



UNIVERSIDAD  
DE MÁLAGA



Hochschule für Angewandte Wissenschaften Hamburg  
*Hamburg University of Applied Sciences*

Universidad de Málaga  
Escuela de Ingenierías Industriales

---

# TESIS DOCTORAL

## NUEVOS ALGORITMOS DE APRENDIZAJE POR REFUERZO DISTRIBUTIVO Y APRENDIZAJE ENSAMBLADOR

En cooperación con HAW Hamburg, Alemania

Autor: M.Sc. Vanya Aziz  
Directores: Dr. Eligius M.T. Hendrix & Dr. Ivo Nowak  
Programa de Doctorado: Programa de Doctorado en Ingeniería Mecatrónica  
Centro: E.T.S. de Ingeniería Informática  
Málaga, Febrero 2025





UNIVERSIDAD  
DE MÁLAGA



Hochschule für Angewandte Wissenschaften Hamburg  
*Hamburg University of Applied Sciences*

University of Málaga  
School of Industrial Engineering

---

# Ph.D. Thesis

## Novel Distributional Reinforcement and Ensemble Learning Algorithms

In cooperation with HAW Hamburg, Germany

**Author:** M.Sc. Vanya Aziz  
**Supervisors:** Dr. Eligius M.T. Hendrix  
Dr. Ivo Nowak

*Málaga, January 2025*





## DECLARACIÓN DE AUTORÍA Y ORIGINALIDAD DE LA TESIS PRESENTADA PARA OBTENER EL TÍTULO DE DOCTOR

D. VANYA AZIZ

Estudiante del programa de doctorado MECATRONICA de la Universidad de Málaga, autor de la tesis, presentada para la obtención del título de doctor por la Universidad de Málaga, titulada: NOVEL DISTRIBUTIONAL REINFORCEMENT AND ENSEMBLE LEARNING ALGORITHMS


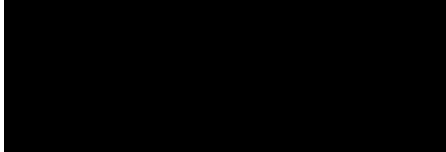

Realizada bajo la tutorización de ELIGIUS MT HENDRIX y dirección de ELIGIUS MT HENDRIX Y IVO NOWAK

DECLARO QUE:

La tesis presentada es una obra original que no infringe los derechos de propiedad intelectual ni los derechos de propiedad industrial u otros, conforme al ordenamiento jurídico vigente (Real Decreto Legislativo 1/1996, de 12 de abril, por el que se aprueba el texto refundido de la Ley de Propiedad Intelectual, regularizando, aclarando y armonizando las disposiciones legales vigentes sobre la materia), modificado por la Ley 2/2019, de 1 de marzo.

Igualmente asumo, ante a la Universidad de Málaga y ante cualquier otra instancia, la responsabilidad que pudiera derivarse en caso de plagio de contenidos en la tesis presentada, conforme al ordenamiento jurídico vigente.

En Málaga, a 25 de ENERO de 2025

 Fdo.: VANYA AZIZ Doctorando	 Fdo.: ELIGIUS MT HENDRIX Tutor
 Fdo.: ELIGIUS MT HENDRIX Y IVO NOWAK Directores de tesis	



A quien corresponda.

Por medio de la presente, Dr. Eligius M. T. Hendrix [REDACTED] Catedrático a la Universidad de Málaga y Dr. Ivo Nowak [REDACTED] Catedrático en la Universidad Aplicada de Hamburgo, como directores y tutor de la tesis titulada "Novel Distributional Reinforcement and Ensemble Learning Algorithms", enmarcada en el programa de doctorado "Ingeniería Mecatrónica" y llevada a cabo por el doctorando Vanya Aziz, dan el visto bueno al depósito y posterior defensa de dicha tesis.

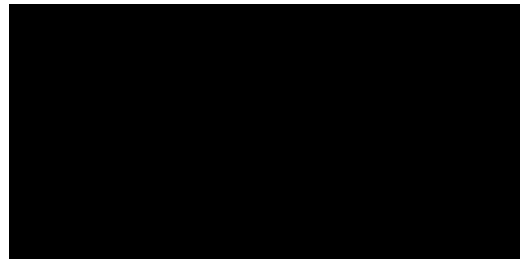
Además, los directores de tesis informan que las publicaciones que avalan la tesis presentado en la página siguiente no han sido utilizadas en tesis anteriores

Málaga a 21 de marzo de 2025.

Firmado (director y tutor de la tesis): Eligius M. T. Hendrix



Firmado (director de la tesis): Ivo Nowak





---

## Abstract

The term “Industry 4.0” describes the fourth industrial revolution and is characterized by the integration of digital technology into manufacturing processes. The transformative concepts in Industry 4.0 enable economic production at radically small lot sizes, requiring unprecedented levels of automation and adaptability. These requirements on production facilities necessitate autonomously acting and self-optimizing systems. The field of machine learning offers promising solutions to achieve the objectives of Industry 4.0, particularly by enabling data-driven decision-making and adaptive control mechanisms. This dissertation focuses on Deep Reinforcement Learning (DRL), a neural network-based approach for solving Markov Decision Processes in high-dimensional spaces with unknown transition dynamics. The main contribution of this thesis is the development of a novel state-of-the-art distributional reinforcement learning algorithm within the maximum-entropy Actor-Critic framework. This algorithm, termed “Cramér-based Soft Distributional Soft Actor-critic” (C-DSAC), demonstrates superior performance to other RL algorithms, especially in environments with high-dimensional spaces and complex dynamics. Its performance is shown to be partly rooted in a phenomenon arising in Cramér-metric-based Distributional Reinforcement Learning, referred to as confidence-driven model updates. This mechanism ensures that the value function approximator is updated more conservatively when confidence in its estimates is low. Theoretical justifications for the algorithm are provided, demonstrating its convergence in the policy evaluation setting and, under widely accepted mild assumptions, in the control setting as well. Beyond foundational algorithmic research, this thesis contributes to the practical application of RL in robotics. Given the crucial role of multi-joint robotic systems in modern production technology, a RL meta-algorithm called “Reinforcement Learning - Inverse Kinematics” (RL-IK) is devised. This approach enhances the applicability of reinforcement learning to robotic control tasks by significantly accelerating convergence to near-optimal policies

---

compared to standard RL methods. An essential prerequisite for real-world RL applications in control systems is machine perception for state identification. To address challenges in this field, this thesis explores novel Supervised Learning (SL) approaches, validated on image classification tasks, with a focus on ensemble learning strategies. To provide a deeper theoretical understanding, the methods are analysed through the lens of classical optimization and compared to heuristic approaches. This investigation offers valuable insights into the difficulties of finding globally optimal classifier ensembles.

---

## Acknowledgements

First and foremost, I thank the living God for His endless grace and blessings, which have illuminated my path through the vast state-space of life's journey. In His intricate environment, I often found myself acting as a curious agent – opting for exploration over mere exploitation, eagerly sampling every state transition, and collecting life's rewards one unexpected discovery at a time. My decisions leading to my current state have been guided by His rewards during pivotal episodes that I now wish to recount.

I am deeply grateful to my parents Avan and Aziz. By continuously feeding my desire for self-improvement, you ensured that I do not become trapped in local optima but embrace exploratory actions to obtain better policies. Through you I learned what dedication and love – be it for another soul or for higher principles – truly means.

To my aunts Nishtiman and Snawbar, your kindness, prayers, and inspiring words have shaped me since I was an inexperienced agent and you remind me every day to be humble and avoid overestimation bias.

I would like to extend my heartfelt gratitude to Ivo, whose guidance has acted as a masterful policy update in my iterative learning process. Your insightful gradients for navigating the high-dimensional space of mathematics significantly refined my approximation models in this field. Like a well-tuned learning rate, your patience and encouragement helped me discover rewarding state transitions, both in this thesis and in my personal growth as a researcher. The discussions we had at the intersection of Machine Learning and classical Optimization were indeed very rewarding actions.

I also want to deeply thank my supervisor, Eligius, whose uncompromising commitment to formalism fined-tuned aspects of this work. You established an environment in Spain where my learning rate could accelerate and exploration could thrive. You have been indeed a good sparring partner.

Thanks goes to Jan Kronqvist as well for his invaluable input on the ensemble learning part of this thesis. Our engaging discussions on novel algorithm

---

concepts consistently guided me to new and exciting states.

Last but certainly not least, I extend my deepest appreciation to my girlfriend Jessica. Thank you for adjusting your own reward function to accommodate my busier days and for tolerating my habit of sneaking machine learning theory into our discussions. Your support has provided me with joy and comfort during the most challenging episodes of this journey.

RL terminology aside, I owe my sincerest thanks to all of you. This achievement would not have been possible without you.

# Contents

<b>Abstract</b>	<b>i</b>
<b>Acknowledgements</b>	<b>iii</b>
<b>1 Introduction</b>	<b>1</b>
1.1 Introduction to Machine Learning . . . . .	1
1.2 Artificial Neural Networks . . . . .	3
1.3 Motivation, Goals and Contributions . . . . .	6
<b>2 Introduction to Reinforcement Learning Algorithms</b>	<b>11</b>
2.1 Reinforcement Learning . . . . .	11
2.2 Maximum-Entropy Reinforcement Learning . . . . .	42
2.3 Distributional Reinforcement Learning . . . . .	49
<b>3 Distributional Reinforcement Learning in the Maximum-Entropy Setting Using the Cramér Metric</b>	<b>67</b>
3.1 Related Work . . . . .	67
3.2 The C-DSAC Framework . . . . .	68
3.3 Implementation . . . . .	75
3.4 Analysis . . . . .	77
3.5 Experiments . . . . .	81
3.6 Conclusion . . . . .	88

---

<b>4</b>	<b>Model-Enhanced Reinforcement Learning for Robotic Manipulation with Inverse Kinematics</b>	<b>91</b>
4.1	Forward Kinematics . . . . .	92
4.2	Inverse Kinematics . . . . .	95
4.3	Separation of Planning and Execution of Movement: RL-IK . . . . .	96
4.4	Reward Engineering for RL-IK . . . . .	98
4.5	Experiments . . . . .	99
<b>5</b>	<b>Introduction to Supervised Learning and Ensemble Strategies</b>	<b>101</b>
5.1	Supervised Learning . . . . .	101
5.2	Ensemble Learning . . . . .	102
5.3	Relevance for Robotics . . . . .	104
<b>6</b>	<b>Optimizing Ensemble Models via Column Generation</b>	<b>107</b>
6.1	Introduction . . . . .	107
6.2	Diverse Ensemble Learning and LPBoost . . . . .	110
6.3	A Data-Reduction Ensemble Learning Algorithm (DDA) . . . . .	115
6.4	Numerical Experiments . . . . .	121
6.5	Discussion and Conclusion . . . . .	127
<b>7</b>	<b>Latent Space Boosting</b>	<b>135</b>
7.1	Introduction . . . . .	136
7.2	Latent Space Boosting . . . . .	136
7.3	Algorithm . . . . .	138
7.4	Results and Conclusion . . . . .	139
<b>8</b>	<b>Conclusions and Future Work</b>	<b>141</b>
8.1	Conclusion . . . . .	141
8.2	Future Work . . . . .	143
<b>A</b>	<b>Publications arising from this Thesis</b>	<b>145</b>
A.1	Journal Publication . . . . .	145
A.2	Submitted Journal Publication . . . . .	145
A.3	Publication in International Conference Proceedings . . . . .	145
8.4	Resumen . . . . .	147
	<b>Resumen en Español</b>	<b>147</b>





This chapter establishes the framework for this thesis. It introduces the concept of machine learning, along with several essential terms and key concepts that are necessary for understanding the more advanced topics explored in subsequent chapters. The chapter includes a brief technical excursion into one of the primary drivers behind the success of contemporary machine learning tools: Powerful function approximators, particularly artificial neural networks, inspired by the operating principles of (biological) neurons. This chapter outlines the key motivations and objectives of the research presented in this thesis.

## **1.1 Introduction to Machine Learning**

Machine learning is a branch of Artificial Intelligence (AI) that “refers to the capacity of computer systems to learn and adapt without following explicit instructions, by using algorithms and statistical models to analyse and draw inferences from patterns in data” [1]. Machine learning is used in leveraging historical data to predict outcomes, classify information, make data-driven decisions, cluster data points, reduce dimensionality, and generate “synthetic” data [2]. The field is broadly divided into four paradigms: supervised learning, unsupervised learning, semi-supervised learning, and reinforcement learning. Unsupervised learning employs algorithms to analyse and cluster unlabelled datasets, uncovering latent patterns or data groupings without human guidance. Supervised learning, in contrast, relies on labelled datasets where each data point is paired with a corresponding target output. The model predictions are evaluated against these target labels using a cost function, which quantifies the difference between predicted and actual values. The cost function serves as a critical optimization met-

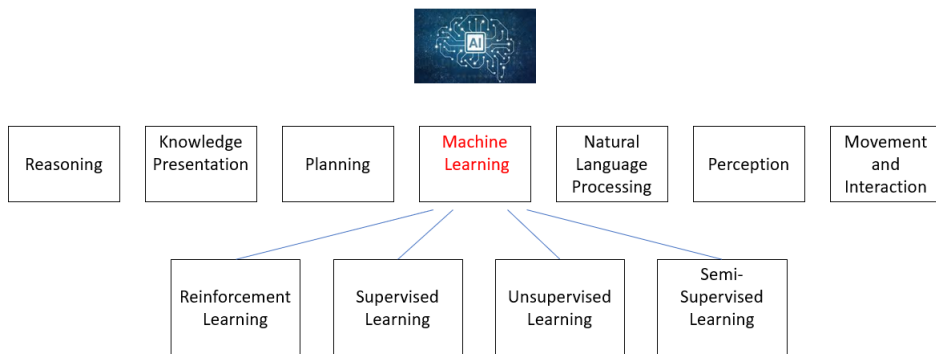
# 1. INTRODUCTION

---

ric, guiding the training process to minimize prediction errors [3, 4]. Semi-supervised learning integrates both labelled and unlabelled data into the learning process, bridging the methodologies of supervised and unsupervised learning [3]. Reinforcement learning addresses sequential decision-making problems typically formulated within the framework of Markov Decision Processes (MDPs) [5]. In reinforcement learning, the exact calculation of the cost function is hindered by the sparsity of the feedback mechanism. The cost must be approximated. Another distinguishing characteristic from other machine learning paradigms is the absence of a predefined dataset; instead, data must be generated through interactions with the environment.

An important concept in all machine learning branches is that of the generalization error. Generalization error, also referred to as out-of-sample error, quantifies model performance on unseen data. By assessing the model on a validation set, one obtains an estimation of this error, providing insight into the model's performance on unseen instances [3]. For example, a control model trained using a reinforcement learning algorithm is said to have a low generalization error if it consistently exhibits the desired control behaviour across a wide range of scenarios similar to those encountered during training.

Function approximators play a central role in machine learning methodologies, with artificial neural networks assuming a distinguished role. These networks are particularly effective for approximating complex functions in high-dimensional and large-scale datasets, making them indispensable tools for minimizing generalization error and enhancing model performance [3].



**Figure 1.1:** Hierarchical representation of the branches of artificial intelligence

## 1.2 Artificial Neural Networks

Artificial neural networks (ANNs) are inspired by biological neural networks (BNNs) and are conceptualized according to the Hebbian principle "neurons that fire together wire together" [6]. This succinct statement encapsulates the concept of plasticity in biological neural networks, where neurons responding to similar stimuli forge associations, leading to heightened responsiveness to these stimuli. In artificial neural networks, this concept is adopted and realized with weightings between the nodes. In addition to their basis in the Hebbian principle, ANNs also integrate mathematical models of learning and optimization, allowing them to simulate complex behaviours and solve intricate problems across various domains, ranging from image recognition to natural language processing. These networks have demonstrated remarkable success in mimicking cognitive functions, advancing fields such as supervised and reinforcement learning.

### 1.2.1 Different kinds of ANNs

The conceptualization of modern ANN elements traces back to 1943, when the perceptron was introduced [7]. In 1958, a hardware-software integrated approach was taken with the creation of the Mark I Perceptron, a machine designed for image classification [8], able to differentiate linearly separable images. The limitations of single-layer perceptrons were elucidated in 1969 with the proof of their inability to learn XOR logical functions [9]. However, the discovery that multi-layer perceptrons could approximate arbitrarily complex functions ignited a resurgence of interest in the field, leading to the proliferation of ANNs with diverse architectures and purposes. Notably, in reinforcement learning, multilayer perceptrons (MLPs) and convolutional neural networks (CNNs) have emerged as effective approximators for control and feature extraction functions, fostering a renewed enthusiasm for RL [10]. As of 2025, the effectiveness of large deep neural networks is evident in the realm of large language models (LLMs), exemplified by "ChatGPT," showcasing adept communication, problem-solving abilities, and broad domain knowledge [11]. Contemporary LLMs predominantly rely on the transformer architecture [12], typically possessing billions of parameters.

### 1.2.2 Description of MLP

The multilayer perceptron is the most commonly used ANN. It consists of nodes that are organised in layers  $\mathbf{I}^{(i)}$  with a maximum number of  $L$  layers. Two special layers are the input layer (index  $i = 0$ ) and the output layer (index  $i = L$ ). The other layers are referred to as hidden layers. The layers have activation functions  $\mathbf{a}^{(i)}(\cdot)$  that are

## 1. INTRODUCTION

---

funneled by  $\sigma$ . The weights  $\mathbf{w}^{(i)}$  and biases  $\mathbf{b}^{(i)}$  constitute the parameters of layer  $i$ . The *activation of the layer*  $\mathbf{a}^{(i)}$  is defined as

$$\mathbf{a}^{(i)}(\mathbf{x}) := \sigma(\mathbf{w}^{(i)}\mathbf{a}^{(i-1)}(\mathbf{x}) + \mathbf{b}^{(i)}),$$

where

$$\mathbf{a}^{(i-1)}(\mathbf{x}) = \sigma(\mathbf{w}^{(i-1)}\mathbf{a}^{(i-2)}(\mathbf{x}) + \mathbf{b}^{(i-1)}),$$

with input layer activation  $\mathbf{a}^{(0)}(\mathbf{x}) := \mathbf{x}$ . With this definition, the output of a MLP can be described by  $f_\theta(\mathbf{x}) := \mathbf{a}^{(L)}(\mathbf{x})$ . The goal is to find the parameters of  $f_\theta$  to approximate a target function.

### 1.2.3 Training and Backpropagation

Fitting of a MLP refers to the process of tuning the parameters of a neural network to approximate a target  $y$  given  $x$ . With a certain amount of nodes per layer spanning the parameter space  $\Theta$ , the goal is to find optimal parameters  $\theta^*$  consisting of weights  $\mathbf{w}$  and biases  $\mathbf{b}$  that minimises the loss  $\mathcal{L}$  by:

$$\theta^* = \underset{\theta \in \Theta}{\operatorname{argmin}} \{ \mathbb{E}[\mathcal{L}(f_\theta(\mathbf{x}), y)] \}.$$

In some machine learning paradigms, the data is given by a finite set of pairs  $T = \{(x_j, y_j)\}_{j \in J}$  that follow an underlying joint probability distribution  $P(x, y)$  over  $\mathcal{X} \times \mathcal{Y}$ . Hence  $\mathbb{E}$  denotes the expectation regarding that joint probability distribution.

The parameters are updated with gradient descent. Let  $\mathbf{z}^{(i)} := \mathbf{w}^{(i)}\mathbf{a}^{(i-1)}(\mathbf{x}) + \mathbf{b}^{(i)}$ , then the gradient for  $f_\theta$  w.r.t. to the parameters of the last layer is obtained by

$$\nabla_{\mathbf{w}^{(L)}} \mathcal{L} = \nabla_{\mathbf{a}^{(L)}} \mathcal{L} \nabla_{\mathbf{z}^{(L)}} \mathbf{a}^{(L)} \nabla_{\mathbf{w}^{(L)}} \mathbf{z}^{(L)}$$

$$\nabla_{\mathbf{b}^{(L)}} \mathcal{L} = \nabla_{\mathbf{a}^{(L)}} \mathcal{L} \nabla_{\mathbf{z}^{(L)}} \mathbf{a}^{(L)} \nabla_{\mathbf{b}^{(L)}} \mathbf{z}^{(L)}$$

With  $\nabla_{\mathbf{w}^{(L)}} \mathbf{z}^{(L)} := \mathbf{a}^{(L-1)}$ ,  $\nabla_{\mathbf{b}^{(L)}} \mathbf{z}^{(L)} := \mathbf{1}$  and  $\sigma'(\mathbf{z}^{(L)}) := \nabla_{\mathbf{z}^{(L)}} \mathbf{a}^{(L)}$ , the updates become

$$\mathbf{w}^{(L)} \leftarrow \mathbf{w}^{(L)} - \alpha \nabla_{\mathbf{a}^{(L)}} \mathcal{L} \sigma'(\mathbf{z}^{(L)}) \mathbf{a}^{(L-1)}$$

and

$$\mathbf{b}^{(L)} \leftarrow \mathbf{b}^{(L)} - \alpha \nabla_{\mathbf{a}^{(L)}} \mathcal{L} \sigma'(\mathbf{z}^{(L)})$$

If  $\mathbf{a}^{(L-1)}$  depends on parameters of deeper layers, the gradients must be further back-tracked:

$$\nabla_{\mathbf{w}^{(L-1)}} \mathcal{L} = \nabla_{\mathbf{z}^{(L)}} \mathcal{L} \nabla_{\mathbf{a}^{(L-1)}} \mathbf{z}^{(L)} \nabla_{\mathbf{w}^{(L-1)}} \mathbf{a}^{(L-1)} \quad (1.1)$$

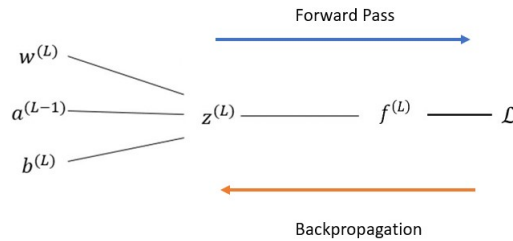
Expanding the terms results in

$$\begin{aligned}\nabla_{\mathbf{z}^{(L)}} \mathcal{L} &= \nabla_{\mathbf{a}^{(L)}} \mathcal{L} \nabla_{\mathbf{z}^{(L)}} \mathbf{a}^{(L)}, \\ \nabla_{\mathbf{a}^{(L-1)}} \mathbf{z}^{(L)} &= \mathbf{w}^{(L)}, \\ \nabla_{\mathbf{w}^{(L-1)}} \mathbf{a}^{(L-1)} &= \nabla_{\mathbf{z}^{(L-1)}} \mathbf{a}^{(L-1)} \nabla_{\mathbf{w}^{(L-1)}} \mathbf{z}^{(L-1)}.\end{aligned}$$

Simplifying the terms and plugging them in Equation (1.1) yields

$$\nabla_{\mathbf{w}^{(L-1)}} \mathcal{L} = \nabla_{\mathbf{a}^{(L)}} \mathcal{L} \nabla_{\mathbf{z}^{(L)}} \mathbf{a}^{(L)} \mathbf{w}^{(L)} \nabla_{\mathbf{z}^{(L-1)}} \mathbf{a}^{(L-1)} \mathbf{a}^{(L-2)}$$

The alterations are traced recursively until reaching the input layer. The process of propagating gradients back to the input layer is termed **backpropagation**. Backpropagation incurs significant computational expense. In practice, multiple pairs are aggregated into minibatches of size  $K$ . The loss is then computed as the average  $\frac{1}{K} \sum_{j=1}^K \mathcal{L}(f_{\theta}(\mathbf{x}_j), y_j)$  over the minibatch, and a single gradient descent step is executed based on the mean loss. The concept of a multi-layer perceptron is summarised in Figure 1.2.



**Figure 1.2:** Visualizing elementary dynamics in multilayer perceptron: Forward pass (for inference) and backpropagation (for learning).

### 1.2.4 Application

Various types of ANNs serve as function approximators across machine learning paradigms. In supervised learning, ANNs are employed for tasks such as classification and regression. For instance, classification tasks encompass discerning between spam and non-spam emails [3] or identifying product irregularities in manufacturing using CNNs [13]. Regression problems may involve estimating the energy consumption in heavy industry [14]. In addition to traditional feedforward neural networks, large visual models (LVMs), notably implemented with transformers similar to the Vision Transformer

## 1. INTRODUCTION

---

(ViT) [15], are pivotal in generative artificial intelligence. Generative models have demonstrated notable success in text-to-image generation [16] and voice cloning [17]. Moreover, ANNs are utilised as MLPs in reinforcement learning to approximate the control function.

### 1.3 Motivation, Goals and Contributions

First introduced at the Hannover Fair in 2011 [18], the concept of "Industry 4.0" represents a paradigm shift in production technology. This shift is characterized by the integration of connectivity, machine cognitivity, innovative human-machine interfaces, and versatile engineering technologies into the production chain. This transformative approach would enable mass production of highly customized products [19]. The proposed concepts of Industry 4.0 aim to revolutionize traditional manufacturing processes by creating intelligent systems capable of operating autonomously and effectively in complex production environments (see also the abstract).

The primary driver of the research presented in this thesis is to contribute to the machine cognitivity aspect of Industry 4.0, addressing the algorithmic challenges associated with generating sophisticated models for decision-making. Reinforcement learning (RL) algorithms, particularly their well-performing distributional variants, are promising avenues for establishing robust and versatile control models. Flexible control models are essential for navigating the complexities of modern industrial systems and are discussed in greater detail in Section 2.3. Distributional reinforcement learning (DistRL) algorithms, in particular, have demonstrated remarkable performance in solving decision-making problems within black-box contexts and under conditions of uncertainty. Their ability to capture the full distribution of potential rewards makes them particularly effective for tasks requiring nuanced decision-making. It is this capability that motivates the following research efforts.

A prerequisite for learning optimal control models is accurate machine perception. Effective decision-making depends on the availability of a clearly defined state. That in turn is possible only if all relevant information is accessible to the acting instance. Falling into the category of machine perception are sensory modalities such as vision, acoustics, and tactility. Machine vision stands out as particularly critical, since it often allows for the inference of additional information from visual data alone. Spatial and contextual cues obtained through vision can provide valuable information that would otherwise require supplementary sensory data, e.g. capacitance measurement for collision detection. By addressing challenges in machine cognitivity, this thesis

seeks to advance the state of the art in reinforcement learning and explores strategies to improve state identification through supervised learning methodologies (see Chapter 5). The ultimate goal is to develop a novel RL algorithm that not only excels in controlled experimental environments but also demonstrates potential for real-world applications in autonomous production facilities. The work’s objectives are described in more detail in the next section.

### 1.3.1 Goals

The objective of this thesis is to contribute to the understanding of DistRL and develop a novel DistRL algorithm that ideally outperforms current state-of-the-art DistRL/RL algorithms. Ensemble learning methods are investigated as a complementary approach for Supervised Learning (SL) tasks. This section outlines the specific objectives that guide the research presented in the following chapters. It is well-established that DistRL algorithms often outperform their base variants. While various hypotheses have been proposed to explain this superior performance, a consensus among researchers has yet to be reached [20]. This work seeks to identify one or more factors responsible for the enhanced capabilities of DistRL algorithms. By addressing this gap in understanding, the operating principles of DistRL are further explored and results stemming from these efforts might form the foundation for further research in this domain. A key metric for improvement in the field of reinforcement learning, as defined in this thesis, is the ability of an algorithm to converge to a desired level of performance in less iterations compared to the current state-of-the-art approaches. To this end, the thesis begins by exploring the characteristics of existing DistRL methods and identifying promising research directions. Based on these insights, a working hypothesis for an improved DistRL algorithm is formulated. The proposed algorithm, called C-DSAC, is rigorously analyzed for its convergence properties, and its performance is thoroughly evaluated through comprehensive experiments in benchmark environments. Ideally, the novel algorithm will constitute an advancement over existing methods. While reinforcement learning is the central focus of this thesis, ensemble learning, a subset of supervised learning, is also explored as part of a tangential research direction. Specifically, ensemble methods are investigated for image classification tasks, with the aim of enhancing robotic vision using low-complexity models. This exploration serves as a precursor and potential auxiliary to state perception in robotics, complementing the primary research on DRL. Readers who wish to bypass the detailed discussions and analyses can find a concise summary of the findings in Chapter 3.6.

The structure of this thesis is organized as follows. Chapter 2 introduces RL,

## 1. INTRODUCTION

---

entropy-augmented RL, and DistRL, while also exploring methods for DistRL implementation. Building on this foundation, Chapter 3 presents the primary contribution of this thesis: a novel entropy-augmented DistRL algorithm based on the Cramér metric, which demonstrates superior performance compared to current RL algorithms. Chapter 4 shifts the focus to model-based RL in robotics, proposing an innovative approach for efficiently generating control using inverse kinematics. The SL counterpart to Chapter 2 is discussed in Chapter 5, which establishes foundational concepts and methods. Chapter 6 explores a novel ensemble learning approach that bridges classical optimization and machine learning through linear programming. In contrast, Chapter 7 details another ensemble learning method, relying on neural networks and non-linear techniques. Finally, Chapter 8 concludes the thesis, summarizing the key findings and providing an outlook on future research directions in reinforcement learning.

### 1.3.2 Contributions

This thesis delivers four principal advances in distributional reinforcement learning (DistRL) and related supervised learning methods, culminating in both theoretical insights and practical algorithms for robotic control and perception. First, a rigorous theoretical framework for entropy-augmented, actor-critic DistRL under the Cramér metric is established. Building on an extensive literature review, it is proved that, under standard assumptions, algorithms developed in this framework converge to the optimal distributional state-action function and policy.

Second, guided by this theory, Cramér-based Distributional Soft Actor-Critic (C-DSAC) is introduced, a novel neural-network implementation that achieves **state-of-the-art performance** in benchmark robotic environments. Analysis reveals that the algorithm possesses a property termed confidence-driven model updates, which contributes to its superior performance. The results of this analysis deepen the understanding of why DistRL algorithms frequently outperform their expectation-based counterparts.

The contributions of this work are extended to specialized robotic applications. "Reinforcement Learning - Inverse Kinematics" (RL-IK) is presented, a model-based approach that decouples task planning from low-level control by integrating a robot's kinematic model into the learning loop. Empirical results show that RL-IK not only accelerates task learning but also produces smoother, more natural robotic movements.

Finally, as a complementary investigation in supervised learning, we propose a novel boosting algorithm that unites column-generation techniques from classical optimization with neural-network-based classifiers. While this method's immediate practical

### **1.3 Motivation, Goals and Contributions**

---

impact is modest, it bridges methodological gaps between optimization theory and modern learning frameworks, paving the way for more research at the intersection of these fields.



## Introduction to Reinforcement Learning Algorithms

### 2.1 Reinforcement Learning

In science and technology, deep reinforcement learning algorithms have been successfully applied to various problems. A notable example is AlphaTensor [21]; a framework/agent exhibiting superhuman performance for automating algorithmic discoveries in matrix multiplication. Reinforcement learning [5] is used in contemporary large language models for fine-tuning and serves as the primary method to eliminate undesired responses from the model [22]. In robotic tasks, this approach has been explored for data-driven generation of optimal control [23]. Nevertheless, applying reinforcement learning presents significant challenges, primarily due to the high sample complexity and the limited amount of data available in real-world scenarios [24]. What has begun as a promising approach to learn discrete control policies [10], and has since evolved to address general problems and mitigate various limitations [25, 26, 27, 28, 29, 30], appears to have reached a developmental standstill. A more recent approach by [20, 31], which expands and formalizes the distributional perspective on reinforcement learning, has revitalized interest in the field, including among researchers from other disciplines [32].

#### 2.1.1 Basic Concepts

##### 2.1.1.1 Introduction to MDP

A Markov Decision Process (MDP) models decision-making in dynamic environments [33]. It consists of a set of states  $\mathcal{S}$ , actions  $\mathcal{A}$ , a reward function  $R$ , and a transition kernel  $P$ , which governs how actions affect the environment [34]. In this work, MDPs are assumed to have the Markov Property, namely that "the probability distribution of

## 2. INTRODUCTION TO REINFORCEMENT LEARNING ALGORITHMS

---

future states of the process conditioned on both the past and present states depends only on the present state” [35]. A policy  $\pi$  selects actions to maximize an objective in this dynamic setting, with the assumption of full observability. This means the resulting next state is always known. The objective is tied to the reward function and can involve maximizing the total reward over a finite horizon, maximizing the cumulative discounted reward, or maximizing the average reward over time. The discounted reward is the most widely studied, as it adjusts the importance of immediate versus future rewards, thus influencing an agents planning behaviour. Note that as the future reward is discounted, variance is reduced at the cost of introducing bias.

