters, et al. “A Survey on Policy Search for Robotics.” In: *Foundations and Trends in Robotics* 2.1-2 (2013), pp. 1–142 (cit. on p. 5).
- [3] Jens Kober and Jan Peters. “Policy Search for Motor Primitives in Robotics.” In: *Advances in Neural Information Processing Systems* 21. Vol. 21. 2008, pp. 849–856 (cit. on p. 10).
- [4] Jan Peters, Katharina Mülling, and Yasemin Altun. “Relative Entropy Policy Search.” In: *AAAI*. AAAI Press, 2010 (cit. on p. 10).
- [5] Jan Peters and Stefan Schaal. “Natural Actor-Critic.” In: *Neuro-computing* 71.7-9 (2008), pp. 1180–1190 (cit. on p. 10).
- [6] Jan Peters and Stefan Schaal. “Reinforcement learning of motor skills with policy gradients.” In: *Neural Networks* 21.4 (2008), pp. 682–697 (cit. on p. 12).
- [7] Frank Sehnke, Christian Osendorfer, Thomas Rückstieß, Alex Graves, Jan Peters, and Jürgen Schmidhuber. “Policy Gradients with Parameter-Based Exploration for Control.” In: *ICANN (1)*. Ed. by Vera Kurková, Roman Neruda, and Jan Koutník. Vol. 5163. Lecture Notes in Computer Science. Springer, 2008, pp. 387–396. ISBN: 978-3-540-87535-2 (cit. on p. 10).
- [8] Richard S. Sutton and Andrew G. Barto. *Introduction to Reinforcement Learning*. 1st. Cambridge, MA, USA: MIT Press, 1998. ISBN: 0262193981 (cit. on p. 3).
- [9] Richard S. Sutton, David A. McAllester, Satinder P. Singh, and Yishay Mansour. “Policy Gradient Methods for Reinforcement Learning with Function Approximation.” In: *Advances in Neural Information Processing Systems 12, [NIPS Conference, Denver, Colorado, USA, November 29 - December 4, 1999]*. Ed. by Sara A. Solla, Todd K. Leen, and Klaus-Robert Müller. The MIT Press, 1999, pp. 1057–1063. ISBN: 0-262-19450-3 (cit. on pp. 6, 10).
- [10] Richard S. Sutton, David A. McAllester, Satinder P. Singh, and Yishay Mansour. “Policy Gradient Methods for Reinforcement Learning with Function Approximation.” In: *Advances in Neural Information Processing Systems 12, [NIPS Conference, Denver, Colorado, USA, November 29 - December 4, 1999]*. 1999, pp. 1057–1063 (cit. on p. 12).

- [11] Ronald J. Williams. "Simple statistical gradient-following algorithms for connectionist reinforcement learning." In: *Machine Learning* 8.3 (1992), pp. 229–256. ISSN: 1573-0565 (cit. on pp. [10](#), [11](#)).
- [12] Tingting Zhao, Hirotaka Hachiya, Gang Niu, and Masashi Sugiyama. "Analysis and Improvement of Policy Gradient Estimation." In: *NIPS*. Ed. by John Shawe-Taylor, Richard S. Zemel, Peter L. Bartlett, Fernando C. N. Pereira, and Kilian Q. Weinberger. 2011, pp. 262–270 (cit. on p. [12](#)).