### 2.1.1.2 Problem Formulation

In reinforcement learning, the interaction of an agent with its environment is modelled as a MDP  $\mathcal{M} = (\mathcal{S}, \mathcal{A}, r, P, \gamma)$ . In this context,  $\mathcal{S}$  is the finite state space,  $\mathcal{A}$  denotes the finite action space,  $r : \mathcal{S} \times \mathcal{A} \rightarrow [r_{min}, r_{max}]$  is the bounded reward function,  $P : \mathcal{S} \times \mathcal{A} \rightarrow Pr(\mathcal{S})$  defines the transition probability dynamics and  $\gamma \in (0, 1)$  denotes the discount factor. At each time step  $t$ , an agent executes an action  $a_t \in \mathcal{A}$  from state  $s_t \in \mathcal{S}$  according to a policy  $\pi$  and receives a reward  $r(s_t, a_t)$ , transitioning into the next state  $s_{t+1} \sim P(\cdot | s_t, a_t)$ . An episodic MDP considers to reach a terminal state at  $T$  within finite time. In this work, the policy is treated as a stochastic quantity, modelled as a probability distribution over the action space  $\pi : \mathcal{S} \rightarrow Pr(\mathcal{A})$ . The action is sampled from the policy, i.e.  $a_t \sim \pi(\cdot | s_t)$ . The general goal is to find an optimal policy  $\pi^*$  that maximizes the expected discounted return along a trajectory of interactions,

$$\begin{aligned}
 J(\pi) &:= \mathbb{E}_{\substack{s_0 \sim d_0 \\ a_t \sim \pi(\cdot | s_t) \forall t=0, \dots, T-1, \\ s_{t+1} \sim P(\cdot | s_t, a_t) \forall t=0, \dots, T-1}} \left[ \sum_{t=0}^{T-1} \gamma^t r(s_t, a_t) \right] \\
 &:= \mathbb{E}_{\pi, P} \left[ \sum_{t=0}^{T-1} \gamma^t r(s_t, a_t) \right],
 \end{aligned} \tag{2.1}$$

where  $s_T$  is an absorbing state.

The initial state is obtained by drawing from the distribution  $s_0 \sim d_0$ . Assuming discrete spaces, equation (2.1) can be written more explicitly as

$$J(\pi) = \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} \sum_{t=0}^{T-1} \Pr(s_t = s | \pi, P, s_0 \sim d_0) \pi(a | s) \gamma^t r(s, a) \tag{2.2}$$

In this setting, the episode has a length of  $T$  and thus  $J$  optimizes a discounted return in a finite MDP. Unlike an infinite-horizon MDP, where an upper bound of  $J$

is the maximum value of the geometric series ad infinitum,  $\max_{\pi} J(\pi) \leq \frac{1}{1-\gamma} r_{max}$ , an upper bound of  $J$  on a finite MDP is  $\max_{\pi} J(\pi) \leq \frac{1-\gamma^T}{1-\gamma} r_{max}$ . This raises a concern, as distant future rewards may be effectively ignored. The problem can be avoided by adjusting the discount factor to align with the average episode length. Additionally, heuristics such as a hyperparameter search or adherence to established values can be employed [36]. In that way, accurate value estimations across different environments can be ensured.

Optimal policies can be approximated by policy iteration, a process that incorporates policy evaluation and improvement using Q-values, see Lemma 2.5.

## 2.1.2 Frameworks

### 2.1.2.1 Critic Only

A Q-value induced by a policy  $\pi$  is denoted by  $Q^{\pi}(s_t, a_t)$ . In critic-only methods, the (greedy) policy is derived from the value function by

$$\hat{\pi}(s) = \operatorname{argmax}_{a \in \mathcal{A}} Q(s, a), \quad (2.3)$$

i.e.,  $\hat{\pi}$  denotes the approximation of the optimal policy (Technically,  $\hat{\pi}$  is a degenerate distribution that puts all the mass on the maximizing action). As a result, the superscript for the policy will be omitted in this work for critic-only approaches. The Q-value is the return conditioned on the state-action pair  $(s_t, a_t)$ ,

$$\begin{aligned} Q(s_k, a_k) &:= \mathbb{E}_{\substack{a_t \sim \hat{\pi}(\cdot | s_t) \forall t=k+1, \dots, T-1 \\ s_{t+1} \sim P(\cdot | s_t, a_t) \forall t=k, \dots, T-1}} \left[ \sum_{t=k}^{T-1} \gamma^{t-k} r(s_t, a_t) \mid s_k, a_k \right] \\ &:= \mathbb{E}_{\hat{\pi}, P} \left[ \sum_{t=k}^{T-1} \gamma^{t-k} r(s_t, a_t) \mid s_k, a_k \right]. \end{aligned} \quad (2.4)$$

To enable tractable computation of the Q-value, the sequential structure of MDPs is exploited by decomposing it into an immediate reward and an expected future return, yielding the Bellman equation. The optimal Q-value can then be expressed as

$$Q(s_t, a_t) = r(s_t, a_t) + \gamma \mathbb{E}_{s_{t+1} \sim P(\cdot | s_t, a_t)} \left[ \max_a Q(s_{t+1}, a) \right]. \quad (2.5)$$

Based on (2.5), we define the Bellman optimality operator  $\mathcal{T}$  by

$$\mathcal{T}Q(s_t, a_t) := r(s_t, a_t) + \gamma \mathbb{E}_{s_{t+1} \sim P(\cdot | s_t, a_t)} \left[ \max_a Q(s_{t+1}, a) \right]. \quad (2.6)$$

## 2. INTRODUCTION TO REINFORCEMENT LEARNING ALGORITHMS

---

The one-step update rule thus becomes

$$Q(s_t, a_t) \leftarrow (1 - \alpha)Q(s_t, a_t) + \alpha(r(s_t, a_t) + \gamma[\max_a Q(s_{t+1}, a)] - Q(s_t, a_t)), \quad (2.7)$$

where  $\alpha$  is the learning rate. The one-step temporal difference error  $\delta := r(s_t, a_t) + \gamma \max_a Q(s_{t+1}, a) - Q(s_t, a_t)$  is relevant for the main algorithm of this work presented in Section 3.2. For optimal policy evaluation, it can be shown that the sequence of value functions  $Q_{k+1} := \mathcal{T}Q_k$ , starting from some arbitrary  $Q_0$ , will exponentially quickly converge to  $Q^*$  as  $k$  increases.

**Lemma 2.1.** Let  $Q_1(s, a)$  and  $Q_2(s, a)$  be two bounded state-action value functions. Assume  $|\mathcal{A}| < \infty$ ,  $|\mathcal{S}| < \infty$  and  $0 \leq \gamma < 1$  (see Section 2.1.1.2). Then  $\mathcal{T}$  according to (2.6) is a  $\gamma$ -contraction, i.e.

$$\|\mathcal{T}Q_1 - \mathcal{T}Q_2\|_\infty \leq \gamma\|Q_1 - Q_2\|_\infty. \quad (2.8)$$

*Proof.* To show that (2.8) holds, the left-hand side is expanded. Set  $\mathbb{E}_P[\cdot] := \mathbb{E}_{s_{t+1} \sim P(\cdot | s_t, a_t)}[\cdot]$ , then

$$\begin{aligned} & \|\mathcal{T}Q_1 - \mathcal{T}Q_2\|_\infty \\ &= \left\| \left( r(s_t, a_t) + \gamma \mathbb{E}_P \left[ \max_a Q_1(s_{t+1}, a) \right] \right) - \left( r(s_t, a_t) + \gamma \mathbb{E}_P \left[ \max_a Q_2(s_{t+1}, a) \right] \right) \right\|_\infty \\ &= \sup_{s_t \in \mathcal{S}, a_t \in \mathcal{A}} \left| \gamma \mathbb{E}_P \left[ \max_a Q_1(s_{t+1}, a) \right] - \gamma \mathbb{E}_P \left[ \max_a Q_2(s_{t+1}, a) \right] \right| \\ &= \gamma \sup_{s_t \in \mathcal{S}, a_t \in \mathcal{A}} \left| \mathbb{E}_P \left[ \max_a Q_1(s_{t+1}, a) \right] - \mathbb{E}_P \left[ \max_a Q_2(s_{t+1}, a) \right] \right| \\ &\leq \gamma \sup_{s_t \in \mathcal{S}, a_t \in \mathcal{A}} \max_a \mathbb{E}_P \left| [Q_1(s_{t+1}, a)] - [Q_2(s_{t+1}, a)] \right| \\ &\leq \gamma \sup_{s \in \mathcal{S}, a \in \mathcal{A}} |Q_1(s, a) - Q_2(s, a)| = \gamma\|Q_1 - Q_2\|_\infty. \end{aligned}$$

Since  $\mathcal{T}$  is a  $\gamma$ -contraction, it follows from the Banach theorem that the series  $\{Q_k\}$  with  $Q_{k+1} := \mathcal{T}Q_k$  converges to a unique fixed point  $Q$ , i.e.  $Q = \mathcal{T}Q$  and  $\|Q_k - Q\|_\infty \rightarrow 0$  as  $k \rightarrow \infty$ .  $\square$

In practice, the optimal function  $Q$  is unknown and must be iteratively improved through exploration, such as with an  $\epsilon$ -greedy approach. With probability  $\epsilon$ , a random action is chosen, allowing updates to  $\max_a Q(s_t, a)$ . In Q-learning (see Section 2.1.6.1), these maximum values are stored in a tabular fashion. Convergence is guaranteed if the state-action values are updated and exploration of the environment occurs under GLIE assumptions, see Section 2.1.2.3.

The state value  $V^{\hat{\pi}}$  is similar to the aforementioned state-action value and can be expressed in terms of the latter as

$$V^{\hat{\pi}}(s_t) = \mathbb{E}_{a_t \sim \hat{\pi}(\cdot|s_t)} [Q(s_t, a_t)]. \quad (2.9)$$

Analogous to the state-action approach, the Bellman equation and operator can be defined for the state value function. In this context, properties and convergence of policy evaluation, improvement and iteration can be proven similarly. The value function is more relevant to on-policy algorithms, see Section 2.1.2.6.

### 2.1.2.2 Temporal Difference Learning and Advantage Estimator

Temporal Difference (TD) learning is a method to learn the value of a state or state-action pair from observations without requiring a model of the environment. TD learning generalizes concepts from Monte Carlo learning and dynamic programming. In  $TD(1)$ , returns are estimated based on the trajectory of an entire episode. Unlike Monte Carlo methods,  $TD(1)$  can be applied to discounted continuing tasks, i.e, learning can occur prior to episode termination if  $\gamma \neq 1$ . Using the Q-function formulation and  $G_t := \sum_{u=t}^T \gamma^{u-t} r(s_u, a_u)$ , the  $TD(1)$  update rule can be expressed as

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha(G_t - Q(s_t, a_t)),$$

where  $\alpha$  denotes the step size. The counterpart to  $TD(1)$  is  $TD(0)$ , where the Q-value is updated considering only the immediate future. An example for  $TD(0)$  is given in the update rule (2.7).  $TD(0)$  is widely used in reinforcement learning and is a foundational method for dynamic programming. Sutton introduced  $TD(\lambda)$  in [5], allowing to strike a balance between  $TD(0)$  and  $TD(1)$ . Defining the off-policy return for  $n$ -steps in the setting of state-action values as

$$G_{t:t+n} := r(s_t, a_t) + \gamma r(s_{t+1}, a_{t+1}) + \dots + \gamma^n \max_a Q(s_{t+n}, a), \quad 0 \leq t \leq T - n, \quad (2.10)$$

and the  $\lambda$ -return as

$$G_t^\lambda := (1 - \lambda) \sum_{n=1}^{T-t} \lambda^{n-1} G_{t:t+n} + \lambda^{T-t} G_t,$$

the  $TD(\lambda)$  update can be written as

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha(G_t^\lambda - Q(s_t, a_t)).$$

Here,  $\lambda$  is the trace-decay parameter,  $0 \leq \lambda \leq 1$ . Equation (2.10) indicates that the update rule (effectively) reduces to  $TD(0)$  and  $TD(1)$  with  $\lambda = 0$  and  $\lambda = 1$ ,

## 2. INTRODUCTION TO REINFORCEMENT LEARNING ALGORITHMS

---

respectively. Note also that a high value for  $\lambda$  reduces bias, but increases variance whereas a low value reduces variance but increases bias. Most notably,  $TD(\lambda)$  was used in 1995 to learn to play Backgammon [37]. Analogous methods for managing the bias-variance trade-off in policy updates exist and are discussed in Section 2.1.2.5.

### 2.1.2.3 Direct Policy Search

In contrast to critic-only methods described in Section 2.1.2.1, the class of actor-only reinforcement learning algorithms aims to directly optimize the parameters of the policy for maximizing the return. Thus, a value function is not needed. In reinforcement learning, where environment dynamics are unknown, direct policy search simplifies the problem to a purely stochastic one. Assuming some kernel  $P$  underlying the transitions and a stochastic policy parameterized by  $\phi$ , the objective to maximize is a function of the policy parameters  $\phi$  only, i.e.

$$J(\phi) = \mathbb{E}_{\pi_\phi, P} \left[ \sum_{t=0}^{T-1} \gamma^t r(s_t, a_t) \right]. \quad (2.11)$$

The above objective optimizes This stochastic problem can be solved by various methods (evolutionary approach, cross-entropy, ..., etc.) of which only the gradient-based method is relevant for this work. The parameters of the policy are updated with gradient update rule [38]

$$\phi \leftarrow \phi + \alpha \nabla_\phi J. \quad (2.12)$$

The policy gradient theorem illustrates the update of a parameterized policy.

**Theorem 2.2** (Policy Gradient Theorem). *Let  $\mathcal{M} = (\mathcal{S}, \mathcal{A}, r, P, \gamma)$  be finite horizon, discrete MDP emitting a bounded reward and let  $\pi_\phi$  be a differentiable policy parameterized by  $\phi \in \mathbb{R}^d$ . Define the parameterized variant of the objective (2.2) to be*

$$J(\phi) := \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} \sum_{t=0}^{T-1} \gamma^t \Pr(s_t = s \mid \pi_\phi, P, s_0 \sim d_0) \pi_\phi(a \mid s) r(s, a),$$

where  $d_0$  is the initial state distribution independent of  $\phi$ ,  $a_t \sim \pi_\phi(\cdot \mid s_t)$ , and  $s_{t+1} \sim P(\cdot \mid s_t, a_t)$ . Then

$$\nabla_\phi J(\phi) = \mathbb{E}_{\pi_\phi, P} \left[ \sum_{u=0}^{T-1} \gamma^u \nabla_\phi \log \pi_\phi(a_u \mid s_u) G_u \right],$$

with the return  $G_u := \sum_{t=u}^{T-1} \gamma^{t-u} r(s_t, a_t)$ .

*Proof.* Finite sums allow summation and gradient interchanges

$$\begin{aligned} J(\phi) &= \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} \sum_{t=0}^{T-1} \gamma^t \Pr(s_t = s \mid \pi_\phi, P, s_0 \sim d_0) \pi_\phi(a \mid s) r(s, a) \\ &= \sum_{t=0}^{T-1} \gamma^t \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} \Pr(s_t = s \mid \pi_\phi, P, s_0 \sim d_0) \pi_\phi(a \mid s) r(s, a). \end{aligned}$$

The probability of the trajectory  $\tau = (s_0, a_0, \dots, a_{T-1}, s_T)$  can be written as

$$\rho_\phi(\tau) = d_0(s_0) \prod_{u=0}^{T-1} \pi_\phi(a_u \mid s_u) P(s_{u+1} \mid a_u, s_u).$$

Note that

$$\begin{aligned} \sum_{\tau} \rho_\phi(\tau) r(s_t, a_t) &= \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} \Pr(s_t = s, a_t = a \mid \pi_\phi, P, s_0 \sim d_0) r(s, a) \\ &= \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} \Pr(s_t = s \mid \pi_\phi, P, d_0) \pi_\phi(a \mid s) r(s, a). \end{aligned}$$

Then the gradient of the objective can be rewritten as

$$\begin{aligned} \nabla_\phi J(\phi) &= \nabla_\phi \left[ \sum_{t=0}^{T-1} \gamma^t \sum_{\tau} \rho_\phi(\tau) r(s_t, a_t) \right] \\ &\text{By finiteness of the trajectory space} \\ &= \nabla_\phi \sum_{\tau} \rho_\phi(\tau) \left[ \sum_{t=0}^{T-1} \gamma^t r(s_t, a_t) \right] \\ &= \sum_{\tau} \nabla_\phi \rho_\phi(\tau) \left[ \sum_{t=0}^{T-1} \gamma^t r(s_t, a_t) \right] \tag{2.13} \end{aligned}$$

Using the identity  $\nabla f = f \nabla \log f$

$$\begin{aligned} &= \sum_{\tau} \rho_\phi(\tau) \nabla_\phi \log \rho_\phi(\tau) \left[ \sum_{t=0}^{T-1} \gamma^t r(s_t, a_t) \right] \\ &= \mathbb{E}_{\pi_\phi, P} \left[ \nabla_\phi \log \rho_\phi(\tau) \left[ \sum_{t=0}^{T-1} \gamma^t r(s_t, a_t) \right] \right]. \end{aligned}$$

Observe that the initial state distribution function and the transition kernel are inde-

## 2. INTRODUCTION TO REINFORCEMENT LEARNING ALGORITHMS

---

pendent of  $\phi$

$$\begin{aligned}\nabla_{\phi} \log \rho_{\phi}(\tau) &= \nabla_{\phi} \log d_0(s_0) + \nabla_{\phi} \sum_{u=0}^{T-1} \log \pi_{\phi}(a_u | s_u) + \nabla_{\phi} \sum_{u=0}^{T-1} \log P(s_{u+1} | a_u, s_u) \\ &= \nabla_{\phi} \sum_{u=0}^{T-1} \log \pi_{\phi}(a_u | s_u) \\ &= \sum_{u=0}^{T-1} \nabla_{\phi} \log \pi_{\phi}(a_u | s_u).\end{aligned}$$

Plugging this expression into Equation (2.13) yields

$$\begin{aligned}\nabla_{\phi} J(\phi) &= \mathbb{E}_{\pi_{\phi}, P} \left[ \sum_{u=0}^{T-1} \nabla_{\phi} \log \pi_{\phi}(a_u | s_u) \left[ \sum_{t=0}^{T-1} \gamma^t r(s_t, a_t) \right] \right] \\ &= \mathbb{E}_{\pi_{\phi}, P} \left[ \sum_{u=0}^{T-1} \nabla_{\phi} \log \pi_{\phi}(a_u | s_u) \left[ \sum_{t=0}^{u-1} \gamma^t r(s_t, a_t) + \gamma^u \sum_{t=u}^{T-1} \gamma^{t-u} r(s_t, a_t) \right] \right] \\ &= \mathbb{E}_{\pi_{\phi}, P} \left[ \sum_{u=0}^{T-1} \nabla_{\phi} \log \pi_{\phi}(a_u | s_u) \sum_{t=0}^{u-1} \gamma^t r(s_t, a_t) + \sum_{u=0}^{T-1} \nabla_{\phi} \log \pi_{\phi}(a_u | s_u) \gamma^u \sum_{t=u}^{T-1} \gamma^{t-u} r(s_t, a_t) \right].\end{aligned}\tag{2.14}$$

Define  $G_0^{(u-1)} := \sum_{t=0}^{u-1} \gamma^t r(s_t, a_t)$ . Note that  $G_0^{(u-1)}$  is measurable up to time  $u$ .

Therefore, it is permissible to write

$$\begin{aligned}E_{\pi_{\phi}, P} \left[ \sum_{u=0}^{T-1} \nabla_{\phi} \log \pi_{\phi}(a_u | s_u) G_0^{(u-1)} \right] &= \sum_{u=0}^{T-1} \mathbb{E}_{\pi_{\phi}, P} \left[ \mathbb{E}_{a_u \sim \pi_{\phi}(\cdot | s_u)} \left[ \nabla_{\phi} \log \pi_{\phi}(a_u | s_u) G_0^{(u-1)} \mid s_u \right] \right] \\ &= \sum_{u=0}^{T-1} \mathbb{E}_{\pi_{\phi}, P} \left[ G_0^{(u-1)} \mathbb{E}_{a_u \sim \pi_{\phi}(\cdot | s_u)} \left[ \nabla_{\phi} \log \pi_{\phi}(a_u | s_u) \mid s_u \right] \right].\end{aligned}$$

By the score function identity

$$= 0.\tag{2.15}$$

The last line follows from  $E_{a_u \sim \pi_{\phi}(\cdot | s_u)} [\nabla_{\phi} \log \pi_{\phi}(a_u | s_u) \mid s_u] = 0$ . Therefore, in the third line of Eq. (2.14), the first term within the expectation can be discarded.

Defining the second term as  $G_u := \sum_{t=u}^{T-1} \gamma^{t-u} r(s_t, a_t)$  and using it in Eq. (2.14), the policy gradient theorem can be written as

$$\nabla_{\phi} J(\phi) = \mathbb{E}_{\pi_{\phi}, P} \left[ \sum_{u=0}^{T-1} \gamma^u \nabla_{\phi} \log \pi_{\phi}(a_u | s_u) G_u \right].\tag{2.16}$$

□

Gradient-based direct policy search has strong convergence properties due to the ones it inherits from gradient descent methods [39]. To guarantee convergence to (a local) maximum,  $J$  must be continuously differentiable,  $\nabla J$  must satisfy Lipschitz conditions and the "Greedy in the Limit with Infinite Exploration (GLIE)" requirements must be met. The GLIE conditions are expressed as

$$\begin{aligned} \sum_{t=0}^{\infty} \alpha_t &= \infty \\ \sum_{t=0}^{\infty} \alpha_t^2 &= \infty. \end{aligned} \tag{2.17}$$

However, the gradient  $\nabla_{\phi} J$  can exhibit high variance depending on the length of the trajectory due to the stochastic nature of  $J$ . A direct policy algorithm will be introduced in Section 2.1.6.3.

#### 2.1.2.4 Actor-Critic

Actor-critic (AC) methods incorporate a policy and a value function, called actor and critic, respectively. The policy and value functions are dependent but separate structures. AC algorithms can handle continuous action spaces and reduce variance (potentially at the cost of higher bias) due to the critic. The critic is updated by evaluating the current policy. The state-action value function induced by a policy  $\pi$  is defined as

$$\begin{aligned} Q^{\pi}(s_k, a_k) &:= \mathbb{E}_{\substack{a_t \sim \pi(\cdot | s_t) \forall t=k+1, \dots, T-1 \\ s_{t+1} \sim P(\cdot | s_t, a_t) \forall t=k, \dots, T-1}} \left[ \sum_{t=k}^{T-1} \gamma^{t-k} r(s_t, a_t) \mid s_k, a_k \right] \\ &:= \mathbb{E}_{\pi, P} \left[ \sum_{t=k}^{T-1} \gamma^{t-k} r(s_t, a_t) \mid s_k, a_k \right], \end{aligned} \tag{2.18}$$

and satisfies the recursive Bellman equation

$$Q^{\pi}(s_t, a_t) = r(s_t, a_t) + \gamma \mathbb{E}_{\substack{a_{t+1} \sim \pi(\cdot | s_{t+1}), \\ s_{t+1} \sim P(\cdot | s_t, a_t)}} [Q^{\pi}(s_{t+1}, a_{t+1})]. \tag{2.19}$$

The policy is updated according to the estimate of the critic and the policy gradient

## 2. INTRODUCTION TO REINFORCEMENT LEARNING ALGORITHMS

---

in (2.16) becomes

$$\begin{aligned}
\nabla_\phi J(\phi) &= \mathbb{E}_{\pi_{\phi,P}} \left[ \sum_{u=0}^{T-1} \gamma^u \nabla_\phi \log \pi_\phi(a_u | s_u) G_u \right] \\
&= \sum_{u=0}^{T-1} \mathbb{E}_{\pi_{\phi,P}} \left[ \mathbb{E}_{\pi_{\phi,P}} [\gamma^u \nabla_\phi \log \pi_\phi(a_u | s_u) G_u | s_u, a_u] \right] \\
&\text{Note that the expression before } G_u \text{ is measurable} \\
&= \sum_{u=0}^{T-1} \mathbb{E}_{\pi_{\phi,P}} \left[ \gamma^u \nabla_\phi \log \pi_\phi(a_u | s_u) \mathbb{E}_{\pi_{\phi,P}} [G_u | s_u, a_u] \right] \tag{2.20}
\end{aligned}$$

From Eq. (2.19)

$$= \sum_{u=0}^{T-1} \mathbb{E}_{\pi_{\phi,P}} [\gamma^u \nabla_\phi \log \pi_\phi(a_u | s_u) Q^{\pi_\phi}(s_u, a_u)].$$

The policy can be optimized independent of the Q-function

$$\begin{aligned}
\nabla_\phi J(\phi) &= \sum_{u=0}^{T-1} \mathbb{E}_{\pi_{\phi,P}} [\gamma^u \nabla_\phi \log \pi_\phi(a_u | s_u) Q^{\pi_\phi}(s_u, a_u)] \\
&= \sum_{u=0}^{T-1} \mathbb{E}_{\pi_{\phi,P}} [\gamma^u [\nabla_\phi \log \pi_\phi(a_u | s_u) Q^{\pi_\phi}(s_u, a_u) + \log \pi_\phi(a_u | s_u) \nabla_\phi Q^{\pi_\phi}]] \tag{2.21} \\
&\approx \sum_{u=0}^{T-1} \mathbb{E}_{\pi_{\phi,P}} [\gamma^u \nabla_\phi \log \pi_\phi(a_u | s_u) Q^{\pi_\phi}(s_u, a_u)].
\end{aligned}$$

The last step is justified by [40] and it can be shown that the approximated gradient still converges to a true local maximum.

Based on (2.19), we define the Bellman operator  $\mathcal{T}^\pi$ ,

$$\mathcal{T}^\pi Q(s_t, a_t) := r(s_t, a_t) + \gamma \mathbb{E}_{\substack{s_{t+1} \sim P(\cdot | s_t, a_t), \\ a_{t+1} \sim \pi(\cdot | s_{t+1})}} [Q(s_{t+1}, a_{t+1})]. \tag{2.22}$$

Analogous to the proof in Section 2.1.2.1, we can show evaluation contraction of policy  $\pi$ .

**Lemma 2.3.** Let  $Q_1(s, a)$  and  $Q_2(s, a)$  be two bounded state-action value functions. Assume  $|\mathcal{A}| < \infty$ ,  $|\mathcal{S}| < \infty$  and  $0 \leq \gamma < 1$ . Then  $\mathcal{T}^\pi$  according to (2.22) is a  $\gamma$ -contraction, i.e.

$$\|\mathcal{T}^\pi Q_1 - \mathcal{T}^\pi Q_2\|_\infty \leq \gamma \|Q_1 - Q_2\|_\infty. \tag{2.23}$$

## 2.1 Reinforcement Learning

*Proof.* To show that (2.23) holds, the left-hand side is expanded. Set  $\mathbb{E}_P[\cdot] := \mathbb{E}_{s_{t+1} \sim P(\cdot | s_t, a_t)}[\cdot]$  and  $\mathbb{E}_{a_t \sim \pi(\cdot | s_t)}[Q(s_t, a_t)] := Q(s_t, \pi(\cdot | s_t))$ , then

$$\begin{aligned}
& \|\mathcal{T}^\pi Q_1 - \mathcal{T}^\pi Q_2\|_\infty \\
&= \left\| \left( r(s_t, a_t) + \gamma \mathbb{E}_P[Q_1(s_{t+1}, \pi(\cdot | s_{t+1}))] \right) - \left( r(s_t, a_t) + \gamma \mathbb{E}_P[Q_2(s_{t+1}, \pi(\cdot | s_{t+1}))] \right) \right\|_\infty \\
&= \sup_{s_t \in \mathcal{S}, a_t \in \mathcal{A}} \left| \gamma \mathbb{E}_P[Q_1(s_{t+1}, \pi(\cdot | s_{t+1}))] - \gamma \mathbb{E}_P[Q_2(s_{t+1}, \pi(\cdot | s_{t+1}))] \right| \\
&= \gamma \sup_{s_t \in \mathcal{S}, a_t \in \mathcal{A}} \left| \mathbb{E}_P[(Q_1(s_{t+1}, \pi(\cdot | s_{t+1})) - Q_2(s_{t+1}, \pi(\cdot | s_{t+1}))) \right| \\
&\leq \gamma \sup_{s_t \in \mathcal{S}, a_t \in \mathcal{A}} \mathbb{E}_P[|Q_1(s_{t+1}, \pi(\cdot | s_{t+1})) - Q_2(s_{t+1}, \pi(\cdot | s_{t+1}))|] \\
&\leq \gamma \sup_{s \in \mathcal{S}} |Q_1(s, \pi(\cdot | s)) - Q_2(s, \pi(\cdot | s))| \\
&\leq \gamma \sup_{s \in \mathcal{S}, a \in \mathcal{A}} |Q_1(s, a) - Q_2(s, a)| = \gamma \|Q_1 - Q_2\|_\infty.
\end{aligned}$$

Since  $\mathcal{T}^\pi$  is a  $\gamma$ -contraction, it follows from the Banach theorem that the series  $\{Q_k\}$  with  $Q_{k+1} := \mathcal{T}^\pi Q_k$  converges to a unique fixed point  $Q^\pi$ , i.e.  $Q^\pi = \mathcal{T}^\pi Q^\pi$  and  $\|Q_k - Q^\pi\|_\infty \rightarrow 0$  as  $k \rightarrow \infty$ .  $\square$

Policy improvement involves exploitation of knowledge about the environment and updating towards higher values, i.e.

$$\pi_{k+1}(s_t) = \operatorname{argmax}_{\pi \in \Pi} \mathbb{E}_{a_t \sim \pi(\cdot | s_t)} Q^{\pi_k}(s_t, a_t). \quad (2.24)$$

**Lemma 2.4.** Given the policy update rule in (2.24) and assuming bounded state-action value functions,  $|\mathcal{A}| < \infty$ ,  $0 \leq \gamma < 1$  then

$$\mathbb{E}_{a_t \sim \pi_k(\cdot | s_t)}[Q^{\pi_k}(s_t, a_t)] \leq \mathbb{E}_{a_t \sim \pi_{k+1}(\cdot | s_t)}[Q^{\pi_{k+1}}(s_t, a_t)] \quad \forall s_t \in \mathcal{S}.$$

*Proof.* The proof relies on expanding the Q-function in the Bellman equation repeatedly with actions of the improved policy. For convenience, the notation  $\pi' := \pi_{k+1}$  and  $\pi := \pi_k$  and  $\mathbb{E}_{\pi', P, k}[\cdot | s_t] := \mathbb{E}_{a_{t:k-1} \sim \pi'(\cdot | s_{t:k-1}), s_{t+1:k} \sim P(\cdot | s_{t:k-1}, a_{t:k-1})}[\cdot | s_t]$  is used. Let  $Q^\pi(s_t, \pi(\cdot | s_t)) := \mathbb{E}_{a_t \sim \pi(\cdot | s_t)}[Q^\pi(s_t, a_t)]$ , then

## 2. INTRODUCTION TO REINFORCEMENT LEARNING ALGORITHMS

---

$$\begin{aligned}
Q^\pi(s_t, \pi(\cdot | s_t)) &\leq Q^\pi(s_t, \pi'(\cdot | s_t)) \\
&= \mathbb{E}_{\pi', P, t+1} [r(s_t, a_t) + \gamma Q^\pi(s_{t+1}, \pi(\cdot | s_{t+1})) | s_t] \\
&\leq \mathbb{E}_{\pi', P, t+1} [r(s_t, a_t) + \gamma Q^\pi(s_{t+1}, \pi'(\cdot | s_{t+1})) | s_t] \\
&= \mathbb{E}_{\pi', P, t+2} [r(s_t, a_t) + \gamma(r(s_{t+1}, a_{t+1}) + \gamma Q^\pi(s_{t+2}, \pi(\cdot | s_{t+2}))) | s_t] \\
&= \mathbb{E}_{\pi', P, t+2} [r(s_t, a_t) + \gamma r(s_{t+1}, a_{t+1}) + \gamma^2 Q^\pi(s_{t+2}, \pi(\cdot | s_{t+2})) | s_t] \\
&\leq \mathbb{E}_{\pi', P, t+3} [r(s_t, a_t) + \gamma r(s_{t+1}, a_{t+1}) + \gamma^2 r(s_{t+2}, a_{t+2}) + \gamma^3 Q^\pi(s_{t+3}, \pi(\cdot | s_{t+3})) | s_t] \\
&\vdots \\
&\leq \mathbb{E}_{\pi', P, T-1} [r(s_t, a_t) + \gamma r(s_{t+1}, a_{t+1}) + \gamma^2 r(s_{t+2}, a_{t+2}) + \dots + \gamma^{T-1} r(s_{T-1}, a_{T-1}) | s_t] \\
&= Q^{\pi'}(s_t, \pi'(\cdot | s_t)).
\end{aligned}$$

Thus for any  $s_t$ ,

$$\mathbb{E}_{a_t \sim \pi_k(\cdot | s_t)} [Q^{\pi_k}(s_t, a_t)] = Q^\pi(s_t, \pi(\cdot | s_t)) \leq Q^{\pi'}(s_t, \pi'(\cdot | s_t)) = \mathbb{E}_{a_t \sim \pi_{k+1}(\cdot | s_t)} [Q^{\pi_{k+1}}(s_t, a_t)].$$

□

In that manner, the policy will strictly be non-decreasing. Note that the finite-horizon assumption permits the use of a non-strictly contractive discount factor  $\gamma \leq 1$ , without compromising the stability of the value functions. The estimation of the state-action value can be improved by better exploration of the state-space, making the inequality generally strict.

The Q-function approximator is parameterized by  $\theta$  and is denoted as  $Q_\theta^\pi$ . Policy iteration is the method of alternating between evaluation and improvement. It can be shown that policy iteration converges towards the optimal value function and an optimal policy. Although its relevance is not exclusive to AC methods, in AC algorithms it can be shown explicitly due to the separation of the acting and evaluating instances.

**Lemma 2.5.** Let the reward be bounded,  $|\mathcal{A}| < \infty$ ,  $|\mathcal{S}| < \infty$ , and  $0 \leq \gamma < 1$ . Let  $\Pi$  be a the set of all stationary (possibly stochastic) policies. In particular,  $\Pi$  contains the point-mass policies. Alternating between exact policy evaluation and policy improvement from some initial policy  $\pi_0 \in \Pi$  and state-action function  $Q^{(0)}$ , the process terminates after finitely many improvement steps at the optimal policy  $\pi^*$  satisfying  $Q^{\pi^*} = Q^*$ .

*Proof.* Fix  $k$ . By Lemma 2.3, the inner evaluation sequence  $\{Q_m\}$  for  $Q_{m+1} := \mathcal{T}^{\pi_k} Q_m$  approaches  $Q^{\pi_k}$ , i.e.  $\|Q_m - Q^{\pi_k}\|_\infty \rightarrow 0$  as  $m \rightarrow \infty$ . It is established by Lemma 2.4 that an improvement of policy  $\pi_k$  by Eq. (2.24), yielding  $\pi_{k+1}$ , satisfies

$$V^{\pi_k}(s_t) = \mathbb{E}_{a_t \sim \pi_k(\cdot|s_t)}[Q^{\pi_k}(s_t, a_t)] \leq \mathbb{E}_{a_t \sim \pi_{k+1}(\cdot|s_t)}[Q^{\pi_{k+1}}(s_t, a_t)] = V^{\pi_{k+1}}(s_t) \forall s_t \in \mathcal{S}.$$

Although policies are treated as stochastic mappings, the argmax in Eq. (2.24) implies that actions are taken deterministically, provided the action set is finite. With finite states and actions and ties being broken deterministically per state (e.g. lexicographic smallest action), the set of greedy policies is also finite. Under these circumstances, a greedy policy for a fixed  $Q^{\pi_k}$  is unique, so  $\pi_{k+1} \neq \pi_k$  implies strict increase in the value  $V$  for at least one state. Eventually  $\pi_{k+1} = \pi_k$  in finite time, implying convergence to a fixed point  $\pi_{k+1} = \pi_k = \pi^*$ . The terminating policy  $\pi^*$  is greedy w.r.t. its own  $Q$ , therefore satisfies the Bellman optimality equation and thus  $Q^{\pi^*} = Q^*$ .  $\square$

### 2.1.2.5 Advantage and Generalized Advantage Estimator

The advantage is defined as the difference between the value obtained by executing a particular action from a state and an (approximated) baseline value for that particular state under policy  $\pi$  [41]. Formally it can be expressed as

$$A^\pi(s_t, a_t) := Q^\pi(s_t, a_t) - V^\pi(s_t). \tag{2.25}$$

The term "advantage" reflects how beneficial a particular action is compared to default behaviour. The selection of the advantage-based policy gradients is simply justified: Policy gradients calculated with Equation 2.26 yield low variance [42]. Replacing the critic-based policy gradient by a advantage-based one yields

$$\nabla_\phi J(\phi) = \mathbb{E}_{\pi_\phi} \left[ \sum_{u=0}^{T-1} \gamma^u \nabla_\phi \log \pi_\phi(a_u | s_u) A^\pi(s_u, a_u) \right]. \tag{2.26}$$

In practice,  $A^\pi(s_t, a_t)$  is not known and when approximated with  $\hat{Q}^\pi$  and  $\hat{V}^\pi$  in a  $TD(0)$ -context, introduces high bias. Let  $\hat{A}(s_t, a_t) = \hat{Q}^\pi(s_t, a_t) - \hat{V}^\pi(s_t)$  denote the approximated advantage and  $\delta^{\hat{V}}(s_t, a_t) = r(s_t, a_t) + \gamma \hat{V}^\pi(s_{t+1}) - \hat{V}^\pi(s_t)$  the TD-error in terms of the approximated value function for the next state. Then analogous to TD-Learning for value functions, the multi-step approximated advantage can be described

## 2. INTRODUCTION TO REINFORCEMENT LEARNING ALGORITHMS

---

by

$$\begin{aligned}
 \hat{A}^{(\ell)}(s_t, a_t) &:= \sum_{k=0}^{\ell-1} \gamma^k \delta^{\hat{V}^\pi}(s_{t+k}, a_{t+k}), & (2.27) \\
 & s_{t+k} \sim P(\cdot | s_{t+k-1}), a_{t+k} \sim \pi(\cdot | s_{t+k}) \\
 &= -\hat{V}^\pi(s_t) + r(s_t, a_t) + \dots + \gamma^{\ell-1} r(s_{t+\ell-1}, a_{t+\ell-1}) + \\
 & \quad \gamma^\ell \hat{V}^\pi(s_{t+\ell}).
 \end{aligned}$$

Note that if  $\ell = T$ , the bias becomes 0; cutting it at  $\ell$  controls the bias-variance tradeoff. Summing over the advantages for  $(1) \dots (T)$  yields the generalized advantage estimator (GAE)

$$\begin{aligned}
 \hat{A}^{GAE}(s_t, a_t) &:= (1 - \lambda) \left( \hat{A}^{(1)}(s_t, a_t) + \lambda \hat{A}^{(2)}(s_t, a_t) \dots \right) & (2.28) \\
 &= \sum_{k=0}^T (\gamma \lambda)^k \delta^{\hat{V}^\pi}(s_{t+k}, a_{t+k}).
 \end{aligned}$$

Then, replacing the theoretical advantage function in Equation (2.26) with the GAE, we obtain the GAE policy update. By introducing the parameter  $\lambda$ , the bias-variance trade-off can be controlled. A high value of  $\lambda$  reduces bias and increases variance and vice-versa. Notably, if  $\lambda = 1$  and the  $\ell$ -step advantage estimator in (2.27) sums either to termination or infinity, it is equal to the GAE. The advantage function is most commonly used in on-policy algorithms, see Section 2.1.2.6. One such algorithm will be discussed in Section 2.1.6.4.

### 2.1.2.6 On- and Off-Policy Learning

Let  $\mu$  denote the **behaviour policy**, which is **any** policy used to generate experiences, and let  $\pi$  denote the policy to be learned. In on-policy algorithms, the behaviour policy and the policy coincide, i.e.,  $\mu = \pi$ . This is the case because the policy improvement theorem and related convergence guarantees depend on the data distribution matching the policy being evaluated.

In off-policy algorithms (distinct from offline RL, which relies on pre-generated experiences),  $\pi$  is optimized using  $\mu$ -experiences.

Off-policy methods may store experiences in a replay buffer, a mutable set of  $\mu$ -experiences and learn from them. Other approaches, such as importance sampling or divergence correction techniques, address distributional mismatches without relying on replay buffers. In the off-policy algorithms described in this work,  $\pi$  is optimized with experiences from previous versions of the policy, making them  $\mu$ -experiences.

### 2.1.2.7 Multi-Agent Reinforcement Learning

Due to its increasing popularity, the reader shall be informed about the existence of multi-agent reinforcement learning (MARL). Multi-agent methods are applied when an acting system can be divided into smaller, independently or autonomously operating agents. Examples include UAV swarm control for mapping and flexible energy grid management, see [43, 44]. Agents interact within a shared environment. In cooperative settings, agents collaborate towards a shared goal. For instance, the agents control the members of a multi-joint robot, acting interdependently to reach the same goal as in coordinating a robot to pick up an object. In a competitive setting, it is required that only one agent reaches the goal and the agents try to outperform each other. Mixed settings can involve agents with varying objectives, such as controlling chess pieces where each agent suggests moves, with an overarching policy selecting the final action. MARL solutions are generally robust due to the optimized behaviour across multiple agents. They support parallel computation, and experience sharing can lead to more thorough state-space exploration. However, MARL also introduces unique challenges, such as coordinating agent communication, assigning credit in cooperative tasks, and scalability issues.

### 2.1.3 Challenges

#### 2.1.3.1 Exploration-Exploitation Dilemma

Reward associated with transitions must be learned by thoroughly searching the state space and requires testing new actions to gain more knowledge about the environment. *Exploration* involves actively seeking out less-known areas of the state space, enabling the learning process to refine the association between actions and resulting rewards. By contrast, *exploitation* uses the present knowledge of the environment to choose actions believed to yield the highest expected reward, based on current, potentially incomplete information. Reinforcement learning faces the challenge of balancing exploitation and exploration [45]. Theoretically, convergence to an optimal policy requires adherence to greedy-in-the-limit (GLIE) conditions (see Equation 2.17), ensuring that all actions in all states are sufficiently sampled over time. However, this approach is often impractical in real-world applications, especially when dealing with high-dimensional, continuous state/action spaces and sparse reward functions (see Section 2.1.4). Brute-force exploration becomes infeasible in such environments, even with policies gradually increasing in greediness.

In practice, reinforcement learning algorithms employ heuristics such as the  $\epsilon$ -greedy

## 2. INTRODUCTION TO REINFORCEMENT LEARNING ALGORITHMS

---

method, widely used in algorithms like DQN (see Section 2.1.6.1), and stochastic policies as used in SAC or PPO (see Sections 2.2.2 and 2.1.6.4). Heuristics add instability to reinforcement learning algorithms as they rely heavily on parameter initialization and offer no convergence guarantee.

### 2.1.3.2 Overestimation Bias in Q-learning

Overestimation of state-action values is a systemic issue in Q-based reinforcement learning algorithms when function approximators are used. First identified by [46] in 1993, overestimation bias remains challenging to mitigate effectively due to the computational expense of current solutions. To illustrate the overestimation bias, consider the RHS of an aggressive update rule where  $\hat{Q}(s_{t+1}, a_{t+1}) := Q^*(s_{t+1}, a_{t+1}) + \epsilon$ ,  $Q^*$  is the true Q value and  $\epsilon$  is a random noise term,

$$\hat{Q}(s_t, a_t) \leftarrow r(s_t, a_t) + \gamma \max_a \hat{Q}(s_{t+1}, a).$$

Then the bias is

$$\begin{aligned} \Delta(s_t, a_t) &= r(s_t, a_t) + \gamma \max_a \hat{Q}(s_{t+1}, a) - \\ &\quad (r(s_t, a_t) + \gamma \max_a Q^*(s_{t+1}, a)) \mid s_{t+1} \sim P(\cdot \mid s_t, a_t) \\ &= \gamma (\max_a \hat{Q}(s_{t+1}, a) - \max_a Q^*(s_{t+1}, a)) \\ &= \gamma (\max_a [Q^*(s_{t+1}, a) + \epsilon] - \max_a Q^*(s_{t+1}, a)). \end{aligned} \quad (2.29)$$

With  $\mathbb{E}[\epsilon] = 0$ , it holds that the bias  $\mathbb{E}_{P,\epsilon}[\Delta(s_t, a_t)] \geq 0$ . W.l.o.g., if  $\epsilon$  is uniformly distributed in  $[-v, v]$ , the max operator tends to select actions with positive noise values, causing overestimation. This effect can lead to premature convergence to a local maximum with inflated value estimates. A common heuristic to counteract overestimation bias called Double Q-Learning [25] is presented in Section 2.1.5.1. This work presents an algorithm that further mitigates the effect of overestimation at low computational cost.

### 2.1.3.3 Large Spaces and Complexity

*Sample complexity* is the number of environment interactions required to learn an  $\epsilon$ -optimal policy with high probability. PAC-style bounds, under standard assumptions on mixing and reward noise, show that achieving  $\epsilon$ -accuracy in a finite MDP with  $|\mathcal{S}|$  states and  $|\mathcal{A}|$  actions requires at least

$$\tilde{\Omega}\left(\frac{|\mathcal{S}||\mathcal{A}|}{\epsilon^2(1-\gamma)^3}\right)$$

samples [47]. Continuous (or very large) state–action spaces are typically handled by function approximation, for which sample-complexity lower bounds often grow exponentially in the intrinsic dimension of the feature space.

*Computational complexity* is viewed typically under worst-case assumptions to define an upper bound on computational requirements. For instance, analysis may be based on assuming discrete spaces, as continuous spaces can be approximated as large discrete ones. The algorithms presented here achieve optimality via (Howardian) policy iteration (PI), which combines policy evaluation and improvement with discounting, see [48]. Scherrer (2013) showed that the upper bound on the number of policy iterations is  $\mathcal{O}(\frac{|\mathcal{A}|}{1-\gamma} \log \frac{1}{1-\gamma})$ , independent of  $\mathcal{S}$  due to sup-norm contraction properties, see [49]. The per-iteration cost is comprised of the cost of evaluation and improvement. For every state, the Bellman condition (2.19) has to be satisfied, yielding a linear system solvable in polynomial time at  $\mathcal{O}(|\mathcal{S}|^3)$ . Consider a constructive form of the policy improvement step (2.24)

$$\pi(s_t) = \operatorname{argmax}_a r(s_t, a) + \gamma \mathbb{E}_{\pi} [Q^{\pi}(s_{t+1}, a_{t+1})].$$

A naive approach might iterate over all state-action pairs, making the complexity  $\mathcal{O}(|\mathcal{S}|^2|\mathcal{A}|)$ . Thus, the upper bound of the per-iteration complexity for PI becomes

$$\mathcal{O}(|\mathcal{S}|^3 + |\mathcal{S}|^2|\mathcal{A}|).$$

An improved version of PI evaluates the policy partially in every  $k$  steps, reducing the complexity to

$$\mathcal{O}(k|\mathcal{S}|^2 + |\mathcal{S}|^2|\mathcal{A}|).$$

In high-dimensional, continuous environments, such as in the humanoid environment with  $\mathcal{S}^d; d = 374$  and  $\mathcal{A}^d, d = 17$  (Section 3.5), a naive approach to PI is computationally prohibitive. Instead, (deep) neural networks with high-generalization capabilities and sophisticated techniques are used in training for managing complexities. A novel neural-based reinforcement learning algorithm based on the sampling approach is presented in this work that can be applied to large and continuous observation and action spaces.

### 2.1.3.4 Sim-to-Real Gap

Training reinforcement learning models in simulations is often preferred, as it allows faster-than-real-time processing and parallelization. A task in a simulated environment can be solved provided that a model is at hand (and reasonably costly to engineer) and the dynamics are implemented accurately. The former is usually not obtainable.

## 2. INTRODUCTION TO REINFORCEMENT LEARNING ALGORITHMS

---

Simulating contact-rich interactions such as encountered in robotic assembly is not a trivial task and requires highly specialized software. These software systems can be used as a reliable surrogate only within rigid boundaries [50] due to insufficient friction modeling, numerical challenges and unknown parameters. The sim-to-real gap refers to the mismatch between simulated and real-world transition dynamics, where the simulation fails to fully replicate the complexities of the physical environment. Model-based reinforcement learning addresses this by iteratively learning both the control (or policy) and a model of the environment. Other strategies to mitigate the sim-to-real gap include domain randomization [51] and Robust RL [52]. Model-based RL is gaining attention, notably as part of the "Alberta Plan", a vision formulated by Richard Sutton for reaching Artificial General Intelligence (AGI) with reinforcement learning principles as integral components [53]. It is important to remember that results in this work, i.e. control models for various simulated tasks, can not be directly employed to hardware without additional transfer methods.

### 2.1.4 Reward Engineering and Value Alignment

Defining a reward function that elicits desired behaviour in real-life settings is a significant challenge and warrants separate discussion due to the critical issues it can present. This section additionally examines the social and moral dimensions of reward design and introduces the concept of value alignment.

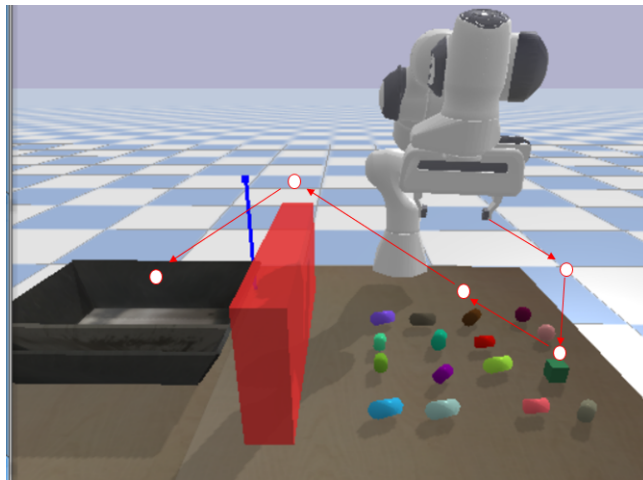
#### 2.1.4.1 Paradigm Challenge

In reinforcement learning, the objective is to identify a policy that maximizes return, Equation (2.11), making the reward function a key determinant of an agent's post-training **behaviour**. In that context, the concept of reward density is introduced. A reward function that provides feedback at each transition is considered a maximal dense reward function, whereas an environment that emits one reward only after the desired behaviour is exhibited has the lowest possible reward density. The density of the reward function significantly impacts learning efficacy. Ideally, an agent learns desired behaviour from sparse feedback, but in practice, sparse reward makes training difficult, often hindering convergence to a local maximum. Crafting a dense reward function requires extensive domain knowledge and human input, which can resemble supervised learning and conflict with autonomous learning goals. Challenges in defining effective guiding signals are ongoing. Solutions include reward-free learning with human guidance through example states and intrinsically motivated reinforcement learning,

where agents explore for the sake of discovery rather than for solving a problem [54, 55]. However, these alternatives are limited by sample inefficiency, dependence on expert knowledge and ambiguity in intrinsic rewards. A promising approach to address reward-based reinforcement learning involves using foundational models, such as LLMs, to automate guidance by either writing reward functions for a specific task or labeling states directly, offering potential for fully automatic reinforcement learning. The next sections discuss issues arising from the paradigm challenge.

### 2.1.4.2 Technical Problem: An example from a Pick and Place Task

Assuming that training with a sparse reward function is computationally tractable, the problem of avoiding unwanted behaviour from a low-resolution guidance signal remains: How to avoid learning unwanted behaviour without over-specifying the feedback? The following example illustrates that the answer to this question is not trivial and simple intuitive solutions can be misleading. Consider a robotic task as illustrated in Figure 2.1.



**Figure 2.1:** PyBullet simulation of a pick and place task with the Franka Emika Panda robot. The robot must pickup a target object and place it inside the tray while avoiding collision with the environment. The robot is allowed to execute 8 collision-free movements before the episode terminates.

The environment includes a target object to be placed in a tray, a robotic agent, and static objects forming the surroundings. The action space consists of spatial points and Eulerian orientation coordinates, see Chapter 4. The agent has up to 8 moves to complete the task, with  $n$  as the current number of moves taken. The episode terminates

## 2. INTRODUCTION TO REINFORCEMENT LEARNING ALGORITHMS

---

prematurely if the agent collides with the environment.

Let the reward function  $r(s, a)$  incorporate penalties and rewards for collision, distance to the object, goal-reaching, and execution time. Define  $v_a(s, a)$  as the robots body position vertices and  $v_e(s, a)$  as static object vertices. Then  $col(s, a) := -c\mathbb{I}(v_a(s, a) \cap v_e(s, a) \neq \emptyset)$ , where  $c$  is a collision penalty and  $\mathbb{I}$  is the indicator function. Let

$$d(\mathbf{p}_a, \mathbf{p}_0; s, a) = \|\mathbf{p}_a(s, a) - \mathbf{p}_0(s)\|$$

denote the distance between the end effector and object position. The full reward function is given by

$$r(s, a) := col(s, a) - d(\mathbf{p}_a, \mathbf{p}_0; s, a) + \mathbb{I}_{s_g}(s, a) + t(s, a)$$

where  $S_g$  is the set of "target states" and  $t(s, a) := -n * (1 - \mathbb{I}_{S_g}(s, a))$  with

$$\mathbb{I}_{s_g}(s, a) = \begin{cases} 1 & \text{if } s \in S_g, \\ 0 & \text{otherwise} \end{cases}$$

can be interpreted as a control cost; it encourages task completion with minimal movements. The function appears intuitive, punishing collisions and excessive moves while rewarding reaching the goal and proximity to the target object. However, it unexpectedly leads the agent, after initial attempts to pick up the object, to collide with the closest edge of the red wall. In attempting to pick up the object, the agent, lacking initial dexterity, frequently collides with the table and incurs a penalty of  $-c$ . Additionally, repeated failed attempts increase the control cost, leading the agent to learn that it can avoid cumulative penalties by terminating the episode quickly. The quickest termination occurs through collision with the nearest object, the wall. Since the agent can not reach the goal state  $s_g \in S_g$  if the target object is not picked up, the reward for  $s_g$  does not affect the agents "suicidal" behaviour. While this issue is easily diagnosed in a simple environment, identifying reward function flaws in more complex settings with numerous objects or agents becomes intractably difficult. Training environments must incorporate guiding signals that align the agent's behaviour with human intentions.

### 2.1.4.3 Moral Problem and Value Alignment

Beyond the technical issues in reinforcement learning, the value alignment problem is critical for ensuring safe human-machine coexistence [56]. Consider a scenario in which a humanoid robot is ordered to fetch a cup of coffee in the kitchen. Let "Fiona" be a pet cat that blocks the entrance to the kitchen. If the robot is solely rewarded for prompt

delivery, it might take harmful actions against Fiona, a behaviour that is impermissible for humans. Then it is said that the machine "ethics" is not aligned with human values [57]. Especially in reinforcement learning, the value alignment problem might be severe, as unwanted behaviour can emerge from seemingly reasonable rewards.

One proposed solution is to enhance guiding signals with a comprehensive, empirically tested ethics library to ensure ethical behaviour when agents interact with humans, animals, or property. Value alignment remains a critical challenge in machine learning, compounded by concerns about the unpredictability of superintelligence – hypothetical entities with intellectual abilities far exceeding human intelligence. Although speculative at the time of writing, prominent researchers are actively exploring methods to control such entities to ensure their goals align with the users intend [58].

### 2.1.5 RL Algorithm Augmentations

In reinforcement learning, different techniques are used to address the challenges described in Section 2.1.3. The following sections detail methods applicable across algorithm classes or all algorithms.

#### 2.1.5.1 Double Q-Learning

The double-Q method addresses the problem discussed in Section 2.1.3.2. The main idea is to decouple the selection of actions from evaluation of the state-action pair [25]. Let  $\theta$  and  $\theta'$  be the parameters of the critic and its target variant (Section 2.1.2.6). Consider the update rule in (2.7) with parameterized approximators and let the estimated label be

$$Y(s_t, a_t) := r(s_t, a_t) + \gamma \max_a Q_\theta(s_{t+1}, a).$$

Then the double-Q method utilizing online and target critics can be expressed as

$$Y^D(s_t, a_t) := r(s_t, a_t) + \gamma Q_{\theta'}(s_{t+1}, \operatorname{argmax}_a Q_\theta(s_{t+1}, a)). \quad (2.30)$$

Note that the  $\operatorname{argmax}$  is used only for target computation; no gradients flow through it.

Within the Actor-Critic (AC) framework (Section 2.1.2.4), one can employ either classical double Q-learning or its clipped variant. The clipped variant maintains two critics and uses  $\min\{Q_\theta, Q_{\theta'}\}$  to define the target Y [29]. This minimum is then used to update the actor. This clipped variant is integrated into C-DSAC (Section 3.2).

## 2. INTRODUCTION TO REINFORCEMENT LEARNING ALGORITHMS

---

### 2.1.5.2 Normalization Techniques

Neural networks perform better with inputs in a consistent range, meaning enhancement in generalization and convergence, see [59, 60]. In reinforcement learning, reward and observation normalization adjust actor and critic inputs to improve stability. One common method is min-max scaling, which rescales feature values  $x_i$  in the vector  $\mathbf{x} \in \mathcal{X}$  (either observation or concatenated observation and action) to  $[0, 1]$ . Let  $x_{min,i} := \min_{\mathbf{x} \in \mathcal{X}} x_i$  and  $x_{max,i} := \max_{\mathbf{x} \in \mathcal{X}} x_i$ , then the rescaled feature value is defined by

$$x'_i = \frac{x_i - x_{min,i}}{x_{max,i} - x_{min,i}}.$$

Min-max normalization does not handle outliers well, since a single extreme value will compress most data into a narrow sub-range. An alternative popular method, Z-Score Standardization, manages outliers by centering features around 0 and adjusting for variability. Let  $\mu_i$  be the mean and  $\sigma_i$  the standard deviation of feature  $i$  over  $\mathcal{X}$ , then the Z-score of a feature value is

$$x'_i = \frac{x_i - \mu_i}{\sigma_i}.$$

When normalization is applied for the activation of neural layers (save the input layer), it is known as Batch Normalization [61]. In reinforcement learning, the available observations in  $\mathcal{X}$  increase with the current iteration number  $N$ . As the number increases and previously unknown parts of the spaces are explored, the data distribution shifts. In the case of max-min scaling, obtaining the new extrema can be achieved by sliding-window min-max or clipped ranges approaches as workarounds. Updating the metrics for the Z-score involves softly updating the mean by  $\mu_{i,N} = (1 - \frac{1}{N})\mu_{i,N-1} + \frac{1}{N}x_i$  or using Welfords methods for the standard deviation. In practice, an exponential moving average,  $\mu_t = (1 - \alpha)\mu_{t-1} + \alpha x_t$ , is used for improved responsiveness and mitigating the fact that  $\frac{1}{N}$  approaches zero. The algorithm presented in Section 3.2 suffers from training instability and application of normalization techniques might mitigate this issue.

### 2.1.5.3 Prioritized Experience Replay

As discussed earlier, off-policy algorithms store the experience tuples in a replay buffer  $\mathcal{B}$ . In the standard case, these experiences are sampled uniformly for learning. Prioritized Experience Replay (PER), however, increases the likelihood of sampling transitions with a higher  $TD$ -error, as these are expected to offer greater learning progress, see [62].

PER defines priority for a transition  $i$  either as  $p_i := |\delta_i| + \epsilon$  or as  $p_i := \frac{1}{\text{rank}(i)}$ , where the function *rank* sorts the experiences according to magnitude of their *TD*-error. Then the sampling probability becomes

$$P(i) = \frac{p_i^\eta}{\sum_k p_k^\eta},$$

where  $\eta \in [0, 1]$  modulates the effect of priority;  $\eta = 0$  reverts to uniform sampling. The algorithm works by sampling a mini-batch according to  $P$  and updating the priorities of the samples in the batch. To make accurate estimates of expected values, estimate updates should match the approximated distribution. This is not the case in PER and thus additional bias is introduced which is adjusted by weighting the *TD*-error with

$$w_i = \left( \frac{1}{|\mathcal{B}|P(i)} \right)^\beta,$$

where  $\beta \in [0, 1]$  determines the degree of bias compensation. It is recommended to anneal  $\beta$  to 1 [62]. PER has shown substantial performance gains in discrete critic-only algorithms and improved performance for the deterministic AC algorithms described in Section 2.1.6.5. The novel reinforcement learning algorithm C-DSAC proposed in this work might benefit from PER as it is relatively unstable in the training phase.

### 2.1.5.4 Initialization and Change of Neural Network Parameters

Neural networks learn by adjusting parameters through backpropagation, see Section 1.2.3. In machine learning, the way these parameters are handled before and during training is crucial. As the name implies, mini-batch stochastic descent is inherently random and might therefore lead to unstable training or prevent convergence altogether. Therefore, neural network parameter initialization matters; the accuracy for the classification tasks described in Section 6.4 can not be achieved if the Glorot uniform initialization scheme [63] is exchanged by a naive Gaussian scheme. Gradient explosion and vanishing are issues where gradients become excessively large or small. The Glorot initialization scheme aims to prevent these phenomena by maintaining the same activation variances across all layers of the neural network, resulting in balanced weights. The forward pass and backpropagation necessitates consideration to the number of both input and output weights. Let  $n_{in}$  and  $n_{out}$  be the number of inputs and outputs of a neural network layer respectively; the mean provides a balanced compromise for both directions. To keep forward and backward signal variances equal for linear activations, the variance should therefore be  $\text{Var}(\mathbf{w}) = \frac{2}{n_{in} + n_{out}}$ . For non-linear activations, [63] suggest either  $\frac{1}{n_{in}}$  or  $\frac{1}{n_{out}}$ , depending on which pass is prioritized. The

## 2. INTRODUCTION TO REINFORCEMENT LEARNING ALGORITHMS

---

variance of a uniform distribution  $\mathcal{U}(a, b)$  is given by  $Var = \frac{1}{12}(b - a)^2$ . Setting one standard deviation as a target and substituting the target in the uniform distribution, the  $k$ -th node of layer  $\ell$  is sampled according to

$$w_{\ell,k} \sim \mathcal{U}\left(-\sqrt{\frac{6}{n_{in} + n_{out}}}, \sqrt{\frac{6}{n_{in} + n_{out}}}\right),$$

with  $\mathbb{E}[\mathcal{U}] = 0$ . As discussed in Section 2.1.5.2, having activations closer to 0 is favorable (if no interfering normalization technique is employed) and therefore the Normal Glorot initialization scheme

$$w_{\ell,k} \sim \mathcal{N}\left(0, \frac{2}{n_{in} + n_{out}}\right)$$

leads to more stable training and higher generalizations. In this work, all neural network-based methods utilize the Normal Glorot scheme for initialization. Besides initialization, the changes of the parameters are also of interest. Especially in RL, where the cost function is approximated, monitoring the gradients is advised. Beyond initial conditions, ongoing stable parameter updates are critical in RL. Some algorithms like PPO (Section 2.1.6.4) have gradient control as their main feature: the advantage function and the ratio of the pre- and post-update action probabilities are used to clip the objective, which indirectly limits the gradient magnitude. Other algorithms like TRPO [27] subject the policy objective per update to a Kullback-Leibler constraint, making sure that the gradients do not steer the model outside of a "trusted region".

### 2.1.5.5 Learning Rate Annealing

The learning rate, or step size for model updates, is an essential yet often underestimated factor in reinforcement learning, with significant implications for convergence [64]. As a matter of fact, the novel confidence-driven concept presented in Section 3.4.1 is related to learning rate control. Learning rate annealing, the systematic reduction of this rate over time, plays a critical role in balancing exploration and exploitation throughout training. At the start of training, a high learning rate promotes exploratory actions and helps the model investigate the state space more broadly, capturing diverse potential strategies. As training progresses and promising regions are identified, a gradual reduction in the learning rate allows for a more careful, refined exploration of these areas, improving convergence to optimal policies. Several annealing strategies facilitate this progression. Linear annealing decreases the learning rate incrementally from an initial value to a minimum threshold, promoting a controlled shift from exploration to exploitation. A simple schedule that decreases the learning rate from  $\beta_0$  to  $\beta_{min}$  over

$H$  steps takes the form

$$\beta_t = \beta_0 - (\beta_0 - \beta_{\min}) \frac{t}{H}.$$

Additionally, periodic annealing methods, such as Cosine Annealing, reset the learning rate to a higher value after it reaches its minimum, allowing the model to periodically reintroduce exploratory behaviour, which can prevent premature convergence to suboptimal policies,

$$\beta_t = \beta_{\min} + \frac{1}{2}(\beta_0 - \beta_{\min})(1 + \cos(\frac{\pi t}{H_{\text{cycle}}})).$$

Exponentially decaying rate,

$$\beta_t = \beta_0 \cdot \exp(-\kappa \frac{t}{H}),$$

where  $\kappa > 0$  sets the decay speed, is often chosen for implementation convenience and to match decaying variance estimates in stochastic gradient methods

Beyond learning rate annealing, modern RL often employs adaptive optimizers (e.g., Adam, RMSProp) to control effective step sizes [65].

### 2.1.6 Relevant Algorithms

To further illustrate the concepts discussed earlier, this section introduces various reinforcement learning algorithms, each emphasizing different aspects of reinforcement learning concepts. Some of these algorithms also serve as the foundation for more advanced methods.

#### 2.1.6.1 Tabular Q-learning

The classical Q-learning algorithm of Watkins is a tabular critic-only (Section 2.1.2.1) reinforcement learning algorithm [66]. This approach involves maintaining a Q-value entry for each possible state-action pair  $(s_t, a_t)$ , representing the expected return for that pair. The Q-value entries are updated iteratively utilizing the Bellman optimality operator (2.6) that performs an  $TD(0)$ -update. The goal is to find a function that satisfies equation (2.5) for all state-action pair entries. During training, the action is selected by an  $\epsilon$ -greedy policy, i.e. at each time step, a random action is selected with probability  $\epsilon$ , and the greedy option (2.3) according to the best current Q-estimate with probability  $(1 - \epsilon)$ . This ensures that each action is eventually taken, allowing the Q-values to converge to optimal values under the GLIE condition described by (2.17). Typically,  $\epsilon$  is initialized at 1, prioritizing exploration in early episodes, and linearly reduced as learning progresses. Tabular Q-learning is computationally feasible only for

## 2. INTRODUCTION TO REINFORCEMENT LEARNING ALGORITHMS

---

small or discrete state-action spaces, as it requires storage for each pair individually. For larger or continuous spaces, function approximation or other methods must be employed.

### 2.1.6.2 Deep Q-Learning

Following a period of stagnation in reinforcement learning research, DeepMind revitalized interest in this field by presenting an AI system that exceeded human performance in Atari 2600 games [10, 67]. This breakthrough was achieved using a deep neural network within the Q-learning framework, with MLPs approximating the Q-function; CNNs processed the raw game pixels for extraction of abstract features from the high-dimensional state-space and were trained end-to-end. In Algorithm 2.1, the CNN component is represented by  $\psi$ . The objective is to minimize the mean squared error between approximated targets and predicted Q-values for the current state-action pairs. Assuming a neural network approximator for the Q-function parameterized by  $\theta$ , the loss function becomes

$$J_Q(\theta) = \mathbb{E}_{\substack{s_t \sim d_\mu^t, \\ a_t \sim \mu(\cdot | s_t)}} \left[ \frac{1}{2} (\mathcal{T}Q(s_t, a_t) - Q_\theta(s_t, a_t))^2 \right], \quad (2.31)$$

where  $d_\mu^t(s) = \frac{1-\gamma}{1-\gamma^T} \sum_{t=0}^{T-1} \gamma^t \Pr(s_t = s | \mu)$  is the state visitation distribution induced by the behavior policy  $\mu$  in a finite horizon.

In practice, this expectation is approximated by sampling mini-batches of experience tuples  $(s_t, a_t, r_t, s_{t+1})$ , with  $r_t := r(s_t, a_t)$ , from a replay buffer  $\mathcal{B}$  (see Section 2.1.2.6) and computing the sample mean. This means also that the discounting in state-visitation is disregarded. For a particular  $s_{t+1}$ , the sample Bellman operator becomes

$$\hat{\mathcal{T}}Q(s_t, a_t) := r(s_t, a_t) + \gamma \max_a Q(s_{t+1}, a).$$

The Q-network parameters,  $\theta$ , are updated based on the negative gradient of the loss

$$\nabla_\theta J_Q(\theta) = \mathbb{E}_{(s_t, a_t, r_t, s_{t+1}) \sim \mathcal{B}} [(\hat{\mathcal{T}}Q(s_t, a_t) - Q_\theta(s_t, a_t)) \nabla_\theta Q_\theta(s_t, a_t)]. \quad (2.32)$$

During training, these tuples containing single transitions are sampled uniformly from the replay buffer, reducing temporal correlations between the experiences. Given that DQN serves as the foundation for many Q-based algorithm in reinforcement learning, the complete DQN procedure is outlined in Algorithm 2.1. To emphasize the separation between policy evaluation and improvement in off-policy algorithms, the two phases are presented as distinct loops in Algorithm 2.1, although in practice each

environment step is typically followed by a single gradient update. This convention is adopted consistently throughout the manuscript. Only small action spaces can be managed by DQN, thus rendering it suitable only for discrete (and small) spaces. Among Atari 2600 benchmarks, *Venture* has one of the largest discrete action sets (18 actions). Section 2.1.6.5 introduces a Q-based AC algorithm handling continuous actions.

---

**Algorithm 2.1** Deep Q-learning (with observation encoder  $\psi$ )

---

```

1: Initialize replay memory  $\mathcal{B}$  with capacity  $|\mathcal{B}|$ 
2: Initialize action-value function  $Q$  with weights  $\theta$ 
3: Initialize observation encoder  $\psi$ 
4: Initialize learning parameter  $\beta_Q$ ,
5: Initialize exploration  $\epsilon = 1$  and decay parameter  $\beta_\epsilon$ 
6: repeat
7:   Receive initial (raw) observation  $s_t = s_1$ 
8:   for each environment step do
9:     With probability  $\epsilon$ , select a random action  $a_t$ 
10:    Otherwise, select  $a_t = \operatorname{argmax}_a Q_\theta(\psi(s_t), a)$ 
11:    Observe transition  $s_{t+1} \sim P(\cdot | s_t, a_t)$ 
12:    Observe reward  $r_t$ 
13:     $\mathcal{B} \leftarrow \mathcal{B} \cup \{(s_t, a_t, r_t, s_{t+1})\}$ 
14:     $s_t \leftarrow s_{t+1}$ 
15:   for each gradient step do
16:     Sample a random minibatch of  $N$  transitions  $(s_j, a_j, r_j, s_{j+1})$  from  $\mathcal{B}$ 
17:     Set  $y_j = \begin{cases} r_j & \text{for terminal } s_{j+1} \\ \hat{\mathcal{T}}Q_\theta(\psi(s_j), a_j) & \text{for non-terminal } s_{j+1} \end{cases} \quad j = 1, \dots, N$ 
18:     Calculate approximated loss  $\hat{J}_Q$  from arithmetic mean  $y$ 
19:      $\theta \leftarrow \theta + \beta_Q \hat{\nabla}_\theta J_Q$  from approximated loss
20:     Update encoder parameters via same gradient
21:      $\epsilon \leftarrow \max(\epsilon - \beta_\epsilon, 0)$ 
22:   until Stopping criterion
23: return  $Q_\theta$ 

```

---

### 2.1.6.3 REINFORCE

REINFORCE, short for "Reward Increment = Nonnegative Factor  $\times$  Offset Reinforcement  $\times$  Characteristic Eligibility" is a foundational on-policy algorithm that performs

## 2. INTRODUCTION TO REINFORCEMENT LEARNING ALGORITHMS

---

a direct policy search by optimizing the policy objective (2.11) for the parameters of a stochastic policy [68]. As one of the earliest policy gradient methods, REINFORCE has played a critical role in establishing techniques used in more advanced policy optimization algorithms. The REINFORCE algorithm operates based on Monte Carlo sampling, where updates occur only after an episode has terminated - mirroring  $TD(1)$ -updates. Let  $G_t := \sum_{k=t}^{T-1} \gamma^{k-t} r(s_t, a_t)$ , then the update gradient is expressed as

$$\nabla_{\phi} J(\phi) = \mathbb{E}_{\pi_{\phi}, P} \left[ \sum_{t=0}^{T-1} \gamma^t \nabla_{\phi} \log \pi_{\phi}(a_t | s_t) G_t \right].$$

The post-update parameters are computed using the gradient

$$\phi_{t+1} = \phi_t + \beta \mathbb{E}_{\pi_{\phi}, P} \left[ \sum_{t=0}^{T-1} \gamma^t \nabla_{\phi} \log \pi_{\phi}(a_t | s_t) G_t \right].$$

The update rule for REINFORCE is obtained by substituting the term inside the expectation into (2.12) for  $\nabla_{\phi} J(\phi)$ . REINFORCE, being a Monte Carlo method, often exhibits high variance in gradient estimates, which can hinder stable learning. A common modification involves introducing a baseline value — typically a state value function  $V^{\pi}$  — to reduce variance, leading to a gradient formulation

$$\nabla_{\phi} J(\phi) = \mathbb{E}_{\pi_{\phi}, P} \left[ \sum_{t=0}^{T-1} \gamma^t \nabla_{\phi} \log \pi_{\phi}(a_t | s_t) (G_t - V^{\pi}(s_t)) \right].$$

With hindsight to (2.27) and assuming the advantage function's independence of  $\phi$ , the above equation can be re-written more compactly in terms of the estimated advantage function

$$\nabla_{\phi} J(\phi) = \mathbb{E}_{\pi_{\phi}, P} \left[ \sum_{t=0}^{T-1} \gamma^t \nabla_{\phi} \log \pi_{\phi}(a_t | s_t) \hat{A}^T(s_t, a_t) \right]. \quad (2.33)$$

”Baselines” are important concepts that are used in various methods. The full procedure with a return-value difference is described in Algorithm 2.2 [69].

In practice, the parameters are repeatedly updated according to the update rule

$$\phi \leftarrow \phi + \beta \gamma^t \nabla_{\phi} \log \pi_{\phi}(a_t | s_t) b_t(s_t),$$

where  $b_t$  is either the return-value difference or the advantage function.

---

**Algorithm 2.2** REINFORCE with Learned V-Baseline

---

- 1: Initialize policy  $\pi_\phi$ , baseline  $V_\theta$ , learning parameters  $\beta_\pi$  and  $\beta_V$
  - 2: **for** each episode **do**
  - 3:     Sample  $\tau = \{s_0, a_0, \dots, s_{T-1}, a_{T-1}\} \sim \pi_\phi$
  - 4:     Compute returns  $G_t$  for  $t = 0, \dots, T - 1$
  - 5:     **for** each  $t$  **do**
  - 6:          $\phi \leftarrow \phi + \beta_\pi \gamma^t \nabla_\phi \log \pi_\phi(a_t | s_t) (G_t - V_\theta(s_t))$
  - 7:          $\theta \leftarrow \theta - \beta_V \nabla_\theta (V_\theta(s_t) - G_t)^2$
- 

**2.1.6.4 PPO**

Proximal Policy Optimization (PPO), an Actor-Critic (AC) algorithm, is a popular on-policy approach that leverages stochastic policies and can incorporate the Generalized Advantage Estimator (GAE) [28]. Noted for its robustness and relative insensitivity to hyperparameter tuning, PPO is straightforward to implement and can solve all environments examined in this work. However, while effective, it is generally less sample-efficient compared to off-policy methods and does not achieve the same level of performance as the novel C-DSAC algorithm presented in Section 3.2. PPO employs a Conservative Policy Iteration (CPI) objective, designed to stabilize learning by implicitly bounding policy updates through loss constraints. Using GAE as discussed in Section 2.1.2.5, the CPI actor objective to be maximized is defined as

$$J^{CPI}(\phi) := \mathbb{E}_{\pi_{\phi_{old}, P}} \left[ \frac{\pi_\phi(a_t | s_t)}{\pi_{\phi_{old}}(a_t | s_t)} \hat{A}^{GAE}(s_t, a_t) \right], \quad (2.34)$$

where  $\kappa_{t, \phi} := \kappa_\phi(s_t, a_t) := \frac{\pi_\phi(a_t | s_t)}{\pi_{\phi_{old}}(a_t | s_t)}$  measures the post-update parameter significance. The central feature of PPO is its clipped objective that balances the incentive to improve the policy with a safeguard against excessive updates towards positive advantages

$$J^{(CLIP)}(\phi) := \mathbb{E}_{\pi_{\phi_{old}, P}} \left[ \min \left( \kappa_{t, \phi} \hat{A}^{GAE}(s_t, a_t), \text{clip}(\kappa_{t, \phi}, 1 - \epsilon, 1 + \epsilon) \hat{A}^{GAE}(s_t, a_t) \right) \right]. \quad (2.35)$$

If the advantage is positive with a high probability ratio, the second *min*-argument is selected for the update; in case of a low probability ratio, the left term is chosen. The mechanism works in the opposite direction for negative advantages: If the probability is low, the clipped term is used and otherwise the left is selected. By taking the minimum of the clipped and unclipped terms, PPO ensures less adaptation to positive experiences

## 2. INTRODUCTION TO REINFORCEMENT LEARNING ALGORITHMS

---

and more adaptation to negative ones, making it a "pessimistic" algorithm. In practice, PPO-agents learn from trajectory segments. This necessitates an  $\ell$ -step GAE

$$\hat{A}^{GAE(\ell)}(s_t, a_t) := \sum_{k=0}^H (\gamma\lambda)^k \delta^{\hat{V}^\pi}(s_{t+k}, a_{t+k})$$

which aggregates the  $\gamma$ - and  $\lambda$ -discounted temporal-difference errors up to the horizon  $H = \min\{\ell - 1, T - t - 1\}$  within each episode. The parameter  $\ell$  manages the bias-variance tradeoff. Furthermore, state-value estimates must be obtained at each time step, approximated by a parameterized critic as introduced in (2.9). The critic itself can be updated with the methods described in Section 2.1.2.2. Because the  $\ell$ -step GAE is employed, the target value can be expressed in terms of

$$V_{tar}(s_t) := V^\pi(s_t) + \hat{A}^{GAE(\ell)}(s_t, a_t).$$

Given the state-visitation distribution  $d_{\pi_{\phi_{old}}}^\gamma(s) = (1 - \gamma) \sum_{t=0}^{\infty} \gamma^t \Pr(s_t = s \mid \pi_{\phi_{old}})$  induced by the old policy, the objective of the critic to be minimized becomes

$$J_V(\theta) := \mathbb{E}_{s_t \sim d_{\pi_{\phi_{old}}}^\gamma} [(V_\theta(s_t) - V_{tar}(s_t))^2]. \quad (2.36)$$

In PPO, exploration is further encouraged by an entropy term, yielding a modified policy objective

$$J^{(CLIP+H)} = J^{(CLIP)} + \alpha \mathbb{E}_{s_t \sim d_{\pi_{\phi_{old}}}^\gamma} [\mathcal{H}(\pi_\phi(\cdot \mid s_t))], \quad (2.37)$$

akin to the entropy-augmented approach in Soft Actor-Critic (SAC) and the newly developed C-DSAC algorithm. The complete PPO framework for multiple actors is outlined in Algorithm 2.3, where  $N$  actors generate trajectories, calculate the  $\ell$ -step GAE, and optimize both actor and critic with the clipping and entropy-augmented losses.

### 2.1.6.5 DDPG

Deep Deterministic Policy Gradient (DDPG) [26] is an off-policy, actor-critic reinforcement learning algorithm that extends Deep-Q Learning (Section 2.1.6.2) for continuous action spaces by leveraging ideas from Deterministic Policy Gradient (DPG) [70]. In DDPG, the policy is deterministic, mapping states directly to actions via a function  $\eta : \mathcal{S} \mapsto \mathcal{A}$ . This deterministic nature separates the exploration problem from the

---

**Algorithm 2.3** PPO, Actor-Critic Style
 

---

- 1: Initialize  $\phi$ ,  $\theta$  and learning rates  $\beta_\pi$ ,  $\beta_V$
  - 2: **for** episode = 1,  $M$  **do**
  - 3:     **for** actor = 1,  $N$  **do**
  - 4:         Run policy  $\pi_{\phi_{old}}$  in environment for  $T$  timesteps
  - 5:         Compute advantage estimates  $\hat{A}(s_t, a_t), \dots, \hat{A}(s_{t+l}, a_{t+1})$
  - 6:          $\phi \leftarrow \phi + \beta_\pi \nabla_\phi J^{CLIP+H}$  # Gradient based on (2.37)
  - 7:          $\theta \leftarrow \theta - \beta_V \nabla_\theta J_V$  # Gradient based on (2.36)
  - 8: **return**  $V_\theta, \pi_\phi$
- 

learning process. In order to encourage exploration, noise is added to the action of the deterministic policy

$$\tilde{\eta}(s_t) := \eta_\phi(s_t) + \epsilon \mid \epsilon \sim \mathcal{N}.$$

The authors in [26] use the Ornstein-Uhlenbeck process to sample temporally correlated noise for exploration. Later implementations have shown that the rather complex noise function can be replaced by a Gaussian distribution without compromising the exploration quality [71]. As an off-policy algorithm, DDPG learns from experiences generated by a behaviour policy, which is stochastic due to the added exploration noise. Analogous to (2.19), the state-action value function  $Q^\eta(s_t, a_t)$  for a deterministic policy satisfies the Bellman equation

$$Q^\eta(s_t, a_t) = r(s_t, a_t) + \gamma \mathbb{E}_{s_{t+1} \sim P(\cdot | s_t, a_t)} [Q^\eta(s_{t+1}, \eta(s_{t+1}))]. \quad (2.38)$$

The critic is updated by minimizing a loss function based on the Bellman operator for the deterministic policy

$$J_Q(\theta) := \mathbb{E}_{\substack{s_t \sim d_\mu^\eta, \\ a_t \sim \mu(\cdot | s_t)}} \left[ \frac{1}{2} (\mathcal{T}^\eta Q(s_t, a_t) - Q_\theta(s_t, a_t))^2 \right], \quad (2.39)$$

where

$$\mathcal{T}^\eta Q(s_t, a_t) := r(s_t, a_t) + \gamma \mathbb{E}_{s_{t+1} \sim P(\cdot | s_t, a_t)} [Q(s_{t+1}, \eta(s_{t+1}))]$$

and  $d_\mu^\eta(s) = \frac{1-\gamma}{1-\gamma^T} \sum_{t=0}^{T-1} \gamma^t \Pr(s_t = s \mid \mu)$  is the state-visitation distribution induced by the stochastic behavior policy  $\mu$ .

Target networks parameterized by  $\bar{\phi}$  and  $\bar{\theta}$  are employed to stabilize training, as described in Section 2.1.2.6. The Bellman updates are performed with estimates of the target networks from a particular  $s_{t+1}$ , thus the Bellman operator used in practice is

$$\hat{\mathcal{T}}^\eta Q(s_t, a_t) := r(s_t, a_t) + \gamma Q_{\bar{\theta}}(s_{t+1}, \eta_{\bar{\phi}}(s_{t+1})).$$

## 2. INTRODUCTION TO REINFORCEMENT LEARNING ALGORITHMS

---

Analogous to the DQN-loss in Eq. (2.31), the expectation is approximated by repeatedly sampling mini-batches of experience tuples  $(s_t, a_t, r_t, s_{t+1})$ , with  $r_t := r(s_t, a_t)$ , from a replay buffer  $\mathcal{B}$ . The negative gradient of the critic loss is thus computed by

$$\nabla_{\theta} J_Q(\theta) = \mathbb{E}_{(s_t, a_t, r_t, s_{t+1}) \sim \mathcal{B}} \left[ (\hat{\mathcal{T}}^{\eta} Q(s_t, a_t) - Q_{\theta}(s_t, a_t)) \nabla_{\theta} Q_{\theta}(s_t, a_t) \right]. \quad (2.40)$$

The actor is optimized using deterministic policy gradients derived by [70].

Assuming uniform sampling from the replay buffer, the state visitation *discounting* is not accounted for,

$$\nabla_{\phi} J_{\pi}(\phi) = \mathbb{E}_{s_t \sim \mathcal{B}} [\nabla_a Q_{\theta}(s_t, \eta_{\phi}(s_t)) \nabla_{\phi} \eta_{\phi}(s_t)]. \quad (2.41)$$

This formulation can be intuitively understood by considering the deterministic policy as a special case of a stochastic policy (2.21) with variance (mass point distribution). In practice, DDPG uses deep neural networks to approximate both the deterministic policy and the Q-function, with parameters  $\phi$  and  $\theta$ , respectively. Due to the limitations of sample-based estimation and the fixed input size of neural networks, gradients  $\nabla_{\phi} J_{\pi}(\phi)$  and  $\nabla_{\theta} J_Q(\theta)$  are approximated from mini-batches of data sampled from the replay buffer. The original DDPG paper proposed an Ornstein-Uhlenbeck process for temporally correlated noise during exploration. However, later implementations demonstrated that simpler Gaussian noise often suffices without degrading exploration performance. The complete DDPG scheme is summarized in Algorithm 2.4

## 2.2 Maximum-Entropy Reinforcement Learning

Novel concepts in this work are rooted in maximum-entropy reinforcement learning, a framework known for its robustness and superior performance in complex, high-dimensional tasks. The notation in this section reflects the theoretical basis in the actor-critic scheme.

### 2.2.1 Problem Formulation

In the maximum-entropy setting, the objective in (2.1) is augmented with the policy entropy to improve exploration

## 2.2 Maximum-Entropy Reinforcement Learning

---

### Algorithm 2.4 DDPG Algorithm

---

- 1: Initialize  $\theta, \bar{\theta}, \phi, \bar{\phi}$
  - 2: Initialize learning parameters  $\beta_Q, \beta_\pi$
  - 3:  $\mathcal{B} \leftarrow \emptyset$
  - 4: **repeat**
  - 5:     Initialize an Ornstein-Uhlenbeck process  $\mathcal{N}$
  - 6:     Receive initial observation state  $s_1$
  - 7:     **for each** environment step **do**
  - 8:          $a_t := \eta_\phi(s_t) + \epsilon \mid \epsilon \sim \mathcal{N}$
  - 9:         Observe transition  $s_{t+1} \sim P(\cdot \mid s_t, a_t)$
  - 10:         Observe reward  $r_t$
  - 11:          $\mathcal{B} \leftarrow \mathcal{B} \cup \{(s_t, a_t, r_t, s_{t+1})\}$
  - 12:     **for each** gradient step **do**
  - 13:         Sample a random minibatch of  $N$  transitions  $(s_i, a_i, r_i, s_{i+1})$  from  $\mathcal{B}$
  - 14:         Calculate approximated gradients  $\hat{\nabla}_\theta J_Q, \hat{\nabla}_\phi J_\pi$  from samples
  - 15:          $\theta \leftarrow \theta + \beta_Q \hat{\nabla}_\theta J_Q$  (Mini-batch variant of Eq. (2.40))
  - 16:          $\phi \leftarrow \phi + \beta_\pi \hat{\nabla}_\phi J_\pi$  (Mini-batch variant of Eq. (2.41))
  - 17:         Update target networks:  $\bar{\theta} \leftarrow \tau\theta + (1 - \tau)\bar{\theta}; \bar{\phi} \leftarrow \tau\phi + (1 - \tau)\bar{\phi}$
  - 18:     **until** Stopping criterion
  - 19: **return**  $Q_\theta, \pi_\phi$
- 

$$\begin{aligned}
 J_h(\pi) &:= \mathbb{E}_{\substack{s_0 \sim d_0 \\ a_t \sim \pi(\cdot \mid s_t) \forall t=0, \dots, T-1, \\ s_{t+1} \sim P(\cdot \mid s_t, a_t) \forall t=0, \dots, T-1}} \left[ \sum_{t=0}^{T-1} \gamma^t (r(s_t, a_t) + \alpha \mathcal{H}(\pi(\cdot \mid s_t))) \right] \\
 &:= \mathbb{E}_{\pi, P} \left[ \sum_{t=0}^{T-1} \gamma^t (r(s_t, a_t) + \alpha \mathcal{H}(\pi(\cdot \mid s_t))) \right],
 \end{aligned} \tag{2.42}$$

and the aim is to find an optimal policy  $\pi^* = \operatorname{argmax}_{\pi \in \Pi} J_h(\pi)$  [30]. The parameter  $\alpha$  represents the impact of policy stochasticity on the state-action value.

Assuming discrete spaces, Eq. (2.42) can be written more explicitly as

$$J_h(\pi) = \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} \sum_{t=0}^{T-1} \Pr(s_t = s \mid \pi, P, s_0 \sim d_0) \pi(a \mid s) \gamma^t (r(s, a) + \alpha \mathcal{H}(\pi(\cdot \mid s_t))). \tag{2.43}$$

The entropy is then calculated over the categorical distribution of the policy with

## 2. INTRODUCTION TO REINFORCEMENT LEARNING ALGORITHMS

---

$\mathcal{H}(\pi(\cdot | s_t)) := -\sum_i^{|\mathcal{A}|} p_i \log p_i$  (and  $\mathcal{H}(\pi(\cdot | s_t)) := \mathbb{E}_{a_t \sim \pi(\cdot | s_t)}[-\log \pi(a_t | s_t)]$  for continuous actions).

### 2.2.2 Soft Actor-Critic Framework

Analogous to the definition of the Bellman operator for standard reinforcement learning, a modified Bellman operator  $\mathcal{T}_h^\pi$  is defined in the maximum entropy setting. The soft Bellman operator  $\mathcal{T}_h^\pi$  is defined as

$$\mathcal{T}_h^\pi Q(s_t, a_t) := r(s_t, a_t) + \gamma \mathbb{E}_{\substack{s_{t+1} \sim P(\cdot | s_t, a_t) \\ a_{t+1} \sim \pi(\cdot | s_{t+1})}} [Q(s_{t+1}, a_{t+1}) + \alpha \mathcal{H}(\pi(\cdot | s_{t+1}))]. \quad (2.44)$$

The soft state-action value is given by

$$V_h^\pi(s_t) := \mathbb{E}_{a_t \sim \pi(\cdot | s_t)} [Q^\pi(s_t, a_t) + \alpha \mathcal{H}(\pi(\cdot | s_t))].$$

As shown in the Lemma 2.6, the operator  $\mathcal{T}_h^\pi$  is a  $\gamma$ -contraction, analogous to  $\mathcal{T}^\pi$ .

Soft policy improvement involves the information projection and update towards the exponential of the new Q-value,

$$\pi_{k+1}(\cdot | s_t) = \operatorname{argmin}_{\pi(\cdot | s_t)} d_{kl} \left( \pi(\cdot | s_t) \parallel \frac{\exp(\frac{1}{\alpha} Q^{\pi_k}(s_t, \cdot))}{W^{\pi_k}(s_t)} \right), \quad (2.45)$$

where  $\alpha > 0$ ,  $d_{kl}$  is the Kullback-Leibler divergence function and  $W^{\pi_k}(s_t) = \sum_i^{|\mathcal{A}|} \exp(\frac{1}{\alpha} Q^{\pi_k}(s_t, a_i))$  denotes the partition function normalizing the distribution.  $W^{\pi_k}(s_t)$  can be ignored since it does not contribute to the new policy.

### 2.2.3 Contraction, Improvement and Iteration in the Maximum-Entropy Setting

Convergence towards the maximum of (2.42) can be shown in soft policy iteration, which consists of soft policy evaluation under  $\mathcal{T}_h^\pi$  and soft policy improvement.

#### 2.2.3.1 Soft Policy Evaluation

**Lemma 2.6.** Let  $Q_1(s, a)$  and  $Q_2(s, a)$  be two bounded state-action value functions. Assume  $|\mathcal{A}| < \infty$ ,  $|\mathcal{S}| < \infty$ ,  $0 \leq \gamma < 1$  and  $\alpha \geq 0$ . Then the soft Bellman operator  $\mathcal{T}_h^\pi$  according to (2.44) is a  $\gamma$ -contraction

$$\|\mathcal{T}_h^\pi Q_1 - \mathcal{T}_h^\pi Q_2\|_\infty \leq \gamma \|Q_1 - Q_2\|_\infty. \quad (2.46)$$

## 2.2 Maximum-Entropy Reinforcement Learning

---

*Proof.* By linearity of expectation, the entropy term can be separated from the Q-function. Define

$$r_h(s_t, a_t) := r(s_t, a_t) + \gamma \mathbb{E}_{s_{t+1} \sim P(\cdot | s_t, a_t)} [\alpha \mathcal{H}(\pi(\cdot | s_{t+1}))].$$

Note that

$$\mathcal{T}_h^\pi Q(s_t, a_t) = r_h(s_t, a_t) + \gamma \mathbb{E}_{\substack{s_{t+1} \sim P(\cdot | s_t, a_t), \\ a_{t+1} \sim \pi(\cdot | s_{t+1})}} [Q(s_{t+1}, a_{t+1})]$$

and since entropy is bounded for discrete actions,  $\mathcal{T}_h^\pi$  acts on a Banach space. Then the proof in Lemma 2.3 can be applied such that

$$\|\mathcal{T}_h^\pi Q_1 - \mathcal{T}_h^\pi Q_2\|_\infty \leq \gamma \|Q_1 - Q_2\|_\infty. \quad (2.47)$$

Thus if  $Q_{k+1} := \mathcal{T}_h^\pi Q_k$ , the series  $\{Q_k\}$  converges to  $Q_h^\pi$ , i.e.  $\|Q_k - Q_h^\pi\|_\infty \rightarrow 0$  as  $k \rightarrow \infty$ .  $\square$

### 2.2.3.2 Soft Policy Improvement

It must be shown that (2.45) improves the policy.

**Lemma 2.7.** (One-step soft improvement). Let  $s_t$  be fixed and

$$Q_h^\pi(s_t, a_t) := r(s_t, a_t) + \gamma \mathbb{E}_{\substack{s_{t+1} \sim P(\cdot | s_t, a_t), \\ a_{t+1} \sim \pi(\cdot | s_{t+1})}} [Q_h^\pi(s_{t+1}, a_{t+1}) - \alpha \log \pi(a_{t+1} | s_{t+1})].$$

If the policy is updated according to (2.45),  $Q_h$  is bounded,  $\alpha > 0$ ,  $0 \leq \gamma < 1$  and  $|\mathcal{A}| < \infty$ , then

$$\mathbb{E}_{a_t \sim \pi_{k+1}(\cdot | s_t)} [Q_h^{\pi_k}(s_t, a_t) - \alpha \log \pi_{k+1}(a_t | s_t)] \geq \mathbb{E}_{a_t \sim \pi_k(\cdot | s_t)} [Q_h^{\pi_k}(s_t, a_t) - \alpha \log \pi_k(a_t | s_t)].$$

*Proof.*

$$\begin{aligned} \pi_{k+1}(\cdot | s_t) &= \operatorname{argmin}_{\pi \in \Pi} D_{KL} \left( \pi(\cdot | s_t) \parallel \frac{\exp(\frac{1}{\alpha} Q_h^{\pi_k}(s_t, \cdot))}{W^{\pi_k}(s_t)} \right) \\ &= \operatorname{argmin}_{\pi \in \Pi} \mathbb{E}_{a_t \sim \pi(\cdot | s_t)} \left[ \log \left( \frac{\pi(a_t | s_t) W^{\pi_k}(s_t)}{\exp \frac{1}{\alpha} Q_h^{\pi_k}(s_t, a_t)} \right) \right] \\ &= \operatorname{argmin}_{\pi \in \Pi} \mathbb{E}_{a_t \sim \pi(\cdot | s_t)} \left[ \log \pi(a_t | s_t) - \frac{1}{\alpha} Q_h^{\pi_k}(s_t, a_t) + \log W^{\pi_k}(s_t) \right]. \end{aligned}$$

Due to minimization, it holds  $\mathbb{E}_{a_t \sim \pi_{k+1}(\cdot | s_t)} [\log \pi_{k+1}(a_t | s_t) - \frac{1}{\alpha} Q_h^{\pi_k}(s_t, a_t) + \log W^{\pi_k}(s_t)] \leq \mathbb{E}_{a_t \sim \pi_k(\cdot | s_t)} [\log \pi_k(a_t | s_t) - \frac{1}{\alpha} Q_h^{\pi_k}(s_t, a_t) + \log W^{\pi_k}(s_t)]$  and since the partition function does not depend on the action, it follows

$$\mathbb{E}_{a_t \sim \pi_{k+1}(\cdot | s_t)} [Q_h^{\pi_k}(s_t, a_t) - \alpha \log \pi_{k+1}(a_t | s_t)] \geq \mathbb{E}_{a_t \sim \pi_k(\cdot | s_t)} [Q_h^{\pi_k}(s_t, a_t) - \alpha \log \pi_k(a_t | s_t)].$$

## 2. INTRODUCTION TO REINFORCEMENT LEARNING ALGORITHMS

---

□

The soft Q-value is expanded to show that  $\pi_{k+1}$  is indeed an improvement on  $\pi_k$

**Lemma 2.8.** (Global soft policy improvement) Assume  $|\mathcal{A}| < \infty$ ,  $\alpha > 0$ ,  $0 \leq \gamma \leq 1$  and bounded  $Q_h^\pi$ . If  $\pi_{k+1}$  is obtained according to the soft policy improvement in (2.45), then

$$Q_h^{\pi_k}(s_t, a_t) \leq Q_h^{\pi_{k+1}}(s_t, a_t) \quad \forall (s_t, a_t) \in \mathcal{S} \times \mathcal{A},$$

when applying Lemma 2.7 pointwise at each successor state.

*Proof.* Analogous to Lemma 2.4, the proof relies on expanding the soft Q-function in the Bellman equation repeatedly with actions of the improved policy. For convenience, define  $\pi' := \pi_{k+1}$ ,  $\pi := \pi_k$  and  $\mathbb{E}_{\pi', P, k}[\cdot | s_t] := \mathbb{E}_{a_{t:k-1} \sim \pi'(\cdot | s_{t:k-1}), s_{t+1:k} \sim P(\cdot | s_{t:k-1}, a_{t:k-1})}[\cdot | s_t]$ . Applying the one-step improvement from Lemma 2.7) at each successor state yields

$$\begin{aligned} V_h^\pi(s_t) &\leq \mathbb{E}_{\pi', P, t+1} [Q_h^\pi(s_t, a_t) - \alpha \log \pi'(a_t | s_t) | s_t] \\ &= \mathbb{E}_{\pi', P, t+1} [r(s_t, a_t) - \alpha \log \pi'(a_t | s_t) + \gamma V_h^\pi(s_{t+1}) | s_t] \\ &\leq \mathbb{E}_{\pi', P, t+2} [r(s_t, a_t) - \alpha \log \pi'(a_t | s_t) + \gamma(Q_h^\pi(s_{t+1}, a_{t+1}) - \alpha \log \pi'(a_{t+1} | s_{t+1})) | s_t] \\ &= \mathbb{E}_{\pi', P, t+2} [r(s_t, a_t) - \alpha \log \pi'(a_t | s_t) + \gamma(r(s_{t+1}, a_{t+1}) - \alpha \log \pi'(a_{t+1} | s_{t+1})) + \\ &\quad \gamma^2 V_h^\pi(s_{t+2}) | s_t] \\ &\vdots \\ &\leq \mathbb{E}_{\pi', P, T-1} [r(s_t, a_t) - \alpha \log \pi'(a_t | s_t) + \gamma(r(s_{t+1}, a_{t+1}) - \alpha \log \pi'(a_{t+1} | s_{t+1})) + \dots \\ &\quad \dots + \gamma^{T-1}(r(s_{T-1}, a_{T-1}) - \alpha \log \pi'(a_{T-1} | s_{T-1})) | s_t = s] \\ &= V_h^{\pi'}(s_t). \end{aligned}$$

Thus for any  $s_t$

$$\begin{aligned} Q_h^{\pi_k}(s_t, a_t) &= \mathbb{E}_{s_{t+1} \sim P(\cdot | s_t, a_t)} [r(s_t, a_t) + \gamma V_h^{\pi_k}(s_{t+1})] \\ &\leq \mathbb{E}_{s_{t+1} \sim P(\cdot | s_t, a_t)} [r(s_t, a_t) + \gamma V_h^{\pi_{k+1}}(s_{t+1})] = Q_h^{\pi_{k+1}}(s_t, a_t). \end{aligned}$$

□

### 2.2.3.3 Soft Policy Iteration

**Lemma 2.9.** Let the reward be bounded,  $|\mathcal{A}| < \infty$ ,  $|\mathcal{S}| < \infty$ ,  $0 \leq \gamma < 1$  and  $\alpha > 0$ . Let  $\Pi$  be a set of all stationary, stochastic policies. Alternating between exact soft

## 2.2 Maximum-Entropy Reinforcement Learning

policy evaluation and global soft policy improvement from some initial policy  $\pi_0 \in \Pi$  and soft state-action function  $Q_h^{(0)}$ , the process converges (in the limit) to the optimal policy  $\pi^*$ , satisfying  $Q_h^{\pi^*} = Q_h^*$ .

*Proof.* At each iteration  $k$ , the policy  $\pi_k$  is fixed during evaluation. By Lemma 2.6,  $\mathcal{T}_h^{\pi_k}$  is a  $\gamma$ -contraction and the evaluation sequence converges to the fixed point  $Q_h^{\pi_k}$ . Therefore,  $Q_h^{\pi_k}$  is well defined at each iteration. By Lemma 2.8, the sequence  $\{Q_h^{\pi_k}(s_t, a_t)\}$  is non-decreasing for each pair  $(s_t, a_t)$ . Because the reward is bounded and  $|\mathcal{S}| < \infty$ ,  $\mathcal{A} < \infty$ , the sequence converges to some point-wise limit  $\bar{Q}$ . Policies are updated by the softmax of the current  $Q_h$  and softmax is continuous, therefore the limit point  $\pi^*$  exists and is the softmax of  $\bar{Q}$ , i.e.  $\pi^* = \text{softmax}(\bar{Q})$ . Since  $\pi^* \in \Pi$ , Lemma 2.6 equally guarantees well-definedness in the limit point  $Q_h^{\pi^*}$ . The map  $\pi \mapsto Q_h^\pi$  is continuous, hence  $\pi_k \rightarrow \pi^*$  implies  $Q_h^{\pi_k} \rightarrow Q_h^{\pi^*}$  and therefore  $\bar{Q} = Q_h^{\pi^*}$ . It follows  $\pi^* = \text{softmax}(Q_h^{\pi^*})$  (the policy is soft-greedy w.r.t.  $Q_h^{\pi^*}$ ) and therefore  $Q_h^{\pi^*}$  satisfies the Bellman optimality equation, thus  $Q_h^{\pi^*} = Q_h^*$ .  $\square$

### 2.2.4 Theoretical Objectives and Practical Gradients

As with other modern reinforcement learning approaches, the actor and critic are approximated by neural networks with parameters  $\phi$  and  $\theta$ , respectively. The update is based on

$$J_{Q_h}(\theta) := \mathbb{E}_{\substack{s_t \sim d_\mu^\gamma, \\ a_t \sim \mu(\cdot | s_t), \\ s_{t+1} \sim P(\cdot | s_t, a_t)}} \left[ \frac{1}{2} (\mathcal{T}_h^\pi Q(s_t, a_t) - Q_\theta(s_t, a_t))^2 \right], \quad (2.48)$$

where  $d_\mu^\gamma(s) = \frac{1-\gamma}{1-\gamma^T} \sum_{t=0}^{T-1} \gamma^t \Pr(s_t = s | \mu)$ .

In practice, the target is approximated from the average of a batch sampled from the replay buffer

$$\hat{\mathcal{T}}_h^\pi Q(s_t, a_t) := r(s_t, a_t) + \gamma \hat{V}_h^\pi(s_{t+1}),$$

with  $\hat{V}_h^\pi(s_{t+1}) = Q(s_{t+1}, a_{t+1}) - \alpha \log \pi(a_{t+1} | s_{t+1})$  for a particular tuple and  $a_{t+1}$ , yielding the approximate critic gradient

$$\hat{\nabla}_\theta J_{Q_h} = \mathbb{E}_{(s_t, a_t, r_t, s_{t+1}) \sim \mathcal{B}} [(\hat{\mathcal{T}}_h^\pi Q(s_t, a_t) - Q_\theta(s_t, a_t)) \nabla_\theta Q_\theta(s_t, a_t)]. \quad (2.49)$$

The policy objective as a function of the actor parameters is

$$J_\pi(\phi) = \mathbb{E}_{\substack{s_t \sim d_\mu^\gamma, \\ \epsilon_t \sim \mathcal{N}(0,1)}} [\alpha \log \pi_\phi(a_t | s_t) - Q_\theta(s_t, a_t)], \quad a_t = f_\phi(s_t; \epsilon_t), \quad (2.50)$$

## 2. INTRODUCTION TO REINFORCEMENT LEARNING ALGORITHMS

---

where  $a_t = f_\phi(s_t; \epsilon_t)$  is the action obtained by a reparameterized policy  $f_\phi$ , an approach for calculating low variance gradients similar to [72]. The approximate policy gradient can be computed from uniform buffer samples,

$$\hat{\nabla}_\phi J_\pi(\phi) = \mathbb{E}_{\substack{s_t \sim \mathcal{B}, \\ \epsilon_t \sim \mathcal{N}(0,1)}} [\alpha \partial_\phi \log \pi_\phi(a_t | s_t)]_{a_t} + (\alpha \nabla_{a_t} \log \pi_\phi(a_t | s_t) - \nabla_{a_t} Q_\theta(s_t, a_t)) \nabla_\phi f_\phi(s_t; \epsilon_t). \quad (2.51)$$

The first term in the expectation arises because the policy depends directly on  $\phi$  and  $a_t$  depends indirectly on  $\phi$ .

### 2.2.5 Algorithm

Algorithm 2.5 presents the SAC algorithm. The original authors proposed a variant where the soft state value  $V_h^\pi(s_t)$  is approximated in addition to the state-action value. In the algorithm described below, the state value approximator is omitted.

---

#### Algorithm 2.5 SAC Algorithm

---

- 1: Initialize  $\theta, \phi$
  - 2: Initialize learning parameters  $\beta_Q, \beta_\pi$
  - 3:  $\mathcal{B} \leftarrow \emptyset$
  - 4: **repeat**
  - 5:   Receive initial observation state  $s_0$
  - 6:   **for each** environment step **do**
  - 7:      $a_t \sim \pi_\phi(\cdot | s_t)$
  - 8:     Observe transition  $s_{t+1} \sim P(\cdot | s_t, a_t)$
  - 9:     Observe reward  $r_t$
  - 10:     $\mathcal{B} \leftarrow \mathcal{B} \cup \{(s_t, a_t, r_t, s_{t+1})\}$
  - 11:   **for each** gradient step **do**
  - 12:     Sample a random mini-batch of  $N$  transitions  $(s_i, a_i, r_i, s_{i+1})$  from  $\mathcal{B}$
  - 13:      $\theta \leftarrow \theta + \beta_Q \hat{\nabla}_\theta J_{Q_h}$  (Eq. (2.49))
  - 14:      $\phi \leftarrow \phi + \beta_\pi \hat{\nabla}_\phi J_\pi$  (Eq. (2.51))
  - 15: **until** Stopping criterion
  - 16: **return**  $Q_\theta, \pi_\phi$
-

## 2.3 Distributional Reinforcement Learning

Distributional reinforcement learning (DistRL) shifts the focus from modeling expected returns, as in traditional value-based methods, to capturing the entire distribution of returns. Therefore, the state-action value function is treated as random variable, e.g. in the actor-critic setting in (2.18), the distribution over returns is

$$Z^\pi(s_k, a_k) := \sum_{t=k}^{T-1} \gamma^{t-k} r(s_t, a_t), (a_t \sim \pi(\cdot | s_t) \forall t = k + 1 \dots T - 1, \quad (2.52)$$

$$s_{t+1} \sim P(\cdot | s_t, a_t) \forall t = k \dots T - 1) | s_k, a_k.$$

If the reward is treated as a deterministic quantity, then stochasticity in the return distribution arises from the transition kernel and the stochastic policy. Although DistRL was devised originally for intrinsic uncertainty [20], i.e. the reward is explicitly treated as a random variable, experiments show that distributional variants outperform their scalar counterparts even in environments with deterministic rewards, see Chapter 3. The expectation recovers the state-action value,

$$Q^\pi(s_k, a_k) = \mathbb{E}[Z^\pi(s_k, a_k)]. \quad (2.53)$$

Empirical data indicate that DistRL algorithms often outperform their non-distributional counterparts. The reasons for their accelerated learning remain a subject of research and several hypotheses and theories have been presented [20, 73, 74].

### 2.3.1 Modelling Distributions and Statistical Metrics

In practice, assumptions about the distribution underlying the random variable  $Z$  matters; modelling choices have important consequences for the optimization dynamics and for convergence. Thus, the choice of statistical metric to quantify the distance between the estimated value distribution and bootstrapped target must be made carefully. An inappropriate choice can hinder convergence to a (local) minimum. Bellemare et al. (2017) [20] formalized and analysed the distributional perspective in reinforcement learning for the first time in context of the  $p$ -Wasserstein metric. For two real-valued random variables  $U$  and  $V$ , the  $p$ -Wasserstein distance is defined via their quantile functions as

$$d_W(U, V) := \left( \int_0^1 |F_U^{-1}(x) - F_V^{-1}(x)|^p dx \right)^{\frac{1}{p}}, \quad (2.54)$$

where  $F$  is the general CDF and  $F^{-1}(\tau) = \inf\{z : F(z) \geq \tau\}$  denotes the quantile function. Note that the distance is formally between two laws of the random variables,

## 2. INTRODUCTION TO REINFORCEMENT LEARNING ALGORITHMS

---

$d_W(\mathcal{L}(U), \mathcal{L}(V))$ , but the common shorthand  $d_W(U, V)$  is used to avoid convoluted notation. The C51 algorithm presented in [20] incorporates a projection step in order to utilise the Wasserstein metric. As highlighted in [75], the empirical Wasserstein loss does not enjoy an “unbiased sample gradients” property, which can make naive stochastic gradient minimisation of the sample Wasserstein problematic. Quantile regression networks and implicit quantile regression networks were introduced in [76, 77] as a solution for the bias problem. They yield unbiased sample gradients calculated from a loss that corresponds to minimizing the 1-Wasserstein. The following sections provide the theoretical framework for DistRL. In the sequel, DistRL is analyzed using the Wasserstein metric.

### 2.3.2 Evaluation Setting

The distributional counterpart of the Bellman operator in (2.22) is used for distributional evaluation. It is defined as

$$\mathcal{T}_D^\pi Z(s_t, a_t) \stackrel{d}{=} r(s_t, a_t) + \gamma Z(s_{t+1}, a_{t+1}) \mid s_{t+1} \sim P(\cdot \mid s_t, a_t), a_{t+1} \sim \pi(\cdot \mid s_{t+1}), \quad (2.55)$$

where  $\stackrel{d}{=}$  indicates that the random variable on both sides are distributed according to the same law. To prove that  $\mathcal{T}_D^\pi$  is a contraction mapping, the following property of the  $p$ -Wasserstein metric is required,

$$d_W(aU, aV) = |a|d_W(U, V), \quad (\text{W1})$$

where  $a$  is an arbitrary constant. The following lemma invokes the linearity property of the Wasserstein metric. If two distributions  $\mu$  and  $\nu$  are composite and have identical weights  $w_i \geq 0$ ,  $\sum_i w_i = 1$ , i.e.  $\mu = \sum_i^k w_i \mu_i$  and  $\nu = \sum_i^k w_i \nu_i$ , then the  $p$ -Wasserstein distance between the two mixtures is upper-bounded by the Wasserstein distance between their weighted components

$$d_W(\mu, \nu) \leq \left( \sum_{i=1}^k w_i d_W(\mu_i, \nu_i)^p \right)^{1/p}. \quad (\text{W2})$$

**Lemma 2.10.** Fix  $1 \leq p < \infty$  and  $(s_t, a_t) \in \mathcal{S} \times \mathcal{A}$ . Let  $\bar{d}_W(Z_1, Z_2) := \sup_{s,a} d_W(Z_1(s, a), Z_2(s, a))$  for two state-action value distributions  $Z_1$  and  $Z_2$ . Assume that for each  $(s, a)$ , the distributions have finite  $p$ -moments to attain completeness of the metric space. Further, assume  $|\mathcal{A}| < \infty$ ,  $|\mathcal{S}| < \infty$  and  $0 \leq \gamma < 1$ . Then  $\mathcal{T}_D^\pi$  is a  $\gamma$ -contraction in the Wasserstein metric, i.e.

$$\bar{d}_W(\mathcal{T}_D^\pi Z_1, \mathcal{T}_D^\pi Z_2) \leq \gamma \bar{d}_W(Z_1, Z_2). \quad (2.56)$$

*Proof.*

$$\begin{aligned}
 & d_W(\mathcal{T}_D^\pi Z_1, \mathcal{T}_D^\pi Z_2) \\
 &= d_W(r(s_t, a_t) + \gamma Z_1(s_{t+1}, a_{t+1}), r(s_t, a_t) + \gamma Z_2(s_{t+1}, a_{t+1})) \\
 &\quad \text{where } s_{t+1} \sim P(\cdot | s_t, a_t), a_{t+1} \sim \pi(\cdot | s_{t+1}). \\
 &\text{By translational invariance and scaling property (W1)} \\
 &= \gamma d_W(Z_1(s_{t+1}, a_{t+1}), Z_2(s_{t+1}, a_{t+1})), \quad s_{t+1} \sim P(\cdot | s_t, a_t), a_{t+1} \sim \pi(\cdot | s_{t+1}). \\
 &\text{By property (W2) with } w_{t+1} := P(s_{t+1} | s_t, a_t)\pi(a_{t+1} | s_{t+1}) \\
 &\leq \gamma \left( \sum_{s_{t+1}, a_{t+1}} w_{t+1} d_W(Z_1(s_{t+1}, a_{t+1}), Z_2(s_{t+1}, a_{t+1}))^p \right)^{1/p} \\
 &\leq \gamma \left( \sum_{s_{t+1}, a_{t+1}} w_{t+1} \left( \sup_{(s,a)} d_W(Z_1(s, a), Z_2(s, a))^p \right) \right)^{1/p} \\
 &= \gamma \left( \sum_{s_{t+1}, a_{t+1}} w_{t+1} \right)^{1/p} \sup_{(s,a)} d_W(Z_1(s, a), Z_2(s, a)). \\
 &\text{Since } \sum_i w_i = 1 \\
 &= \gamma \bar{d}_W(Z_1, Z_2).
 \end{aligned}$$

Taking the supremum over  $(s, a) \in \mathcal{S} \times \mathcal{A}$  on both sides yields the contraction

$$\bar{d}_W(\mathcal{T}_D^\pi Z_1, \mathcal{T}_D^\pi Z_2) \leq \gamma \bar{d}_W(Z_1, Z_2).$$

Since  $\mathcal{T}_D^\pi$  is a  $\gamma$ -contraction, it follows from the Banach theorem that the series  $\{Z_k\}$  with  $Z_{k+1} := \mathcal{T}_D^\pi Z_k$  converges to a unique fixed point  $Z^\pi$ , i.e.  $Z = \mathcal{T}_D^\pi Z$  and  $\bar{d}_W(Z_k, Z^\pi) \rightarrow 0$  as  $k \rightarrow \infty$ .  $\square$

Contraction mapping can also be shown in the expectation.

**Lemma 2.11.** Assume  $\sup_{s \in \mathcal{S}, a \in \mathcal{A}} \mathbb{E}[|Z(s, a)|] < \infty$ ,  $|\mathcal{A}| < \infty$  and  $|\mathcal{S}| < \infty$  then

$$\|\mathbb{E} \mathcal{T}_D^\pi Z_1 - \mathbb{E} \mathcal{T}_D^\pi Z_2\|_\infty \leq \gamma \|\mathbb{E} Z_1 - \mathbb{E} Z_2\|_\infty. \quad (2.57)$$

*Proof.* If the first moment and spaces are finite, linearity of expectation ensues,

$$\|\mathbb{E} \mathcal{T}_D^\pi Z_1 - \mathbb{E} \mathcal{T}_D^\pi Z_2\|_\infty = \|\mathcal{T}^\pi \mathbb{E} Z_1 - \mathcal{T}^\pi \mathbb{E} Z_2\|_\infty.$$

From Lemma (2.3) and Eq. (2.53),

$$\|\mathcal{T}^\pi \mathbb{E} Z_1 - \mathcal{T}^\pi \mathbb{E} Z_2\|_\infty \leq \gamma \|\mathbb{E} Z_1 - \mathbb{E} Z_2\|_\infty,$$

## 2. INTRODUCTION TO REINFORCEMENT LEARNING ALGORITHMS

---

thus

$$\|\mathbb{E} \mathcal{T}_D^\pi Z_1 - \mathbb{E} \mathcal{T}_D^\pi Z_2\|_\infty \leq \gamma \|\mathbb{E} Z_1 - \mathbb{E} Z_2\|_\infty.$$

□

### 2.3.3 Control Setting

Recall that in the control setting, the goal is to find a policy that maximizes the value function. This principle is central to reinforcement learning and has been formalized in the classical context: for the critic-only framework, this was detailed in Equation (2.3), while for the actor-critic (AC) framework, it was introduced in Equation (2.24). In critic-only settings, a distributional variant of the Bellman optimality operator can be defined as

$$\begin{aligned} \mathcal{T}_D Z(s_t, a_t) : & \stackrel{d}{=} r(s_t, a_t) + \gamma Z(s_{t+1}, a_{t+1}) \mid s_{t+1} \sim P(\cdot \mid s_t, a_t), \\ & a_{t+1} \in \operatorname{argmax}_a \mathbb{E}[Z(s_{t+1}, a)]. \end{aligned}$$

However, as shown in [20], this operator lacks a contraction property in the general sense. It may admit multiple fixed points and even if a fixed point  $Z^* = \mathcal{T}_D Z^*$  exists, convergence to that fixed point from iterative application is not guaranteed under any metric [20]. These challenges cast the distributional Bellman optimality operator in a less favorable light for practical applications. However, if there is a unique optimal policy  $\pi^*$  or if there is a total ordering on the set of optimal policies, the optimal distributional Bellman operator is still guaranteed to have a fixed point.

### 2.3.4 Moment-Agnostic Distributions

In this work, free-moment distributions refer to probability distributions that are not constrained in their moments, meaning that they can adopt arbitrary moments. The following sections introduce free-moment algorithms in the critic-only framework (Section 2.1.2.1).

#### 2.3.4.1 Fixed Supports - C51

The C51 algorithm, introduced in [20], models the return distribution using a discrete representation with  $N$  atoms. These fixed supports are defined by a discretization of the interval  $[V_{min}, V_{max}]$  as  $\{z_i = V_{min} + i\Delta z : 0 \leq i < N\}$ , where  $V_{min}, V_{max} \in \mathbb{R}$  are assumed value limits and  $\Delta z := \frac{V_{max} - V_{min}}{N-1}$  represents the spacing between the atoms.

## 2.3 Distributional Reinforcement Learning

The return distribution is approximated by

$$Z_\theta = z_i, \quad \text{with probability } p_i(s_t, a_t) := \frac{e^{\ell_{\theta_i}(s_t, a_t)}}{\sum_j e^{\ell_{\theta_j}(s_t, a_t)}},$$

where  $\ell_{\theta_i}(s_t, a_t)$  is the  $i$ -th logit of the approximator. Note that  $i$  is the atom of the prediction and  $j$  is the atom of the target. Although theoretically attractive, challenges arise in directly minimizing the Wasserstein distance between the estimated distribution  $Z$  and the target  $\mathcal{T}_D$  in a mini-batch. Lemma 7 in [20] shows that stochastic gradient descent generally cannot achieve this metric. In addition, the support of  $\mathcal{T}_D Z$  generally do not coincide with that of  $Z$ . To address this, C51 projects the distribution of  $\mathcal{T}_D Z$  onto  $Z$  through a projection  $\Phi$ . This effectively reduces the problem to a multi-class classification. In its closed-form, the  $i$ -th component of the projected target is defined by

$$m_i = (\Phi \mathcal{T}_D Z(s_t, a_t))_i := \sum_{j=0}^{N-1} \left[ 1 - \frac{|[\mathcal{T}_D z_j]_{V_{min}^{max}} - z_i|}{\Delta z} \right] p_j(s_{t+1}, a_{t+1}) \quad (2.58)$$

and the loss becomes

$$J_C(\theta) = \mathbb{E}_{\substack{s_t \sim d_\mu^\gamma, \\ a_t \sim \mu(\cdot | s_t)}} \left[ - \sum_i m_i \log p_i(s_t, a_t) \right]. \quad (2.59)$$

The expression  $[\mathcal{T} z_j]_{V_{min}^{max}}$  denotes the actual target bounded by the limits. The projected categorical loss is related to minimizing the 1-Cramér distance [78].

Another limitation of C51 lies in the reliance on domain knowledge to set the parameters  $V_{min}$  and  $V_{max}$ , which define the limits of the supports. This practical challenge highlights a gap between theory and implementation, later addressed by quantile regression networks that can operate under the Wasserstein metric and eliminate the need for such predefined bounds.

During training, mini-batches of experience tuples  $(s_t, a_t, r_t, s_{t+1})$ , with  $r_t = r(s_t, a_t)$ , are uniformly sampled from the replay buffer  $\mathcal{B}$  (thus disregarding state visitation discounting). Analogous to the scalar variant of the sample Bellman optimality operator, the distributional Bellman optimality operator for a particular  $s_{t+1}$  (assuming a unique optimal policy) is

$$\hat{\mathcal{T}}_D Z(s_t, a_t) \stackrel{d}{=} r(s_t, a_t) + \gamma Z(s_{t+1}, a_{t+1}), \quad a_{t+1} = \operatorname{argmax}_a \mathbb{E}[Z(s_{t+1}, a)].$$

The  $i$ -th projected sample target in its closed-form becomes

$$\hat{m}_i = (\Phi \hat{\mathcal{T}}_D Z(s_t, a_t))_i := \sum_{j=0}^{N-1} \left[ 1 - \frac{|[\hat{\mathcal{T}}_D z_j]_{V_{min}^{max}} - z_i|}{\Delta z} \right]^1 p_j(s_{t+1}, a_{t+1}),$$

## 2. INTRODUCTION TO REINFORCEMENT LEARNING ALGORITHMS

---

where  $\hat{\mathcal{T}}_D z_j \big|_{V_{min}^{max}}$  is the moved  $j$ -th atom of the target.

Algorithm 2.6 describes the C51 procedure. For improved readability, the gradient descent step is exemplified using one experience tuple and it is assumed that the target index is always a non-integer. Consequently, for a particular tuple, the sample loss is described as

$$\hat{J}_C(\theta) = - \sum_i \hat{m}_i \log p_i(s_t, a_t). \quad (2.60)$$

As highlighted in the description of the previous off-policy algorithms, this practical choice disregards state-visitation discounting.

### 2.3.4.2 Quantile Regression

Quantile Regression in DistRL offers a novel approach to represent the distribution over returns using a set of quantiles. This approach enables the direct computation of sample gradients over the Wasserstein metric, bridging the gap between theory and practical application described in [20]. The algorithm derived from this framework, Quantile Regression Deep Q-Learning (QR-DQN), employs a new parameterization of the approximator as well as a quantile regression loss [76]. In QR-DQN, the parameterization of C51 is transposed in so far that the model approximates quantile distributions that map state-action pairs to uniform probability distributions supported on the set  $\{q_i(s_t, a_t)\}_{i=1}^N$ . The quantiles are fixed and the values vary. The model is defined as

$$Z_\theta(s_t, a_t) := \frac{1}{N} \sum_{i=1}^N \delta_{q_\theta^i(s_t, a_t)}, \quad (2.61)$$

where  $\delta_{q_\theta^i(s_t, a_t)}$  denotes a delta distribution at  $q_\theta^i(s_t, a_t)$  and  $N$  is the number of uniform masses. For policy evaluation, the Wasserstein-1 metric between an arbitrary target and the parameterized model is minimized,

$$\Pi_{W_1} Z := \operatorname{argmin}_{Z_\theta} d_{W_1}(Z, Z_\theta),$$

where

$$d_{W_1}(U, V) := \sum_{i=1}^N \int_{\tau_{i-1}}^{\tau_i} |F_U^{-1}(\omega) - q^i| d\omega.$$

Lemma 2 of [76] shows that the set

$$\left\{ q^i \in \mathbb{R} \mid F(q^i) = \left( \frac{\tau_{i-1} + \tau_i}{2} \right) \right\} \forall i \in I$$

### 2.3 Distributional Reinforcement Learning

minimizes  $\int_{\tau_{i-1}}^{\tau_i} |F_U^{-1}(\omega) - q_i| d\omega$ . Therefore it is sufficient for the approximator to output  $q_\theta^i$  that corresponds to  $\hat{\tau}_i = (\tau_{i-1} + \tau_i)/2$ . The quantile regression loss is defined as

$$J_{QR}(\theta) := \sum_i^N \mathbb{E}_{\substack{s_t \sim d_{\mu}^{\gamma}, \\ a_t \sim \mu(\cdot|s_t)}} [\rho_{\hat{\tau}_i}(\mathcal{T}_D Z(s_t, a_t) - q_\theta^i(s_t, a_t))] \quad \text{where} \quad (2.62)$$

$$\rho_{\hat{\tau}_i}(u) = u(\hat{\tau}_i - \delta_{\{u < 0\}}), \forall u \in \mathbb{R}$$

To keep a simple notation for the subsequent analysis, the objective is reformulated as a direct function of the quantile value  $q$ . Let  $\hat{Z} = \mathcal{L}(\mathcal{T}_D Z(s_t, a_t))$ , then

$$J_{QR}(q) := \sum_i^N \mathbb{E}_{Z \sim \hat{Z}} [\rho_{\hat{\tau}_i}(Z - q^i)] \quad \text{where}$$

$$\rho_{\hat{\tau}_i}(u) = u(\hat{\tau}_i - \delta_{\{u < 0\}}), \forall u \in \mathbb{R}.$$

Note that the multiplier introduces an asymmetric loss. Overestimation in lower quantiles is more severely punished (and in the right proportions as shown in the following paragraphs). The quantile regression objective moves  $q^i$  to the true corresponding quantile in expectation. To show this, consider an objective with one quantile to be approximated from samples

$$\frac{\partial \rho_\tau}{\partial u} = \tau - \mathbb{I}_{\{u < 0\}}.$$

Since  $\frac{\partial u}{\partial q} = -1$ , by chain rule

$$\frac{\partial \rho_\tau}{\partial q} = \mathbb{I}_{\{Z - q < 0\}} - 1.$$

The expectation over the gradients is

$$\frac{\partial J_{QR}}{\partial q} = \mathbb{E}_{Z \sim \hat{Z}} [\mathbb{I}_{\{Z - q < 0\}} - \tau]. \quad (2.63)$$

At the solution, the gradient must be zero (albeit not the loss):

$$0 = \mathbb{E}_{Z \sim \hat{Z}} [\mathbb{I}_{\{Z - q < 0\}} - \tau]. \quad (2.64)$$

Because  $\mathbb{E}[\mathbb{I}_A] = \Pr(A)$ , Equation (2.64) yields

$$\Pr(Z < q) = \tau.$$

That is, in expectation the approximator outputs a  $q$  such that the CDF at  $q$  corresponds to  $\tau$ . Further, Equation (2.63) shows that the gradient  $g(Z; q) = \mathbb{I}_{\{Z - q < 0\}} - \tau$  is unbiased.

## 2. INTRODUCTION TO REINFORCEMENT LEARNING ALGORITHMS

---

To smoothen the gradients as  $u \rightarrow 0$ , Equation (2.62) is augmented with the Huber loss

$$J_{QHR}(\theta) := \sum_i^N \mathbb{E}_{\substack{s_t \sim d_\mu^\gamma \\ a_t \sim \mu(\cdot|s_t)}} [\rho_{\hat{\tau}_i}^\kappa(\mathcal{T}_D Z(s_t, a_t) - q_\theta^i(s_t, a_t))], \quad \text{where} \quad (2.65)$$

$$\rho_{\hat{\tau}}^\kappa(u) := |\hat{\tau}_i - \delta_{\{u < 0\}}| H(u; \kappa)$$

with

$$H^\kappa(u) := \begin{cases} \frac{1}{2}u^2 & \text{if } |u| < \kappa \\ \kappa(|u| - \frac{1}{2}\kappa) & \text{otherwise.} \end{cases}$$

In the the policy evaluation setting, the authors in [76] showed contraction of the projected Bellman operator only in the  $\infty$ -norm, i.e.

$$\bar{d}_W(\Pi_{W_1} \mathcal{T}_D^\pi Z_1, \Pi_{W_1} \mathcal{T}_D^\pi Z_2) \leq \bar{d}_W(\Pi_{W_1} Z_1, \Pi_{W_1} Z_2),$$

where the common shorthand  $d_W(U, V)$  is used for distance between two laws of random variables. Recall from Section 2.3.3 that the optimal distributional Bellman operator is guaranteed to have a fixed point iff there exists a unique optimal policy or if a total ordering on the set of optimal policies is imposed. Analogous to the optimal distributional Bellman operator, if the greedy action is unique, then the projected distributional Bellman optimality operator admits a fixed point. QR-DQN is derived from the theoretical framework outlined in this Section and described in Algorithm 2.7. To improve readability, the algorithm describes a gradient update for one single transition. Note that if the reward is treated as a deterministic quantity (s. Section 2.3), the target distribution for a particular  $s_{t+1}$  can be written quantile-wise in the following notation:

$$\mathcal{T}_D q_j(s_t, a_t) = r(s_t, a_t) + \gamma q_j(s_{t+1}, a_{t+1}), \quad (2.66)$$

where

$$a_{t+1} = \operatorname{argmax}_a \mathbb{E}[Z(s_{t+1}, a)] \iff a_{t+1} = \operatorname{argmax}_a \frac{1}{N} \sum_i^N q_j(s_{t+1}, a).$$

### 2.3.4.3 Implicit Quantile Networks

Implicit Quantile Networks (IQN) extend the Distributional Reinforcement Learning paradigm by improving upon the limitations of QR-DQN, primarily through a more flexible representation of the value distribution. In contrast to QR-DQN, which relies on

### 2.3 Distributional Reinforcement Learning

---

a fixed set of quantile fractions, IQN models the value distribution more expressively by learning a continuous quantile function [77]. The goal is to model the random variable  $Z_\tau := F_Z^{-1}(\tau)$ , where  $\tau$  is uniformly sampled, i.e.  $\tau \sim \mathcal{U}([0, 1])$ . This approach allows the return distribution to be modeled by inference from an approximator at different quantiles. This modeling choice captures the variability in returns across quantiles without explicitly parameterizing the distribution. In the control setting, IQN incorporates risk-sensitive learning by considering how a utility function  $U$  influences decision-making in the presence of risk. A concave  $U$  reflects risk-averse behavior, a convex  $U$  corresponds to risk-seeking behaviour, and the identity function reflects risk neutrality. By sampling  $\tau$  non-uniformly, the policy can emphasize different parts of the return distribution. For example, sampling  $\tau$  near 1 focuses on upper quantiles, reflecting optimism and risk-seeking behaviour. This non-uniform sampling can be thought of as  $\tau$  being re-weighted by some distortion risk measure  $\xi : [0, 1] \mapsto [0, 1]$ , which transforms the baseline uniform sampling. The distorted expectation of the value distribution is given by

$$Q_\xi(s_t, a_t) := \mathbb{E}_{\tau \sim \mathcal{U}([0,1])} [Z_{\xi(\tau)}(s_t, a_t)]. \quad (2.67)$$

The corresponding policy under  $\xi$  is derived through

$$\pi_\xi(s_t) = \operatorname{argmax}_{a \in \mathcal{A}} Q_\xi(s_t, a). \quad (2.68)$$

Policy evaluation is conducted with uniformly sampled quantiles, with  $N$  and  $N'$  samples for the current estimate and Bellman update, respectively. The same quantile-huber regression loss as in QR-DQN, Eq. (2.65), is employed. In practice, the target distribution represented by a parameterized model for a particular  $s_{t+1}$  is computed from  $\tau' \sim U([0, 1])$  by

$$\mathcal{T}_D Z_{\tau'}(s_t, a_t) := r(s_t, a_t) + \gamma Z_{\tau'}(s_{t+1}, a_{t+1}), \quad (2.69)$$

where

$$a_{t+1} = \operatorname{argmax}_a \frac{1}{K} \sum_{j=1}^K Z_{\xi(\tau'_j)}(s_{t+1}, a).$$

IQN employs a novel network architecture for approximating the quantile function. The sampled  $\tau$  is transformed into a high-dimensional embedding vector via a cosine encoding  $g(\tau) := \cos(\pi k \tau)$ ,  $k = 1, \dots, K$ . This ensures that the network is sensitive to small changes in  $\tau$  due to the larger shifts in higher frequencies. The cosine basis functions enable the construction of complex, non-linear approximations through learned

## 2. INTRODUCTION TO REINFORCEMENT LEARNING ALGORITHMS

---

linear combinations,  $\phi(\tau) = \text{ReLU}(W_\phi g(\tau))$ . Another motivation for the embedding is to facilitate Feature-Wise Linear Modulation (FiLM). Specifically, the Hadamard product of the quantile embedding and latent state representation generated by a network  $\psi$  is computed and passed through the last fully connected layers to produce the quantile value estimate, i.e.  $Z_\tau^\psi(s_t, a_t) = f(\psi(s_t) \odot \phi(\tau))_{a_t}$ . As IQN represents a functional generalization of QR-DQN, Algorithm 2.8 isolates the primary divergence, namely the stochastic sampling of quantiles and the resulting integration of the quantile embedding function.

### 2.3.5 Fixed-Moment Distributions

Fixed-moment distributions represent a specific class of approximations within the broader field of distributional reinforcement learning. These distributions are fully characterized by their first two moments. Each return distribution is obtained by shifting and rescaling a single standardized density (e.g. the standard Gauss). The model changes in location and scale, but holds higher-order shape characteristics constant across states and actions. By constraining the return distribution to this simplified representation, fixed-moment distributional reinforcement learning algorithms operate under strong assumptions about the underlying true distributions of returns. The fixed-moment approach offers computational simplicity, but it comes at the potential cost of accuracy. The rigid assumptions can (and generally will) lead to biased approximations of the true return distribution. However, this error introduces an intriguing tradeoff between variance and bias in learning. By reducing the complexity of the distribution - requiring fewer parameters for approximation - fixed-moment algorithms can yield lower gradient variance during optimization. The reduction in variance can enhance learning stability and efficiency, making these algorithms appealing in scenarios where computational resources are limited, as was the case for executing experiments for this work. High variance might also impede effective policy updates.

#### 2.3.5.1 Parameterized Distributions - Gauss

Representing return distributions via Gaussian parameterization offers a computationally efficient and straightforward baseline. The Gaussian probability density function is defined as

$$\varphi_{\mu,\sigma}(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2}, \quad (2.70)$$

where  $\mu$  represents the mean, and  $\sigma$  denotes the standard deviation. In the context of reinforcement learning, since  $Q = \mathbb{E}[Z]$ , the mean  $\mu$  of the Gaussian can be directly

## 2.3 Distributional Reinforcement Learning

---

replaced with  $Q$ . This makes the Gaussian distribution particularly straightforward to integrate into existing reinforcement learning frameworks. The Gaussian distribution has been utilized for modeling the policy, yielding strong empirical performance in algorithms such as PPO (Section 2.1.6.4) and SAC (Section 2.2.3). Modeling returns as Gaussian distributions is hypothesized to enhance algorithmic stability while maintaining high performance. Empirical evaluations, detailed in Sections 3.3 and 3.5, corroborate this assumption. The use of Gaussian distributions for fixed-moment representations is advantageous due to their well-understood mathematical properties and computational efficiency. While the assumption of a symmetric and unimodal distribution may not always align with the true return distribution, the Gaussian provides a practical tradeoff between simplicity and performance.

### 2.3.5.2 Parameterized Distributions - Gauss Mixture Models

To allow for a more expressive representation, the return in distributional reinforcement learning can be modeled as a mixture distribution. One promising approach is to use Gaussian Mixture Models (GMMs). GMMs have been used for population modelling by Pearson as early as 1894 [79]. Gaussian mixture models represent return distributions as weighted sums of Gaussian components and, with sufficiently many components, could approximate a wide class of probability densities to arbitrary accuracy. The GMM return distribution is expressed as:

$$q_{\mu,\sigma}(x) = \sum_{k=1}^K w_k \varphi_{\mu_k, \sigma_k}(x),$$

where  $w_k$  is the weight of the  $k$ -th kernel and  $\varphi_{\mu_k, \sigma_k}$  denotes the Gaussian kernel with mean  $\mu_k$  and variance  $\sigma_k^2$ . GMMs are particularly appealing for DistRL because they can inherently model multimodality, allowing the return distribution to capture diverse structural properties of the environment. This could facilitate deducing insights about the dynamics of the environment. To assess the adequacy of GMMs in deep DistRL, experiments have been conducted simulating a  $TD(0)$ -step (Section 2.1.2.2) with a noisy target. The regularized Cramér metric, Equation (3.1), was employed as a loss function, measuring the distance between the predicted and target distributions. Figure 2.2 demonstrates the results of fitting a simple two-kernel GMM neural network approximator with a single hidden layer containing 10 nodes. In this setup, the learnable parameters are the kernel weights, means, and variances, each represented by an output node. To reduce the impact of noise on the results, the approximators are trained to output baseline predictions first. Figures 2.2b and 2.2d show the resulting

## 2. INTRODUCTION TO REINFORCEMENT LEARNING ALGORITHMS

---

target fits from high- and low-variance baselines. The simple approximator was able to fit the target distributions almost perfectly. The left column of figure 2.2 illustrates the CDF approximations from low (blue) and high (orange) variance baselines to the targets (green) at the end of the training. Figure 2.2a shows fitting to the high-variance target while Figure 2.2c shows fitting to the low variance target. Since the target is approximated well, the curves coincide and only the target curve is visible. The right column shows the change of the kernel means as a function of training iterations. Figures 2.2b and 2.2d reveal that while the kernel means converge to the target means, the fitting process is slow. The experimental configurations are detailed in Table 2.1.

In contrast, as the complexity of the GMM neural network approximator increases, its performance in fitting the target deteriorates. Subsequent experiments were conducted with more complex model architectures, employing a two-hidden-layer MLP with 256 units per layer. Note that the target distribution means are different from the simple model experiments. Details are listed in Table 2.2. Figure 2.3 shows that a more sophisticated approximator fails to generalize well to the target distributions. Figure 2.3b shows that the model does not sufficiently learn the high-variance target means. Figures 2.3a and 2.3c indicate significant discrepancies in the target and predicted variances after training is concluded. These results suggest that overparameterization in GMM neural network approximators hinder their ability to capture the target distribution effectively. In reinforcement learning, however, overparameterized models achieve better generalization in diverse, procedurally generated settings, compared to simpler models [80, 81]. The substantial reduction in parameter dimensionality required for GMM compatibility does not compensate for the resulting degradation in performance that higher-capacity models generally offer. Therefore, GMMs will not be employed in this work. Future research may focus on identifying techniques that enable the use of large GMMs.

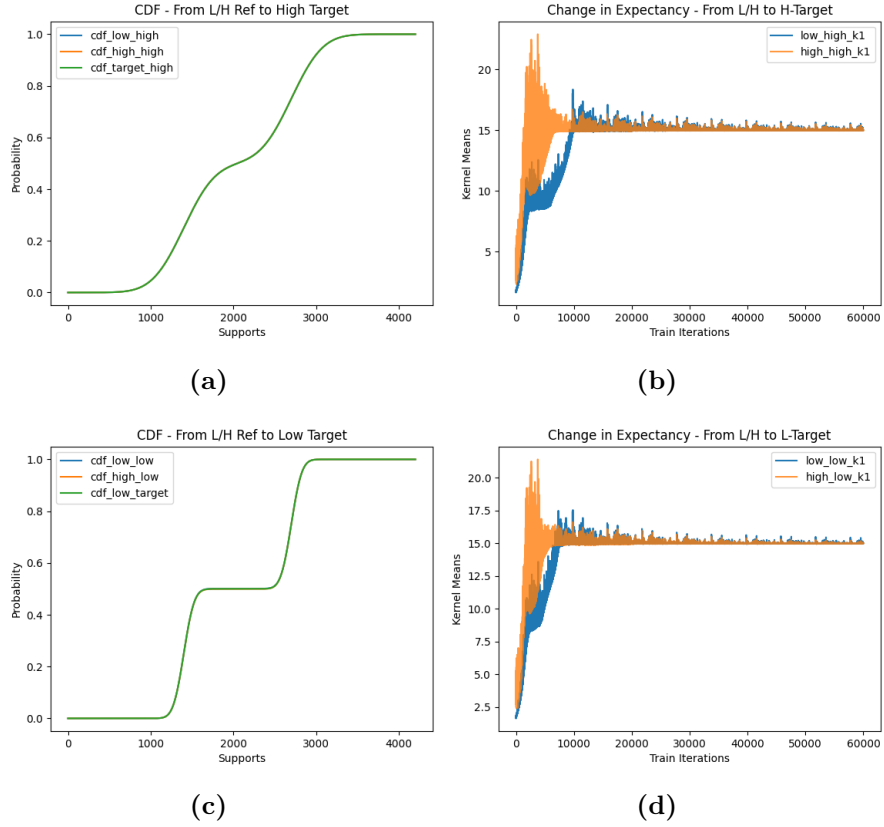
Parameter	Baseline		Target	
	Low	High	Low	High
Means $[\mu_1, \mu_2]$	[0.0, 1.0]	[0.0, 1.0]	[13.0, 17.0]	[13.0, 17.0]
Std. Dev. $[\sigma_1, \sigma_2]$	[1.0, 1.0]	[3.0, 3.0]	[1.0, 1.0]	[3.0, 3.0]
Weights $[w_1, w_2]$	[0.5, 0.5]	[0.5, 0.5]	[0.5, 0.5]	[0.5, 0.5]
MLP-Architecture	$(inp \rightarrow 10 \rightarrow out)$		N.A.	

**Table 2.1:** Parameters for baseline and target configurations with simple GMM approximators for experiments shown in Figure 2.2

## 2.3 Distributional Reinforcement Learning

Parameter	Baseline		Target	
	Low	High	Low	High
Means $[\mu_1, \mu_2]$	[0.0, 1.0]	[0.0, 1.0]	[2.0, 15.0]	[2.0, 15.0]
Std. Dev. $[\sigma_1, \sigma_2]$	[1.0, 1.0]	[3.0, 3.0]	[1.0, 1.0]	[3.0, 3.0]
Weights $[w_1, w_2]$	[0.5, 0.5]	[0.5, 0.5]	[0.5, 0.5]	[0.5, 0.5]
MLP-Architecture	$(inp \rightarrow 256 \rightarrow 256 \rightarrow out)$		N.A.	

**Table 2.2:** Parameters for baseline and target configurations with complex GMM approximators for experiments shown in Figure 2.3.



**Figure 2.2:** Plots (a) and (c) show fitting of a simple neural network GMM approximator with  $K = 2$ , consisting of one hidden layer with 10 nodes, to a high- and low-variance target from low (blue) and high (orange) reference distributions. All curves coincide, indicating sufficient fitting to the target. The neural network consists of one hidden layer containing 10 nodes. Figures (b) and (d) depict the change of the mean of one kernel per iteration in the neural network fitting process.

## 2. INTRODUCTION TO REINFORCEMENT LEARNING ALGORITHMS

---

**Algorithm 2.6** C51 Algorithm (with observation encoder  $\psi$ )

---

```

1: Initialize replay memory  $\mathcal{B}$  with capacity  $|\mathcal{B}|$ 
2: Initialize distributional state-action function  $Z$  with weights  $\theta$ 
3: Initialize observation encoder  $\psi$ 
4: Initialize learning parameter  $\beta_z$ 
5: Initialize exploration  $\epsilon = 1$  and decay parameter  $\beta_\epsilon$ 
6: repeat
7:   Receive initial (raw) observation state  $s_1$ 
8:   for each environment step do
9:     With probability  $\epsilon$  select random action  $a_t \sim \mathcal{U}(\{1, \dots, |\mathcal{A}|\})$ 
10:    Otherwise, select action  $a_t = \operatorname{argmax}_a \sum_i z_i p_i(s_t, a)$ 
11:    Observe transition  $s_{t+1} \sim P(\cdot | s_t, a_t)$ 
12:    Observe reward  $r_t$ 
13:     $\mathcal{B} \leftarrow \mathcal{B} \cup \{(s_t, a_t, r_t, s_{t+1})\}$ 
14:     $s_t \leftarrow s_{t+1}$ 
15:   for each gradient step do
16:     Sample random transition  $(s_t, a_t, r_t, s_{t+1})$  from  $\mathcal{B}$ 
17:      $a_{t+1} = \operatorname{argmax}_a \sum_i z_i p_i(\psi(s_{t+1}), a)$ 
18:      $m_i = 0, \quad i \in 0, \dots, N - 1$ 
19:     for  $j \in 0, \dots, N - 1$  do
20:       Set  $y_j \leftarrow \begin{cases} r_t |_{V_{min}^{max}} & \text{if done} \\ [\hat{\mathcal{T}}_D Z_{\theta_j}(\psi(s_t), a_t)]_{V_{min}^{max}} & \text{otherwise} \end{cases}$ 
21:        $b_j \leftarrow (y_j - V_{min}) / \Delta z$  # non-integer atom index of target
22:        $l \leftarrow \lfloor b_j \rfloor, u \leftarrow \lceil b_j \rceil$  # ( $l \neq u$  for simplicity)
23:       if done then
24:          $m_u \leftarrow m_u + (b_j - l)$ 
25:          $m_l \leftarrow m_l + (u - b_j)$ 
26:       break
27:     else
28:        $m_u \leftarrow m_u + p_j(\psi(s_{t+1}), a_{t+1})(b_j - l)$ 
29:        $m_l \leftarrow m_l + p_j(\psi(s_{t+1}), a_{t+1})(u - b_j)$ 
30:      $\theta \leftarrow \theta - \beta_Z \hat{\nabla}_\theta J_C$  from (2.60)
31:     Update encoder parameters via same gradient
32:      $\epsilon \leftarrow \max(\epsilon - \beta_\epsilon, 0)$ 
33:   until Stopping criterion
34: return  $Z_\theta$ 

```

---

---

**Algorithm 2.7** Quantile Regression Deep Q-Learning (with observation encoder  $\psi$ )

---

- 1: Initialize replay memory  $\mathcal{B}$  with capacity  $|\mathcal{B}|$
  - 2: Initialize quantile model  $Z$  with parameters  $\theta$  and  $N$  quantile weights
  - 3: Initialize  $\kappa$
  - 4: Initialize observation encoder  $\psi$
  - 5: Initialize learning parameter  $\beta_z$
  - 6: Initialize exploration  $\epsilon = 1$  and decay parameter  $\beta_\epsilon$
  - 7: **repeat**
  - 8:     Receive initial (raw) observation state  $s_t = s_1$
  - 9:     **for each** environment step **do**
  - 10:         With probability  $\epsilon$  select random action  $a_t \sim \mathcal{U}(\{1, \dots, |\mathcal{A}|\})$
  - 11:         Otherwise, select action  $a_t = \operatorname{argmax}_a \sum_i \frac{1}{N} q_\theta^i(\psi(s_t), a)$
  - 12:         Observe transition  $s_{t+1} \sim P(\cdot | s_t, a_t)$
  - 13:         Observe reward  $r_t$
  - 14:          $\mathcal{B} \leftarrow \mathcal{B} \cup \{(s_t, a_t, r_t, s_{t+1})\}$
  - 15:          $s_t \leftarrow s_{t+1}$
  - 16:     **for each** gradient step **do**
  - 17:         Sample random transition  $(s_t, a_t, r_t, s_{t+1})$  from  $\mathcal{B}$
  - 18:          $J_{QHR} \leftarrow 0$
  - 19:         **for**  $j = 1, \dots, N$  **do**
  - 20:             Set  $y_j = \begin{cases} r_t & \text{if done} \\ \mathcal{T}q_j(\psi(s_t), a_t) & \text{Eq.(2.66) otherwise} \end{cases}$
  - 21:              $J_{QHL} \leftarrow J_{QHL} + \frac{1}{N^2} \sum_{i=1}^N \rho_{\tau_i}^\kappa (y_j - q_\theta^i(\psi(s_t), a_t)),$    Eq (2.65) normalized
  - 22:              $\theta \leftarrow \theta - \beta_Z \hat{\nabla}_\theta J_{QHL}$
  - 23:             Update encoder parameters via same gradient
  - 24:              $\epsilon \leftarrow \max(\epsilon - \beta_\epsilon, 0)$
  - 25:     **until** Stopping criterion
  - 26: **return**  $Z_\theta$
-

## 2. INTRODUCTION TO REINFORCEMENT LEARNING ALGORITHMS

---



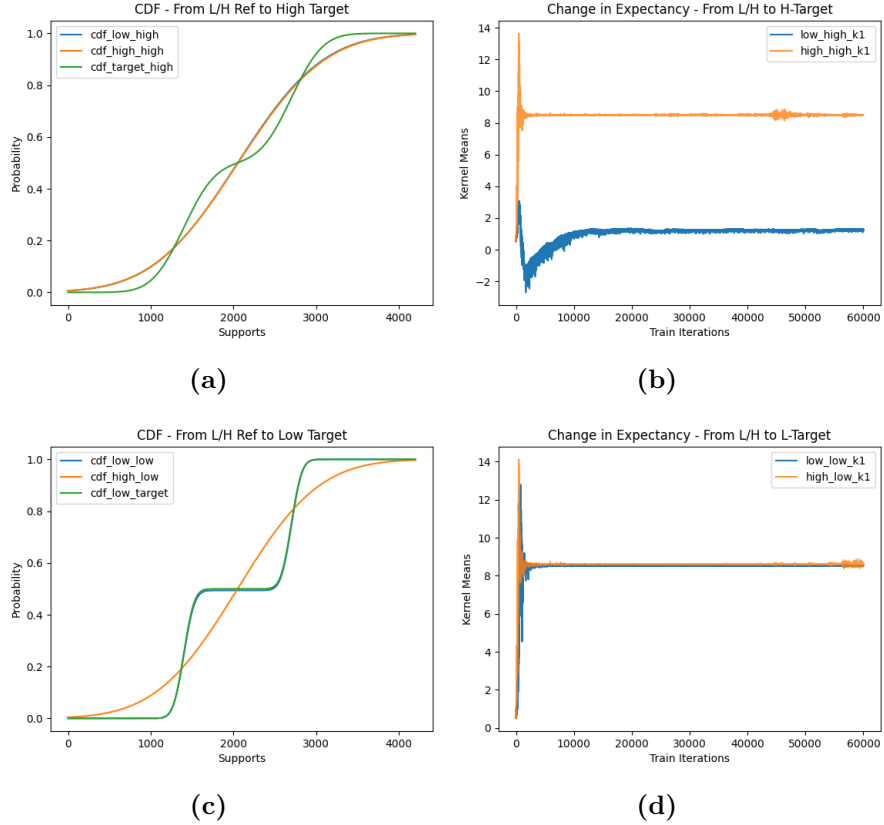
---

### Algorithm 2.8 Implicit quantile networks (with FiLM)

---

- 1: **for each** gradient step **do**
  - 2:   Sample random transition  $(s_t, a_t, r_t, s_{t+1})$  from  $\mathcal{B}$
  - 3:    $J_{IQN} \leftarrow 0$
  - 4:    $\tau[1, \dots, N] \sim \mathcal{U}([0, 1])$
  - 5:    $\tau'[1 \dots N'] \sim \mathcal{U}([0, 1])$
  - 6:   **for**  $j = 1, \dots, N'$  **do**
  - 7:     Set  $y_j = \begin{cases} r_t & \text{if done} \\ \mathcal{T}_D Z^\psi(s_t, a_t; \tau'_j) & \text{Eq.(2.69) otherwise} \end{cases}$
  - 8:      $J_{IQN} \leftarrow J_{IQN} + \frac{1}{NN'} \sum_{i=1}^N \rho_{\tau_i}^\kappa (y_j - Z_\theta^\psi(s_t, a_t; \tau_i), \quad \text{Eq (2.65) normalized}$
  - 9:      $\theta \leftarrow \theta - \beta_Z \hat{\nabla}_\theta J_{IQN}$
  - 10:    Update encoder parameters via same gradient
  - 11:     $\epsilon \leftarrow \max(\epsilon - \beta_\epsilon, 0)$
  - 12: **return**  $Z_\theta$
-

## 2.3 Distributional Reinforcement Learning



**Figure 2.3:** In this run, the neural network GMM approximator consists of two hidden layers, each containing 256 nodes. Plots (a) and (c) show fitting of the neural network GMM approximator with  $K = 2$  to a high- and low-variance target from low (blue) and high (orange) baseline distributions. In (a), the orange and blue curves coincide but both deviate from the target CDF, indicating insufficient fitting to the target. However, in (c), only the model initialized with the low variance baseline distribution (blue) is able to fit the target, while the other model (orange) fails to capture the target CDF. Figures (b) and (d) depict the change of the kernel means per iteration in the neural network fitting process.



## Distributional Reinforcement Learning in the Maximum-Entropy Setting Using the Cramér Metric

This chapter introduces a novel, state-of-the-art, Q-based (Section 2.1.6.1), distributional reinforcement learning (Section 2.3) algorithm with an entropy maximization objective (Section 2.2) in the actor-critic (Section 2.1.2.4) framework. The algorithm advances state-of-the-art performance in reinforcement learning and has a rigorous theoretical framework that ensures convergence under standard assumptions. It provides stable model updates, even in states with high epistemic uncertainty.

### 3.1 Related Work

Bellemare et al. [20] provided the first formal analysis of the distributional approach in reinforcement learning, in which the objective is to learn return distributions rather than their expectations. They employed the Wasserstein metric to measure distances between return distributions, established convergence guarantees for policy evaluation, and analyzed distributional policy improvement. It was shown that under mild assumptions, distributional policy evaluation (Section 2.3.2) and improvement (Section 2.3.3) converged to the optimal policy. This work culminated in the C51 algorithm, a critic-only method for discrete action spaces. However, as outlined in Section 2.3.4.1, C51 employs a projection step not aligned with theory.

In [82] and [83], the concepts of SAC (Section 2.2) and DistRL were applied simultaneously for the first time. The former method relies on the Kullback-Leibler divergence as a probability metric. This conflicts with the properties required for distributional policy evaluation, as the distance becomes infinite if the distributions have dissimilar supports. To address this issue, the authors applied aggressive clipping, which nega-

### 3. DISTRIBUTIONAL REINFORCEMENT LEARNING IN THE MAXIMUM-ENTROPY SETTING USING THE CRAMÉR METRIC

---

tively impacts performance. Moreover, the author of this work was unable to reproduce the results using their provided software. Ma et al. [83] extend the SAC framework by incorporating quantile-regression-based distributional reinforcement learning (Section 2.3.4.2) This integration spans several quantile-based approaches, ranging from standard QR-DQN to more complex architectures like IQN (Section 2.3.4.3) and "Fully Parameterized Quantile Function" [84]. The authors reported superior performance compared to common reinforcement learning algorithms in the MuJoCo benchmark environments with the IQN configuration. While the reliance on IQN entails an expanded parameter set via quantile embedding, it is shown that such complexity is redundant for the novel approach Cramér-based Distributional Soft Actor-Critic (C-DSAC). The results in this chapter indicate that C-DSAC's inherent risk-aversion enables superior performance through a more streamlined architecture. The work of [85] focuses on the Cramér distance "in the setting of fixed quantile levels" and does not present an algorithm for continuous control. Nam et al. [86] similarly utilize the energy distance within a DistRL framework, primarily to enhance the stability of on-policy algorithms. Their core contribution, the Sample-Replacement ( $SR(\lambda)$ ) algorithm, establishes a principled distributional generalization of the multi-step  $\lambda$ -return. In contrast, the present work pivots from on-policy stability to the distinct challenge of overestimation bias inherent in off-policy  $TD(0)$  learning.

## 3.2 The C-DSAC Framework

In the subsequent sections, the distributional reinforcement learning paradigm is extended by introducing an adaptive update rule for Q-value distributions via the regularized Cramér distance. The C-DSAC framework is developed by synthesizing the distributional reinforcement learning principles discussed in Section 2.3, specifically under the assumption of deterministic rewards, with the maximum-entropy paradigm. Furthermore, the gradients for neural network integration is derived and empirical evidence demonstrating that the novel approach achieves state-of-the-art performance across a diverse suite of benchmark environments is presented. In order to explain the efficiency of the new approach, a thorough analysis of the algorithm's dynamics is conducted. The focus is on the algorithm's behavior in presence of **epistemic uncertainty**: The critic gradient equation is derived showing that adaptation to Q-values at state action pairs with high variance, i.e. low confidence in its value, is slower. Furthermore, it is emphasized how these equations reveal the algorithm's inherent mechanism for mitigating overestimation bias.

### 3.2.1 The Squared Cramér Distance

Choosing an appropriate metric to quantify the distance between random variables  $U$  and  $V$  is critical for policy evaluation and improvement in distributional reinforcement learning.

Bellemare et al. [75] noted that a metric must be both "ideal" according to Zolotarev et al. [87] and suitable for gradient descent methods to be effective in machine learning with distributions. The squared  $L_2$  Cramér distance possesses such ideal properties and is suited for machine learning.

To formalize this, consider the proper  $L_2$  Cramér metric. For two random variables  $U$  and  $V$  it is defined as

$$d_c(U, V) := \left( \int_{-\infty}^{\infty} (F_U(x) - F_V(x))^2 dx \right)^{\frac{1}{2}}, \quad (3.1)$$

where  $F_U$  and  $F_V$  are the cumulative distribution functions of  $U$  and  $V$  respectively. Note that the distance is formally between two laws of the random variables,  $d_c(\mathcal{L}(U), \mathcal{L}(V))$ , but the common shorthand  $d_c(U, V)$  is used to avoid convoluted notation. Ideal metrics require to be sum invariant and scale sensitive. The Cramér metric is sum invariant, therefore it holds

$$d_c(A + U, A + V) \leq d_c(U, V), \quad (C1)$$

where  $A$  is a random variable independent of  $U$  and  $V$ . Treating  $A$  as a degenerate random variable with  $P(A = a) = 1$ , this property can be also stated in terms of a constant

$$d_c(a + U, a + V) \leq d_c(U, V). \quad (3.2)$$

Scale sensitivity, the second required property, ensures that for  $k > 0$ ,

$$d_c(kU, kV) = |k|^{\frac{1}{2}} d_c(U, V). \quad (C2)$$

**Lemma 3.1.** The Cramér metric  $d_c$  satisfies properties (C1) and (C2).

*Proof.* See Theorem 2 of [75]. □

A third property is needed to guarantee unbiased sample gradients. Within the class of  $L_p$  Cramér metrics, only the squared Cramér distance ( $p = 2$ ), also called the energy distance, satisfies this requirement. The energy distance is expressed as

$$d_e(U, V) := d_c^2 = \int_{-\infty}^{\infty} (F_U(x) - F_V(x))^2 dx, \quad (3.3)$$

### 3. DISTRIBUTIONAL REINFORCEMENT LEARNING IN THE MAXIMUM-ENTROPY SETTING USING THE CRAMÉR METRIC

---

and has the unbiased sample gradient property

$$\mathbb{E}_{\mathbf{X}_m \sim \Gamma} \nabla_{\theta} d_e(\hat{\Gamma}_m, Z_{\theta}) = \nabla_{\theta} d_e(\Gamma, Z_{\theta}), \quad (\text{C3})$$

where  $\mathbf{X}_m$  is a vector with elements  $X_1, \dots, X_m$  drawn from a Bernoulli distribution  $\Gamma$ , and  $\hat{\Gamma}_m := \frac{1}{m} \sum_{i=1}^m \delta_{X_i}$  is an approximate distribution formed by the samples with  $\delta$  being Dirac functions at values  $X_i$ .

**Lemma 3.2.** The energy distance  $d_e$  satisfies properties (C1),(C3), and

$$d_e(kU, kV) = |k|d_e(U, V), \quad (\text{C4})$$

with  $k > 0$ .

*Proof.* For (C3) and (C4), see Appendix A3 of [75]. Showing that the energy distance possesses the sum invariant property is rather trivial. Assume that  $A$  is independent of  $U$  and  $V$ , then

$$\begin{aligned} d_e(U + A, V + A) &= \int_{-\infty}^{\infty} (F_{U+A}(x) - F_{V+A}(x))^2 dx \\ &= \int_{-\infty}^{\infty} (\mathbb{E}_A[F_U(x - A)] - \mathbb{E}_A[F_V(x - A)])^2 dx. \end{aligned}$$

Jensen dictates that

$$\leq \int_{-\infty}^{\infty} \mathbb{E}_A[(F_U(x - A) - F_V(x - A))^2] dx.$$

Since the integrand is non-negative, Tonelli gives

$$= \mathbb{E}_A \left[ \int_{-\infty}^{\infty} (F_U(x - A) - F_V(x - A))^2 dx \right].$$

Substituting  $y = x - A$  yields

$$\begin{aligned} &= \mathbb{E}_A \left[ \int_{-\infty}^{\infty} (F_U(y) - F_V(y))^2 dy \right] \\ &= \mathbb{E}_A[d_e(U, V)] \\ &= d_e(U, V). \end{aligned}$$

□

By parity of reasoning with the derivation of Equation (3.2), this identity extends to the case of constant values.

The following Lemma provides the last property necessary to build the algorithm's theoretical framework.

**Lemma 3.3.** Let  $d_e$  be the energy distance as defined in (3.3). For any collection of distributions  $\{\mu_i, \nu_i\}_{i=1}^k$  and non-negative weights  $\{w_i\}$  such that  $\mu = \sum_i w_i \mu_i$ ,  $\nu = \sum_i w_i \nu_i$  and  $\sum_i w_i = 1$ , the following inequality holds, where for notational consistency the energy distance is written in terms of variables  $U \sim \mu$ ,  $V \sim \nu$  and  $U_i \sim \mu_i$ ,  $V_i \sim \nu_i$ , respectively

$$d_e(U, V) \leq \sum_{i=1}^k w_i d_e(U_i, V_i). \quad (\text{C5})$$

*Proof.* By Theorem 22 of Sejdinovic et al. [88], the energy distance equals a constant multiple of the RKHS norm of kernel mean embeddings. Using linearity and standard results in the RKHS yields (C5).  $\square$

### 3.2.2 Objective

The C-DSAC objective is equivalent to (2.42) if  $\mathcal{H}(\pi(\cdot | s_t)) = -\log \pi(a_t | s_t)$ , with the distributional soft state-action value made explicit,

$$J_h(\pi) = \mathbb{E}_{s_0 \sim d_0} [V_h^\pi(s_0)], \quad (3.4)$$

where

$$V_h^\pi(s_t) := \mathbb{E}_{a_t \sim \pi(\cdot | s_t)} [\mathbb{E}[Z_h^\pi(s_t, a_t)] - \alpha \log \pi(a_t | s_t)]$$

and

$$\begin{aligned} Z_h^\pi(s_t, a_t) &= r(s_t, a_t) + \gamma(Z_h^\pi(s_{t+1}, a_{t+1}) - \alpha \log \pi(a_{t+1} | s_{t+1})) \\ &\quad | s_{t+1} \sim P(\cdot | s_t, a_t), a_{t+1} \sim \pi(\cdot | s_{t+1}). \end{aligned}$$

The aim is to find a policy  $\pi^* = \operatorname{argmax}_{\pi \in \Pi} J_H(\pi)$ .

### 3.2.3 Policy Evaluation

Consider a distributional reinforcement learning setting (Section 2.3) under the squared  $L_2$  Cramér distance where the distribution over returns  $Z$  is utilized. Then the soft distributional Bellman operator is defined as

$$\begin{aligned} \mathcal{T}_H^\pi Z(s_t, a_t) &: \stackrel{d}{=} r(s_t, a_t) + \gamma(Z(s_{t+1}, a_{t+1}) - \alpha \log \pi(a_{t+1} | s_{t+1})) \\ &\quad | s_{t+1} \sim P(\cdot | s_t, a_t), a_{t+1} \sim \pi(\cdot | s_{t+1}). \end{aligned} \quad (3.5)$$

**Theorem 3.4.** Fix  $1 \leq p < \infty$  and  $(s_t, a_t) \in \mathcal{S} \times \mathcal{A}$ . Let  $\bar{d}_e(Z_1, Z_2) := \sup_{s,a} d_e(Z_1(s, a), Z_2(s, a))$  for two state-action value distributions  $Z_1$  and  $Z_2$ . Assume that for each  $(s, a)$ , the distributions have finite first moments to attain completeness of the metric space. Further,

### 3. DISTRIBUTIONAL REINFORCEMENT LEARNING IN THE MAXIMUM-ENTROPY SETTING USING THE CRAMÉR METRIC

---

assume  $|\mathcal{A}| < \infty$ ,  $|\mathcal{S}| < \infty$ ,  $0 \leq \gamma < 1$  and  $\alpha \geq 0$ . Then  $\mathcal{T}_H^\pi$  is a  $\gamma$ -contraction in the energy distance, i.e.

$$\bar{d}_e(\mathcal{T}_H^\pi Z_1, \mathcal{T}_H^\pi Z_2) \leq \gamma \bar{d}_e(Z_1, Z_2). \quad (3.6)$$

*Proof.* The reward can be augmented with the entropy term, making it a random variable. Let

$$R_h(s_t, a_t) := r(s_t, a_t) - \gamma \alpha \log \pi(a_{t+1} | s_{t+1}) | s_{t+1} \sim P(s_t, a_t), a_{t+1} \sim \pi(\cdot | s_{t+1}),$$

then

$$\begin{aligned} & d_e(\mathcal{T}_H^\pi Z_1, \mathcal{T}_H^\pi Z_2) \\ &= d_e(R_h(s_t, a_t) + \gamma Z_1(s_{t+1}, a_{t+1}), R_h(s_t, a_t) + \gamma Z_2(s_{t+1}, a_{t+1})) \\ &\quad \text{where } s_{t+1} \sim P(\cdot | s_t, a_t), a_{t+1} \sim \pi(\cdot | s_{t+1}). \end{aligned}$$

By Property (C5) with  $w_{t+1} := P(s_{t+1} | s_t, a_t) \pi(a_{t+1} | s_{t+1})$

$$\leq \sum_{s_{t+1}, a_{t+1}} w_{t+1} d_e(R_h(s_t, a_t) + \gamma Z_1(s_{t+1}, a_{t+1}), R_h(s_t, a_t) + \gamma Z_2(s_{t+1}, a_{t+1})).$$

Property (C1) holds for  $d_e$  (Lemma 3.2) and each state-action pair is fixed

$$\leq \sum_{s_{t+1}, a_{t+1}} w_{t+1} d_e(\gamma Z_1(s_{t+1}, a_{t+1}), \gamma Z_2(s_{t+1}, a_{t+1})).$$

By Property (C4)

$$= \gamma \sum_{s_{t+1}, a_{t+1}} w_{t+1} d_e(Z_1(s_{t+1}, a_{t+1}), Z_2(s_{t+1}, a_{t+1})).$$

Since  $\sum_i w_i = 1$ ,

$$\begin{aligned} & \leq \gamma \sup_{s, a} d_e(Z_1(s, a), Z_2(s, a)) \\ &= \gamma \bar{d}_e(Z_1, Z_2). \end{aligned}$$

Taking the supremum over  $(s, a) \in \mathcal{S} \times \mathcal{A}$  on both sides yields the contraction

$$\bar{d}_e(\mathcal{T}_H^\pi Z_1, \mathcal{T}_H^\pi Z_2) \leq \gamma \bar{d}_e(Z_1, Z_2). \quad \square$$

Since  $\mathcal{T}_H^\pi$  is a  $\gamma$ -contraction in maximal form, it follows from the Banach fixed point theorem:

**Corollary 3.1.** *Given  $Z_{k+1} := {}^d \mathcal{T}_H^\pi Z_k$ , the series  $\{Z_k\}$  converges to the fixed point (distribution)  $Z_h^\pi$  in the energy distance, i.e.,*

$$\lim_{k \rightarrow \infty} \bar{d}_e(Z_k, Z_h^\pi) = 0 \quad \text{with} \quad \mathcal{T}_H^\pi Z_h^\pi \stackrel{d}{=} Z_h^\pi. \quad (3.7)$$

### 3.2.4 Policy Improvement

Soft policy improvement is achieved via information projection, where the policy is updated toward the Boltzmann distribution defined by the expectation of the distributional soft state-action function  $Z_h^\pi$ ,

$$\pi_{k+1}(\cdot | s_t) = \operatorname{argmin}_{\pi(\cdot | s_t)} d_{kl} \left( \pi(\cdot | s_t) \parallel \frac{\exp(\frac{1}{\alpha} \mathbb{E}[Z_h^{\pi_k}(s_t, \cdot)])}{W^{\pi_k}(s_t)} \right). \quad (3.8)$$

It must be shown that Eq. (3.8) improves the policy.

**Lemma 3.5.** (One-step soft policy improvement)

Let  $s_t$  be fixed. If the policy is updated according to eq. (3.8),  $Z_h^\pi$  has bounded first moments,  $\alpha > 0$ ,  $0 \leq \gamma < 1$  and  $|\mathcal{A}| < \infty$ , then

$$\mathbb{E}_{a_t \sim \pi_{k+1}(\cdot | s_t)} [\mathbb{E}[Z_h^{\pi_{k+1}}(s_t, a_t)] - \alpha \log \pi_{k+1}(a_t | s_t)] \geq \mathbb{E}_{a_t \sim \pi_k(\cdot | s_t)} [\mathbb{E}[Z_h^{\pi_k}(s_t, a_t)] - \alpha \log \pi_k(a_t | s_t)].$$

*Proof.* Define

$$Q_h^\pi(s_t, a_t) := \mathbb{E}[Z_h^\pi(s_t, a_t)],$$

then

$$\begin{aligned} \pi_{k+1}(\cdot | s_t) &= \operatorname{argmin}_{\pi \in \Pi} D_{KL} \left( \pi(\cdot | s_t) \parallel \frac{\exp(\frac{1}{\alpha} Q_h^{\pi_k}(s_t, \cdot))}{W^{\pi_k}(s_t)} \right) \\ &= \operatorname{argmin}_{\pi \in \Pi} \mathbb{E}_{a_t \sim \pi(\cdot | s_t)} \left[ \log \left( \frac{\pi(a_t | s_t) W^{\pi_k}(s_t)}{\exp \frac{1}{\alpha} Q_h^{\pi_k}(s_t, a_t)} \right) \right] \\ &= \operatorname{argmin}_{\pi \in \Pi} \mathbb{E}_{a_t \sim \pi(\cdot | s_t)} \left[ \log \pi(a_t | s_t) - \frac{1}{\alpha} Q_h^{\pi_k}(s_t, a_t) + \log W^{\pi_k}(s_t) \right]. \end{aligned}$$

Due to minimization, it holds  $\mathbb{E}_{a_t \sim \pi_{k+1}(\cdot | s_t)} [\log \pi_{k+1}(a_t | s_t) - \frac{1}{\alpha} Q_h^{\pi_k}(s_t, a_t) + \log W^{\pi_k}(s_t)] \leq \mathbb{E}_{a_t \sim \pi_k(\cdot | s_t)} [\log \pi_k(a_t | s_t) - \frac{1}{\alpha} Q_h^{\pi_k}(s_t, a_t) + \log W^{\pi_k}(s_t)]$  and since the partition function does not depend on the action, it follows

$$\mathbb{E}_{a_t \sim \pi_{k+1}(\cdot | s_t)} [Q_h^{\pi_{k+1}}(s_t, a_t) - \alpha \log \pi_{k+1}(a_t | s_t)] \geq \mathbb{E}_{a_t \sim \pi_k(\cdot | s_t)} [Q_h^{\pi_k}(s_t, a_t) - \alpha \log \pi_k(a_t | s_t)].$$

Substituting the definition  $Q_h^{\pi_k}(s_t, a_t) = \mathbb{E}[Z_h^{\pi_k}(s_t, a_t)]$  recovers the statement of the lemma.  $\square$

The expectation of the distributional state-action value is expanded to show that  $\pi_{k+1}$  is indeed an improvement on  $\pi_k$

### 3. DISTRIBUTIONAL REINFORCEMENT LEARNING IN THE MAXIMUM-ENTROPY SETTING USING THE CRAMÉR METRIC

---

**Lemma 3.6.** (Global soft policy improvement)

Assume  $|\mathcal{A}| < \infty$ ,  $\alpha > 0$ ,  $0 \leq \gamma < 1$  and bounded first moments of  $Z_h^\pi$ . If  $\pi_{k+1}$  is obtained according to the soft policy improvement in (2.45), then

$$\mathbb{E}[Z_h^{\pi_k}(s_t, a_t)] \leq \mathbb{E}[Z_h^{\pi_{k+1}}(s_t, a_t)] \quad \forall (s_t, a_t) \in \mathcal{S} \times \mathcal{A},$$

when applying Lemma 3.5 pointwise at each successor state.

*Proof.* The proof relies on expanding the expectation of the soft Z-function in the Bellman equation repeatedly with actions of the improved policy. For convenience, define  $\pi' := \pi_{k+1}$ ,  $\pi := \pi_k$  and  $\mathbb{E}_{\pi', P, k}[\cdot | s_t] := \mathbb{E}_{a_{t:k-1} \sim \pi'(\cdot | s_{t:k-1}), s_{t+1:k} \sim P(\cdot | s_{t:k-1}, a_{t:k-1})}[\cdot | s_t]$ . Define

$$Q_h^\pi(s_t, a_t) := \mathbb{E}[Z_h^\pi(s_t, a_t)],$$

then applying the one-step improvement from Lemma 3.5 at each successor state yields

$$\begin{aligned} V_h^\pi(s_t) &\leq \mathbb{E}_{\pi', P, t+1}[\mathbb{E}[Z_h^\pi(s_t, a_t)] - \alpha \log \pi'(a_t | s_t) | s_t] \\ &= \mathbb{E}_{\pi', P, t+1}[r(s_t, a_t) - \alpha \log \pi'(a_t | s_t) + \gamma V_h^\pi(s_{t+1}) | s_t] \\ &\leq \mathbb{E}_{\pi', P, t+2}[r(s_t, a_t) - \alpha \log \pi'(a_t | s_t) + \gamma(\mathbb{E}[Z_h^\pi(s_{t+1}, a_{t+1})] - \alpha \log \pi'(a_{t+1} | s_{t+1})) | s_t] \\ &= \mathbb{E}_{\pi', P, t+2}[r(s_t, a_t) - \alpha \log \pi'(a_t | s_t) + \gamma(r(s_{t+1}, a_{t+1}) - \alpha \log \pi'(a_{t+1} | s_{t+1})) + \\ &\quad \gamma^2 V_h^\pi(s_{t+2}) | s_t] \\ &\vdots \\ &\leq \mathbb{E}_{\pi', P, T-1}[r(s_t, a_t) - \alpha \log \pi'(a_t | s_t) + \gamma(r(s_{t+1}, a_{t+1}) - \alpha \log \pi'(a_{t+1} | s_{t+1})) + \dots \\ &\quad \dots + \gamma^{T-1}(r(s_{T-1}, a_{T-1}) - \alpha \log \pi'(a_{T-1} | s_{T-1})) | s_t = s] \\ &= V_h^{\pi'}(s_t). \end{aligned}$$

Thus for any  $s_t$

$$\begin{aligned} \mathbb{E}[Z_h^\pi(s_t, a_t)] &= \mathbb{E}_{s_{t+1} \sim P(\cdot | s_t, a_t)}[r(s_t, a_t) + \gamma V_h^{\pi_k}(s_{t+1})] \\ &\leq \mathbb{E}_{s_{t+1} \sim P(\cdot | s_t, a_t)}[r(s_t, a_t) + \gamma V_h^{\pi_{k+1}}(s_{t+1})] = \mathbb{E}[Z_h^{\pi_{k+1}}(s_t, a_t)]. \end{aligned}$$

□

#### 3.2.5 Policy Iteration

The following proof mirrors the one of Lemma 2.9 but accounts for the distributional nature of the variables through the expectation operator. Note that the following result does not claim distributional optimality.

**Theorem 3.7.** *Let the reward be bounded,  $|\mathcal{A}| < \infty$ ,  $|\mathcal{S}| < \infty$ ,  $0 \leq \gamma < 1$  and  $\alpha > 0$ . Let  $\Pi$  be a set of all stationary, stochastic policies. Alternating between exact distributional soft policy evaluation and global soft policy improvement from some initial policy  $\pi_0 \in \Pi$  and soft state-action function  $\mathbb{E}[Z_h^{(0)}]$ , the process converges (in the limit) to the optimal policy  $\pi^*$ , satisfying  $\mathbb{E}[Z_h^{\pi^*}] = \mathbb{E}[Z_h^*]$ .*

*Proof.* At each iteration  $k$ , the policy  $\pi_k$  is fixed during evaluation. By Theorem 3.4 and its Corollary,  $\mathcal{T}_H^{\pi_k}$  is a  $\gamma$ -contraction in the energy distance and the evaluation sequence converges to the fixed point  $Z_h^{\pi_k}$ . Therefore,  $\mathbb{E}[Z_h^{\pi_k}]$  is well defined at each iteration. By Lemma 3.6, the sequence  $\{\mathbb{E}[Z_h^{\pi_k}(s_t, a_t)]\}$  is non-decreasing for each pair  $(s_t, a_t)$ . Because the reward is bounded and  $|\mathcal{S}| < \infty$ ,  $|\mathcal{A}| < \infty$ , the sequence converges to some point-wise limit  $\bar{L}$ . Policies are updated by the softmax of the current  $\mathbb{E}[Z_h]$  and softmax is continuous, therefore the limit point  $\pi^*$  exists and is the softmax of  $\bar{L}$ , i.e.  $\pi^* = \text{softmax}(\bar{L})$ . Since  $\pi^* \in \Pi$ , Theorem 3.4 equally guarantees well-definedness in the limit point  $\mathbb{E}[Z_h^{\pi^*}]$ . The map  $\pi \mapsto \mathbb{E}[Z_h^\pi]$  is continuous, hence  $\pi_k \rightarrow \pi^*$  implies  $\mathbb{E}[Z_h^{\pi_k}] \rightarrow \mathbb{E}[Z_h^{\pi^*}]$  and therefore  $\bar{L} = \mathbb{E}[Z_h^{\pi^*}]$ . It follows  $\pi^* = \text{softmax}(\mathbb{E}[Z_h^{\pi^*}])$  and therefore  $\mathbb{E}[Z_h^{\pi^*}]$  satisfies the Bellman optimality equation, thus  $\mathbb{E}[Z_h^{\pi^*}] = \mathbb{E}[Z_h^*]$ .  $\square$

The theorems above provide a theoretical basis for a Cramér-based distributional SAC approach (C-DSAC). In Section 3.3 we use these insights to construct the C-DSAC algorithm.

### 3.3 Implementation

In this section, the theoretical results in the previous sections are utilized to derive parameterized objective functions suitable for implementation. To handle the large state and action spaces in complex tasks, neural networks are used as parameterized function approximators; the aim is to optimize their parameters w.r.t the objective. The parameters for the random variable and policy will be denoted by  $\theta$  and  $\phi$ , respectively. It is assumed that the random variable obeys a Gaussian distribution with expectation  $Q_{\theta_1}$  and standard deviation  $\sigma_{\theta_2}$ , i.e. C-DSAC is implemented as a fixed-moment algorithm as described in Section 2.3.5.1.

#### 3.3.1 Objective Functions and Algorithm

Formulas are derived for implementation with parameterized function approximators. Neural networks are utilized for their high generalization capabilities. Previously observed states and actions are stored in a replay buffer  $\mathcal{B}$ . Implementing the Double-Q

### 3. DISTRIBUTIONAL REINFORCEMENT LEARNING IN THE MAXIMUM-ENTROPY SETTING USING THE CRAMÉR METRIC

---

method described in Section 2.1.5.1, the target network  $\bar{\theta}$  is employed for the critic to further reduce approximation errors.

**Evaluation** The neural network parameters of  $Z_\theta$  are optimized by minimizing the energy distance relative to the soft distributional Bellman target

$$J_Z(\theta) := \mathbb{E}_{\substack{s_t \sim d_\mu^\gamma \\ a_t \sim \mu(\cdot | s_t) \\ s_{t+1} \sim P(\cdot | s_t, a_t)}} [d_e(\mathcal{T}_H^\pi Z(s_t, a_t), Z_\theta(s_t, a_t))], \quad (3.9)$$

where  $d_\mu^\gamma = \frac{1-\gamma}{1-\gamma^T} \sum_{t=0}^{T-1} \gamma^t Pr(s_t = s | \mu)$  is the state-visitation distribution induced by the stochastic behavior policy  $\mu$ .

In practice, gradients are estimated using a mini-batch  $\hat{\mathcal{B}}$  sampled uniformly from the replay buffer

$$\hat{\nabla}_\theta J_Z(\theta) = \frac{1}{|\hat{\mathcal{B}}|} \sum_{(s_t, a_t, r_t, s_{t+1}) \in \hat{\mathcal{B}}} [\nabla_\theta d_e(\hat{\mathcal{T}}_H^\pi Z_{\bar{\theta}}(s_t, a_t), Z_\theta(s_t, a_t))], \quad (3.10)$$

where the empirical soft distributional Bellman operator for a fixed  $s_{t+1}$  is defined as

$$\hat{\mathcal{T}}_H^\pi Z_{\bar{\theta}}(s_t, a_t) := r(s_t, a_t) + \gamma(Z_{\bar{\theta}}(s_{t+1}, a_{t+1}) - \alpha \log \pi(a_{t+1} | s_{t+1})) | a_{t+1} \sim \pi(\cdot | s_{t+1})$$

and  $\bar{\theta}$  are frozen parameters (for clarity,  $\bar{\theta}$  will be omitted in the subsequent sections). The target network parameters are not updated via gradient descent; instead, it tracks the online network through Polyak averaging. In order to simplify the computation of the value distribution, it is assumed that its probability density function is Gaussian denoted by

$$\varphi_{Q,\sigma}(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{x-Q}{\sigma}\right)^2}$$

regarding a mean value  $Q \in [Q, \bar{Q}] \subset \mathbb{R}$  and a variance  $\sigma \in [\underline{\sigma}, \bar{\sigma}] \subset \mathbb{R}_+$ . For terminal states, the standard deviation of the target is set to 0. The cumulative distribution function of  $\varphi_{Q,\sigma}$  is denoted by

$$F_{Q,\sigma}(x) = \int_{-\infty}^x \varphi_{Q,\sigma}(t) dt.$$

Let  $Q_{\theta_1}(s_t, a_t)$  and  $\sigma_{\theta_2}(s_t, a_t)$  be neural network models of  $Q$  and  $\sigma$  over  $\mathcal{S} \times \mathcal{A}$  respectively. Then  $Z_\theta(s_t, a_t)$  denotes a parameterized Gaussian value distribution with parameters  $\theta = (\theta_1, \theta_2)$ , defined by the distribution function

$$f_{Z_\theta(s_t, a_t)} = \varphi_{Q_{\theta_1}(s_t, a_t), \sigma_{\theta_2}(s_t, a_t)}.$$

**Improvement** The policy improvement follows from the information projection equation (3.8). Similar to [30] and [72], a reparameterized policy is used to obtain the actions,  $a_t = g_\phi(s_t; \xi_t)$ , where  $\xi_t \sim \mathcal{N}(0, 1)$ . The policy  $\pi_\phi(\cdot | s_t)$  is defined as the pushforward of  $g_\phi(s_t; \cdot)$  under the law of  $\xi_t$ . The policy loss can then be defined as

$$J_\pi(\phi) := \mathbb{E}_{\substack{s_t \sim d_\mu^\gamma, \\ \xi_t \sim \mathcal{N}(0,1) \\ a_t = g_\phi(s_t; \xi_t)}} [\alpha \log \pi_\phi(a_t | s_t) - \mathbb{E}[Z_\theta(s_t, a_t)]], \quad (3.11)$$

where  $d_\mu^\gamma = \frac{1-\gamma}{1-\gamma^T} \sum_{t=0}^{T-1} \gamma^t Pr(s_t = s | \mu)$  is the state-visitation distribution induced by the stochastic behavior policy  $\mu$ . The gradient of the empirical loss regarding a mini-batch  $\hat{\mathcal{B}} \subset \mathcal{B}$  is thus

$$\hat{\nabla}_\phi J_\pi(\phi) = \frac{1}{|\hat{\mathcal{B}}|} \sum_{\substack{s_t \in \hat{\mathcal{B}}, \\ \xi \sim \mathcal{N}(0,1)}} \alpha \partial_\phi \log \pi_\phi(a_t | s_t)|_{a_t} + (\alpha \nabla_{a_t} \log \pi_\phi(a_t | s_t) - \nabla_{a_t} \mathbb{E}[Z_\theta(s_t, a_t)]) \nabla_\phi g_\phi(s_t; \xi). \quad (3.12)$$

Similar to the distributional value function, the action distribution is modeled as a Gaussian. Additionally, the law of  $\xi$  for each action component is set to be a standard Gauss distribution.

**Algorithm** The C-DSAC algorithm is presented in Algorithm 3.1, with the Cramér gradient update described in line 13.

## 3.4 Analysis

The effect of the energy distance is investigated on the algorithm dynamics and the aim is to provide a possible explanation for its observed state-of-the-art performance. As discussed in the following sections, C-DSAC has an inherent mechanism to counteract overestimation bias. It is shown that the update of the neural network  $Q_{\theta_1}(s_t, a_t)$  using gradients  $\nabla_{\theta_1} J_Z(\theta)$  of the Cramér loss function  $J_Z(\theta)$  improves the  $Q$ -model most where the variance  $\sigma_{\theta_2}(s_t, a_t)$ , which is a measure of the uncertainty of  $Z_\theta$ , is small.

### 3.4.1 Confidence-Driven Model Updates

Fix a target distribution  $\varphi_{Q', \sigma'}$ , where  $Q'$  and  $\sigma'$  represent the Gaussian parameters. The cost function  $C(Q, \sigma)$  measures the energy distance between a candidate distribution and this target

$$C(Q, \sigma) := d_e(\varphi_{Q, \sigma}, \varphi_{Q', \sigma'}) = \int_{-\infty}^{\infty} (F_{Q, \sigma}(x) - F_{Q', \sigma'}(x))^2 dx.$$

### 3. DISTRIBUTIONAL REINFORCEMENT LEARNING IN THE MAXIMUM-ENTROPY SETTING USING THE CRAMÉR METRIC

---



---

**Algorithm 3.1** C-DSAC (Cramér-based Distributional Soft Actor-Critic)

---

- 1: Initialize parameters  $\theta, \bar{\theta}, \phi$
  - 2: Initialize learning rates  $\beta_Z, \beta_\pi$
  - 3:  $\mathcal{B} \leftarrow \emptyset$
  - 4: **repeat**
  - 5: Receive initial observation state  $s_0$
  - 6: **for each** environment step **do**
  - 7: Select action  $a_t \sim \pi_\phi(\cdot | s_t)$
  - 8: Observe transition  $s_{t+1} \sim P(\cdot | s_t, a_t)$
  - 9: Observe reward  $r_t$
  - 10:  $\mathcal{B} \leftarrow \mathcal{B} \cup \{(s_t, a_t, r_t, s_{t+1})\}$
  - 11: **for each** gradient step **do**
  - 12: Uniformly sample a mini-batch of  $N$  transitions  $(s_i, a_i, r_i, s_{i+1})$
  - 13:  $\theta \leftarrow \theta + \beta_Z \hat{\nabla}_\theta J_Z(\theta)$ , (Eq. (3.10))
  - 14:  $\phi \leftarrow \phi + \beta_\pi \hat{\nabla}_\phi J_\pi(\phi)$ , (Eq. (3.12))
  - 15: Update target network:  $\bar{\theta} \leftarrow \tau\theta + (1 - \tau)\bar{\theta}$
  - 16: **until** Stopping criterion
  - 17: **return**  $Z_\theta, \pi_\phi$
-

**Lemma 3.8.** Let  $\varphi_{Q,\sigma}$  and  $\varphi_{Q',\sigma'}$  be a current and a target distribution, respectively. Then

$$\frac{\partial}{\partial Q} C(Q, \sigma) = -\frac{2}{\sigma} B(Q, \sigma), \quad (3.13)$$

where

$$B(Q, \sigma) := \int_{-\infty}^{\infty} (F_{Q,\sigma}(x) - F_{Q',\sigma'}(x)) \varphi_{Q,\sigma}(x) dx. \quad (3.14)$$

*Proof.* Let  $f(x) = \varphi_{(0,1)}$  be the standard normal density function and  $\Phi = F_{0,1}$  the standard normal cumulative distribution function. Applying the chain rule yields

$$\frac{\partial}{\partial Q} F_{Q,\sigma} = \frac{\partial}{\partial Q} \Phi\left(\frac{x-Q}{\sigma}\right) = -\frac{1}{\sigma} f\left(\frac{x-Q}{\sigma}\right) = -\frac{\varphi_{Q,\sigma}}{\sigma}.$$

Hence,

$$\begin{aligned} \frac{\partial}{\partial Q} C(Q, \sigma) &= \int_{-\infty}^{\infty} 2(F_{Q,\sigma}(x) - F_{Q',\sigma'}(x)) \frac{\partial}{\partial Q} F_{Q,\sigma}(x) dx \\ &= \int_{-\infty}^{\infty} -2(F_{Q,\sigma}(x) - F_{Q',\sigma'}(x)) \frac{\varphi_{Q,\sigma}(x)}{\sigma} dx \\ &= -\frac{2}{\sigma} \int_{-\infty}^{\infty} (F_{Q,\sigma}(x) - F_{Q',\sigma'}(x)) \varphi_{Q,\sigma}(x) dx. \end{aligned}$$

□

The following Lemma and Proposition establish the core result of the analysis.

**Lemma 3.9.** Let  $\varphi_{Q,\sigma}$  and  $\varphi_{Q',\sigma'}$  be a current and a target distribution, respectively. Then

$$\lim_{\sigma \rightarrow \infty} \frac{\partial}{\partial Q} C(Q, \sigma) = 0.$$

*Proof.* Since  $F_{Q,\sigma}(x) \in [0, 1]$  and  $\|\varphi_{Q,\sigma}\|_1 = 1$ , it holds from the Hölder inequality

$$|B(Q, \sigma)| \leq \|(F_{Q,\sigma} - F_{Q',\sigma'}) \cdot \varphi_{Q,\sigma}\|_1 \leq \|F_{Q,\sigma} - F_{Q',\sigma'}\|_{\infty} \cdot \|\varphi_{Q,\sigma}\|_1 \leq 1.$$

This gives

$$\lim_{\sigma \rightarrow \infty} \left| \frac{\partial}{\partial Q} C(Q, \sigma) \right| = \lim_{\sigma \rightarrow \infty} \frac{2}{\sigma} \cdot |B(Q, \sigma)| \leq \lim_{\sigma \rightarrow \infty} \frac{2}{\sigma} = 0.$$

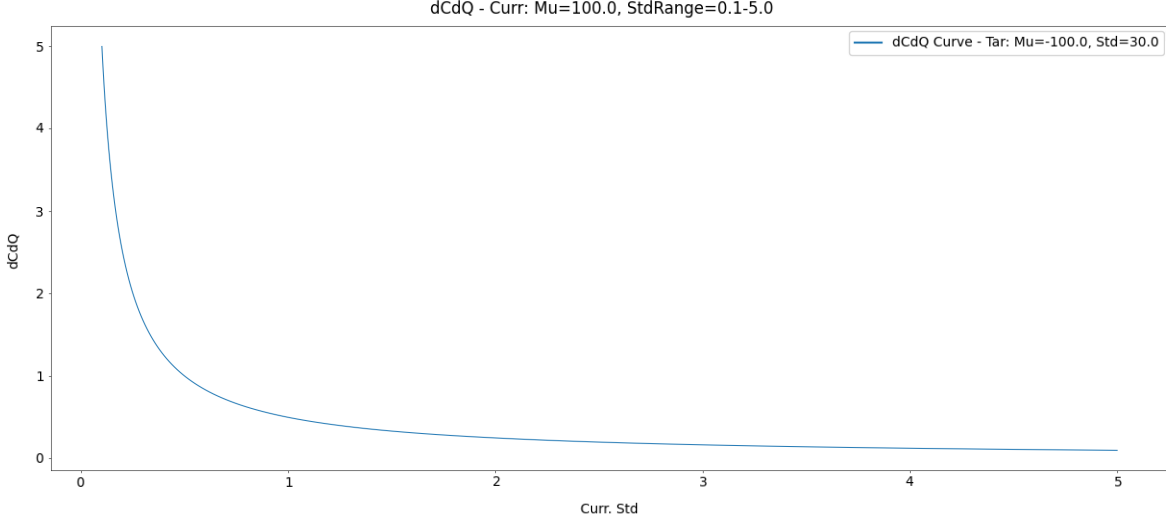
□

**Proposition 3.10.** Denote by  $\Psi_{\theta}(s_t, a_t) := \frac{\partial}{\partial \theta_{\beta_1}(s_t, a_t)} d_e(\hat{\mathcal{T}}_H^{\pi} Z(s_t, a_t), Z_{\theta}(s_t, a_t))$  a gradient weight of Eq. (3.10) such that

$\hat{\nabla}_{\theta_1} J_Z(\theta) = \frac{1}{|\beta|} \sum_{\beta} [\Psi_{\theta}(s_t, a_t) \nabla_{\theta_1} Q_{\theta_1}(s_t, a_t)]$ . Then

$$\lim_{\sigma_{\theta_2}(s_t, a_t) \rightarrow \infty} \Psi_{\theta}(s_t, a_t) = 0. \quad (3.15)$$

### 3. DISTRIBUTIONAL REINFORCEMENT LEARNING IN THE MAXIMUM-ENTROPY SETTING USING THE CRAMÉR METRIC



**Figure 3.1:** Value of  $\frac{\partial}{\partial Q}C(Q, \sigma)$  for varying  $\sigma \in [\underline{\sigma}, \bar{\sigma}]$

*Proof.* From Lemma 3.8 it follows

$$\Psi_{\theta}(s_t, a_t) = \frac{\partial}{\partial Q_{\theta_1}(s_t, a_t)} d_e(\hat{T}_H^{\pi} Z(s_t, a_t), Z_{\theta}(s_t, a_t)) = -\frac{2B(Q_{\theta_1}(s_t, a_t), \sigma_{\theta_2}(s_t, a_t))}{\sigma_{\theta_2}(s_t, a_t)}.$$

The statement follows from Lemma 3.9.  $\square$

Figure 3.1 shows that the gradient weight  $\Psi_{\theta}(s_t, a_t)$  is quickly reducing, if the variance  $\sigma_{\theta_2}(s_t, a_t)$  is increased. The analysis reveals that for a particular pair  $(s_t, a_t)$ , the magnitude of the  $Q_{\theta_1}$  update is reversely scaled by the variance  $\sigma_{\theta_2}$ , such that the update vanishes as  $\sigma_{\theta_2} \rightarrow \infty$ . Hence,  $Q_{\theta_1}(s_t, a_t)$  is updated more aggressively at state-action pairs  $(s_t, a_t)$  where  $\sigma_{\theta_2}(s_t, a_t)$  is small, reflecting temporal stability of the target. This update type is therefore called *confidence-driven*.

#### 3.4.2 Effect on Overestimation Bias

A pervasive challenge in  $Q$ -based reinforcement learning is the overestimation bias that arises from function approximation in noisy environments [46]. In this section, it is shown that the influence of the overestimation error at  $(s_t, a_t)$  on the gradient of the energy loss function (3.10) is small, if the variance  $\sigma_{\theta_2}(s_t, a_t)$  of  $Z_{\theta}(s_t, a_t)$  is large. Denote by  $\tilde{Q}$  and  $Q^*$  a noisy and an exact mean ( $Q$ -value), respectively. Consider a noisy target

$$\tilde{Y}(s_t, a_t) := r(s_t, a_t) + \gamma Z(s_{t+1}, a_{t+1}), \quad a_{t+1} = \operatorname{argmax}_{a \in \mathcal{A}} \tilde{Q}(s_{t+1}, a)$$

and an exact target

$$Y^*(s_t, a_t) := r(s_t, a_t) + \gamma Z(s_{t+1}, a_{t+1}), \quad a_{t+1} = \operatorname{argmax}_{a \in \mathcal{A}} Q^*(s_{t+1}, a).$$

**Proposition 3.11.** Let  $\Delta\Psi_\theta(s_t, a_t) := |\tilde{\Psi}(s_t, a_t) - \Psi^*(s_t, a_t)|$  be the gradient weight error of the energy distance (3.10) regarding noisy and exact gradient weights, where

$$\tilde{\Psi}(s_t, a_t) := \frac{\partial}{\partial Q_{\theta_1}} d_c(Z_\theta(s_t, a_t), \tilde{Y}(s_t, a_t)), \quad \Psi^*(s_t, a_t) := \frac{\partial}{\partial Q_{\theta_1}} d_c(Z_\theta(s_t, a_t), Y^*(s_t, a_t)).$$

Then

$$\lim_{\sigma_{\theta_2}(s_t, a_t) \rightarrow \infty} \Delta\Psi_\theta(s_t, a_t) = 0. \quad (3.16)$$

*Proof.* Define  $B^*$  and  $\tilde{B}$  via (3.14) using  $Y^*$  and  $\tilde{Y}$ , respectively. From Lemma 3.8 it follows

$$\left| \tilde{\Psi}(s_t, a_t) - \Psi^*(s_t, a_t) \right| = 2 \left| \frac{\tilde{B}_\theta(s_t, a_t)}{\sigma_{\theta_2}(s_t, a_t)} - \frac{B_\theta^*(s_t, a_t)}{\sigma_{\theta_2}(s_t, a_t)} \right| \leq \frac{2}{\sigma_{\theta_2}(s_t, a_t)},$$

since  $|B_\theta^*(s_t, a_t)| \leq 1$  and  $|\tilde{B}_\theta(s_t, a_t)| \leq 1$ . This proves (3.16).  $\square$

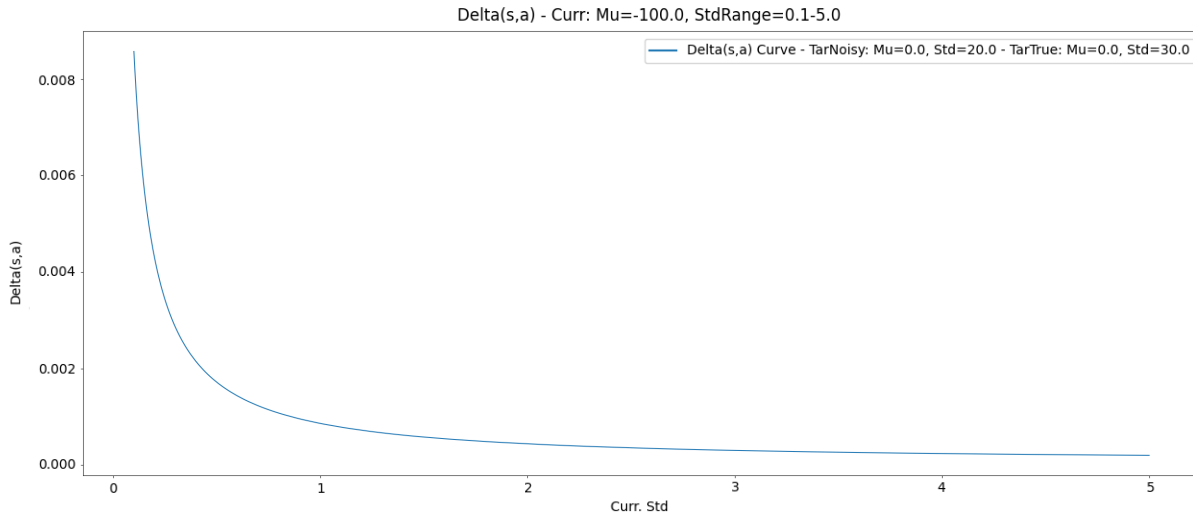
Figure 3.2 shows that the gradient weight error  $\Delta\Psi_\theta(s_t, a_t)$  is quickly reducing if the variance  $\sigma_{\theta_2}(s_t, a_t)$  is increased. The variance is elevated at state-action pairs exhibiting high return stochasticity, coinciding with the regions most susceptible to value overestimation [89].

## 3.5 Experiments

The experiments are conducted with the same hyperparameter values as in [30], except for C-DSAC-specific ones (e.g, limits on standard deviation, ... etc. ), see Table 3.1. Implementation-level nuances and low-level optimizations frequently exert a substantial influence on an algorithm’s empirical performance [90]. It is assumed that the original implementation used for the baseline results in [30] did not employ additional optimization techniques such as value function clipping, orthogonal initialization, layer scaling, learning rate annealing, reward/observation clipping, or observation and reward normalization. To be consistent with these assumptions, these code-level optimization techniques is not employed in the C-DSAC implementation. Employing the same network architecture as the SAC experiments would grant C-DSAC a disproportionate advantage, owing to the structural requirements of its dual-output critic. Therefore, the last hidden layer is reduced by one neuron. Consequently, the C-DSAC critic has

### 3. DISTRIBUTIONAL REINFORCEMENT LEARNING IN THE MAXIMUM-ENTROPY SETTING USING THE CRAMÉR METRIC

---



**Figure 3.2:** Value of  $\Delta\Psi$  for varying  $\sigma \in [\underline{\sigma}, \bar{\sigma}]$

slightly fewer parameters than those reported in the original SAC configuration, see Appendix 3.5.1. Although [83] and [82] claim superior performance over SAC baselines, it is important to note that the former uses more complex neural network architectures with batch normalization for the value approximator.

#### 3.5.1 Hyperparameters

Hyperparameters for the C-DSAC experiments were systematically tuned to maintain parity with the SAC baselines, ensuring an equitable comparison.

**Table 3.1:** C-DSAC Hyperparameters

Parameter	Value
<i>Shared</i>	
Optimizer	Adam [65]
Learning rate	$3 \cdot 10^{-4}$
Discount factor ( $\gamma$ )	0.99
Replay buffer size	$10^6$
Actor architecture	(256, 256)
Minibatch size	256
Layer activation	ReLU
Target smoothing coefficient ( $\tau$ )	0.005
Target update interval	1
Gradient steps	1
<i>SAC</i>	
Critic architecture	(256, 256)
<i>C-DSAC</i>	
Critic architecture	(256, 255)
Min., Max. Std.	0.01, 1000
Gauss grid points	31
Gauss bounds	$15\sigma$

**Table 3.2:** Environment-specific Entropy

Env. Name	$\alpha$ -value
Hopper-v4	0.2
Ant-v4	0.2
HalfCheetah-v4	0.2
Walker2d-v4	0.2
Humanoid-v4	0.05

### 3.5.2 Comparable Evaluation with SAC Baselines

The goal of the experiments is to evaluate the impact of the energy distance on the entropy-augmented actor-critic framework across diverse and complex environments [91], see Figure 3.3: Hopper-v4,  $(\mathcal{S} \times \mathcal{A}) \in \mathbb{R}^{11} \times \mathbb{R}^3$ ; Ant-v4  $(\mathcal{S} \times \mathcal{A}) \in \mathbb{R}^{111} \times \mathbb{R}^8$ ; Humanoid-v4  $(\mathcal{S} \times \mathcal{A}) \in \mathbb{R}^{376} \times \mathbb{R}^{17}$ ; HalfCheetah-v4  $(\mathcal{S} \times \mathcal{A}) \in \mathbb{R}^{17} \times \mathbb{R}^6$ ; Walker2d-v4

### 3. DISTRIBUTIONAL REINFORCEMENT LEARNING IN THE MAXIMUM-ENTROPY SETTING USING THE CRAMÉR METRIC

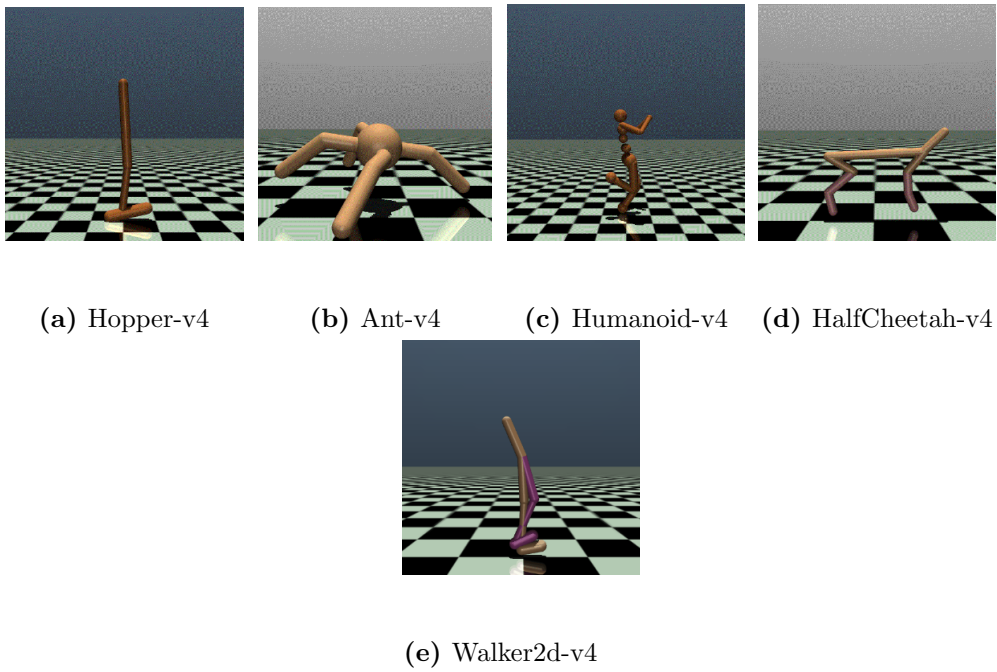
---

$(\mathcal{S} \times \mathcal{A}) \in \mathbb{R}^{17} \times \mathbb{R}^6$ . The entropy coefficient is constant for each environment, following the configuration in [30]. Environment-specific values for  $\alpha$  are given in Table 3.2. Following the recommendations of [92], evaluation rollouts are performed every  $10^3$  iterations. During evaluation, agents follow a deterministic policy derived from the distribution mean. The agent is trained for one million iterations in each environment. The results are averaged over five random seeds. The graphs in Figure 3.4 showcase (among others) the evaluation performances of C-DSAC and SAC in the benchmark environments. Especially in the most difficult environment, Humanoid-v4, C-DSAC clearly outperforms SAC and reaches an average reward of about 5730, a value that is achieved by SAC only after approximately four million iterations. Notably, C-DSAC performance has not yet reached a plateau by the conclusion of training, suggesting further potential for improvement with extended computation. Although C-DSAC and SAC reach the same average rollout in the Hopper-v4 environment, C-DSAC does so considerably faster. The instability after  $4 \times 10^5$  iterations can be due to the occurrence of catastrophic forgetting in one of the runs. However, it manages to recover. In the Walker2d-v4 environment, C-DSAC attains higher mean rewards, but the training profile reveals greater stochasticity relative to the baseline. The least stable training occurs in the Ant-v4 environment, where high stochasticity can be observed. Still, C-DSAC outperforms SAC in Ant-v4. In HalfCheetah-v4, C-DSAC obtains an average reward of less than 10,000 in the end of training and is therefore outperformed by SAC in this environment.

To provide more details about the performance, Table 3.3 presents the best model performances across all environments, reporting average episodic rewards over 100 runs.

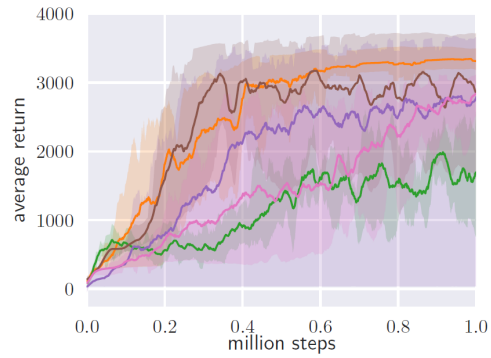
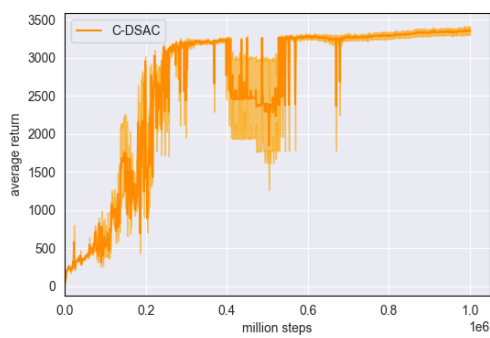
**Table 3.3:** Average performance over 100 runs of the best models in the training runs

Environment	Best Performance
Hopper-v4	Target: $3352 \pm 97$
Ant-4	Target: $4375 \pm 299$ Online: $5376 \pm 79$
Humanoid-v4	Target: $5715 \pm 255$
HalfCheetah-4	Target: $9124 \pm 309$ Online: $10023 \pm 228$
Walker-v4	Target: $4808 \pm 223$

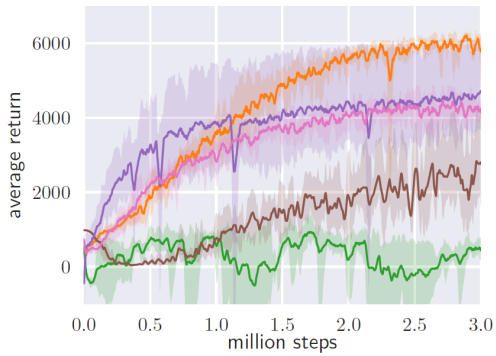
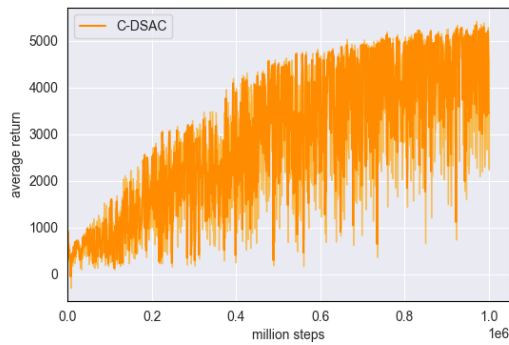


**Figure 3.3:** OpenAI Gym robotic testing environments

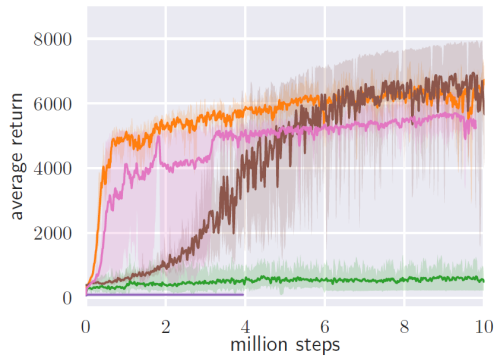
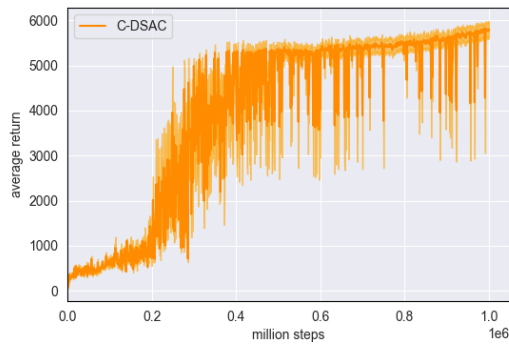
### 3. DISTRIBUTIONAL REINFORCEMENT LEARNING IN THE MAXIMUM-ENTROPY SETTING USING THE CRAMÉR METRIC



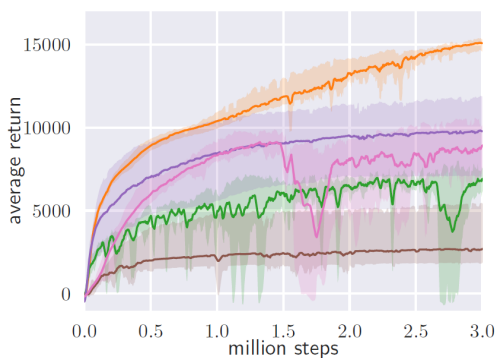
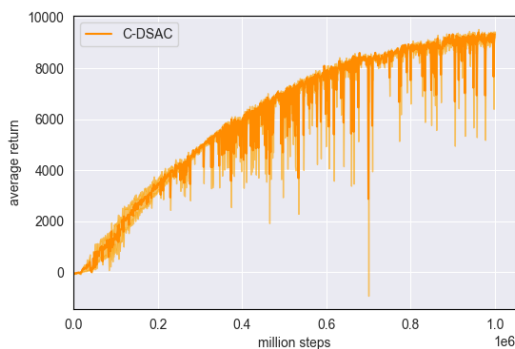
(a) Hopper



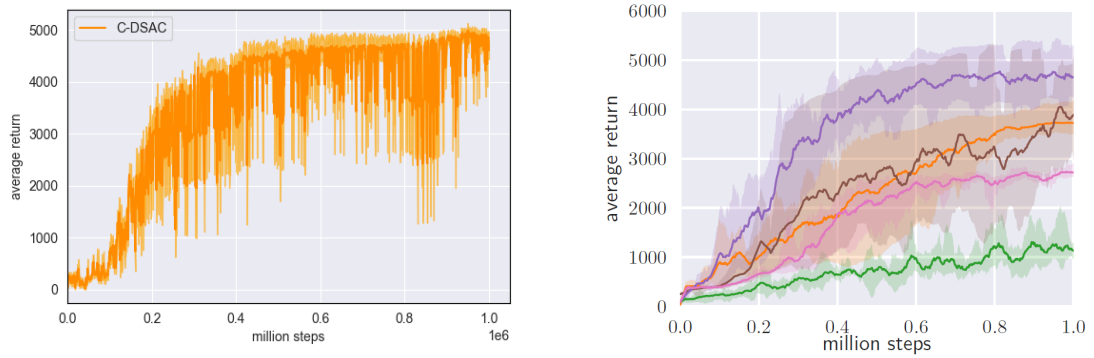
(b) Ant (left one million, right three million steps)



(c) Humanoid (left one million, right ten million steps)



(d) HalfCheetah (left one million, right three million steps)



(e) Walker

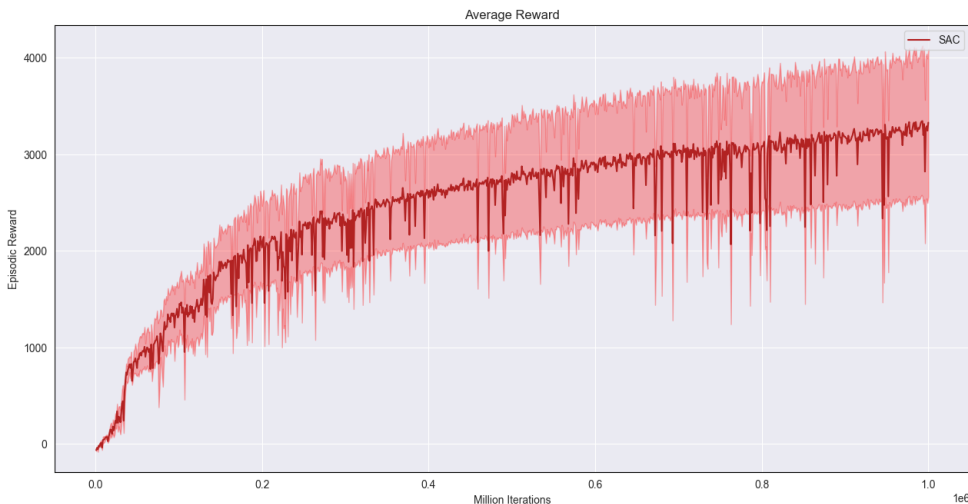
**Figure 3.4:** Orange curves represent C-DSAC (left) and SAC (right, from [30]). Other curves show performances of DDPG (green), PPO (brown), SQL (pink), and TD3 (violet). Shaded regions indicate  $\pm 1$  standard deviation across 5 independent trials.

### 3. DISTRIBUTIONAL REINFORCEMENT LEARNING IN THE MAXIMUM-ENTROPY SETTING USING THE CRAMÉR METRIC

---

#### 3.5.3 Comparative Implementation Evaluation against a Standard Baseline

To assess the quality of the C-DSAC implementation - and therefore its effect on the performance, its code is modified to revert it to SAC by adjusting the loss function and critic network architecture. The experiments are repeated on HalfCheetah-v4 with five runs, Figure 3.5. The SAC hyperparameter are set to values given in Appendix 3.5.1. The custom SAC implementation based on the C-DSAC code achieved an average reward of slightly above 3000, compared to over 10000 in the original SAC implementation after one million iterations. The observed performance gap potentially stems from unoptimized implementation details, suggesting that C-DSAC has yet to reach its theoretical performance ceiling.



**Figure 3.5:** Performance of SAC based on C-DSAC implementation on HalfCheetah-v4

## 3.6 Conclusion

This chapter presented C-DSAC, a distributional maximum-entropy deep reinforcement learning algorithm. The objective function incorporates an energy distance-based metric to estimate the expected Q-values and their associated variances. It has been formally established that distributional soft policy iteration, utilizing the energy distance metric for evaluation alongside standard policy improvement, converges to an optimal

policy. Numerical experiments show that C-DSAC outperforms SAC in most benchmark environments despite a sub-optimal implementation. The performance seems to be superior especially for complex environments, like the Humanoid-v4 environment.

The presented analysis suggests that the superior performance of C-DSAC may among others be attributed to a confidence-driven update of Q-values. It was shown that gradient weights of the critic loss function are small at state-action pairs with high variance. This results in a more conservative update for high-variance targets, reducing adaptation to potentially overestimated values. In this context, high variance indicates low confidence in the target, which improves only after repeated observations. As a consequence, C-DSAC learns Q-values in a two-phase process: First, building confidence in the target, then adapting to its expected value. While this confidence-building phase may imply slow and unstable learning, the numerical results indicate that the reduction in the effect of overestimation outweighs the additional cost. The C-DSAC framework and the current implementation remain open to various optimizations. Experiments with different initialization values can be carried out, particularly for the critic’s variance head. For a fair comparison with SAC baselines, automatic temperature adjustment was omitted, but could be integrated to enhance performance. Additionally, an exploration strategy based on the value distribution can be developed. C-DSAC models the return with a fixed-moment distribution. Efforts can be made to efficiently apply GMMs with complex neural network approximators to allow for a more expressive return distribution. Another promising direction for research is the utilization of task and reward decomposition based on the moments of the value distribution.



## Model-Enhanced Reinforcement Learning for Robotic Manipulation with Inverse Kinematics

This chapter explores the application of RL algorithms in the field of robotics and presents a novel approach designed specifically for robotic manipulation tasks. The proposed RL algorithm, termed "Reinforcement Learning - Inverse Kinematics" (RL-IK) [93], integrates the knowledge of a robot's kinematics and can be interpreted as a model-based RL algorithm. This allows for more efficient learning by leveraging the underlying physical model of the robot, particularly because movement of the robot and planning the task are separated and assigned to different modules. The focus of this research is on multi-joint robots in  $3D$  space, controlled by policies that are learned within the actor-critic (AC) framework, see Section 2.1.2.4. The algorithm aims to enhance the performance of robots in tasks involving manipulation and positioning by combining inverse kinematics with reinforcement learning for planning and decision-making at a higher abstraction level. According to ISO 8373, a robot is defined as an autonomous, programmed actuator ensemble capable of locomotion, manipulation, or tasks involving positioning. An industrial robot, specifically, is described as an autonomous, reprogrammable, multi-purpose manipulator with three or more programmable axes, which may be stationary or mounted on a mobile platform [94]. In this context, RL-IK is particularly suited for industrial robots performing complex manipulation tasks. Although the experiments in this chapter are conducted using the Franka Emika Panda robot [95], Figure 4.1, RL-IK is generalizable and can be applied to a wide range of robotic platforms, making it a versatile tool for various robotic manipulation applications.

## 4. MODEL-ENHANCED REINFORCEMENT LEARNING FOR ROBOTIC MANIPULATION WITH INVERSE KINEMATICS

---



**Figure 4.1:** Franka Emika Panda robot with a gripper end-effector

### 4.1 Forward Kinematics

Robotic systems consist of rigid links connected by joints, forming a kinematic chain that defines the robot’s geometry and degrees of freedom. Each link and joint is associated with a reference coordinate frame to ensure consistency in the description of joint configurations and the kinematic state. Forward Kinematics focuses on determining the position and orientation of the end-effector – the last link of the robot used for manipulation and other tasks – as a function of the robot’s joint configuration. Examples of end-effector tasks include gripping, tool operation and interaction with the environment. The end-effector can move freely, thus rendering the kinematic chain as open. The robot’s configuration in joint space can be described as a vector of  $n$  joint parameters:

$$\mathbf{v}_t = [\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n],$$

where  $\mathbf{v}_i$  represents joint-specific parameters such as angles for revolute joints or displacement for prismatic joints. In robotics, the goal is to find a sequence

$$\{\mathbf{v}_{t_1}, \dots, \mathbf{v}_{t_T}\}$$

to complete a given task without violating dynamic and environmental constraints. To accurately describe the influence of joint parameters on the end-effector, a consistent application of transformations between coordinate frames, starting from the base link to the end-effector, is required. A widely adopted convention in mechatronics is the Denavit-Hartenberg (DH) convention, which standardizes the assignment of coordinate frames along the kinematic chain. A transformation between frames consists of rotational and translational components. Rotations are represented using a  $3 \times 3$  rotational matrix  $\mathbf{R}$ , such that the orthogonality property  $\mathbf{R}\mathbf{R}^T = \mathbf{I}$  holds. Further, the

determinant is constrained to be 1 ensuring preservation of link lengths and reference frame orientation. Rotations are defined around the  $z$ -axis (yaw). Displacements are represented by a  $3 \times 1$  translation vector  $\mathbf{d}$ .

Under the Denavit-Hartenberg (DH) convention, illustrated in Figure 4.2, the assignment of coordinate frames follows a standardized methodology to describe transformations between links. The  $z$ -axis is aligned with the axis of rotation or translation for each joint, while the  $x$ -axis for the next frame is defined as the common normal between the  $z$ -axes, with the new origin positioned at the intersection between the subsequent  $z$ -axis and the common normal. The transformations between coordinate frames are characterized by four parameters. The angle  $\vartheta_{k-1}$  specifies the rotation required to align the base frame  $x$ -axis with that of the next frame. The parameter  $r_{k-1,k}$  represents the length of the common normal between the  $z$ -axes of consecutive joints. The parameter  $d_k$  denotes the distance along the previous  $z$ -axis to the intersection with the new  $x$ -axis. Finally, the angle  $\alpha_{k-1,k}$  defines the rotation around the new  $x$ -axis necessary to align the two  $z$ -axes. With these four parameters, the frame transformations can be described in terms of screw displacements consisting of a translation **TRA** and rotation function **ROT** around the  $z$ - and  $x$ -axis. The transformation matrix from frame  $k-1$  to frame  $k$  is given by

$${}^{k-1}\mathbf{T}_k(d_{k-1}, \vartheta_{k-1}, r_{k-1,k}, \alpha_{k-1,k}) = \mathbf{Z}_{k-1}(d_{k-1}, \vartheta_{k-1}) \cdot \mathbf{X}_k(r_{k-1,k}, \alpha_{k-1,k}), \quad (4.1)$$

where

$$\mathbf{Z}_{k-1}(d_{k-1}, \vartheta_{k-1}) = \mathbf{TRA}_{\mathbf{Z}_{k-1}}(d_{k-1}) \cdot \mathbf{ROT}_{\mathbf{Z}_{k-1}}(\vartheta_{k-1})$$

and

$$\mathbf{X}_k(r_{k-1,k}, \alpha_{k-1,k}) = \mathbf{TRA}_{\mathbf{X}_k}(r_{k-1,k}) \mathbf{ROT}_{\mathbf{X}_k}(\alpha_{k-1,k}).$$

The matrices associated with translation and (counterclockwise) rotation on these axes are explicitly defined by

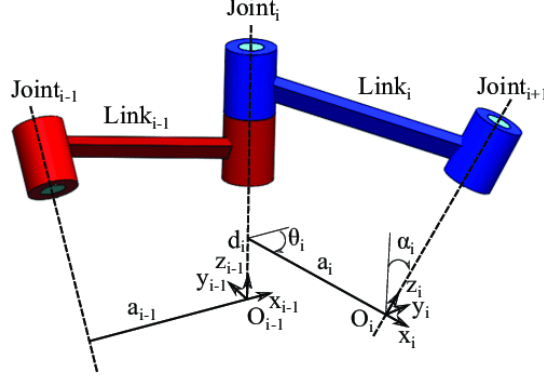
$$\mathbf{TRA}_{\mathbf{Z}_{k-1}}(d_{k-1}) = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & d_{k-1} \\ 0 & 0 & 0 & 1 \end{bmatrix}, \quad \mathbf{ROT}_{\mathbf{Z}_{k-1}}(\vartheta_{k-1}) = \begin{bmatrix} \cos(\vartheta_{k-1}) & -\sin(\vartheta_{k-1}) & 0 & 0 \\ \sin(\vartheta_{k-1}) & \cos(\vartheta_{k-1}) & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}$$

and

$$\mathbf{TRA}_{\mathbf{X}_k}(r_{k-1,k}) = \begin{bmatrix} 1 & 0 & 0 & r_{k-1,k} \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}, \quad \mathbf{ROT}_{\mathbf{X}_k}(\alpha_{k-1,k}) = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & \cos(\alpha_{k-1,k}) & -\sin(\alpha_{k-1,k}) & 0 \\ 0 & \sin(\alpha_{k-1,k}) & \cos(\alpha_{k-1,k}) & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}$$

#### 4. MODEL-ENHANCED REINFORCEMENT LEARNING FOR ROBOTIC MANIPULATION WITH INVERSE KINEMATICS

---



**Figure 4.2:** Illustration of the Denavit-Hartenberg convention [96]. Reproduced under CC BY 4.0 license (<https://creativecommons.org/licenses/by/4.0/>).

Performing the sequential matrix operations between two links in the kinematic chain as described in (4.1) yields

$${}^{k-1}\mathbf{T}_k = \begin{bmatrix} \cos(\vartheta_{k-1}) & -\sin(\vartheta_{k-1}) \cos(\alpha_{k-1,k}) & \sin(\vartheta_{k-1}) \sin(\alpha_{k-1,k}) & r_{k-1,k} \cos(\vartheta_{k-1}) \\ \sin(\vartheta_{k-1}) & \cos(\vartheta_{k-1}) \cos(\alpha_{k-1,k}) & -\cos(\vartheta_{k-1}) \sin(\alpha_{k-1,k}) & r_{k-1,k} \sin(\vartheta_{k-1}) \\ 0 & \sin(\alpha_{k-1,k}) & \cos(\alpha_{k-1,k}) & d_{k-1} \\ 0 & 0 & 0 & 1 \end{bmatrix}. \quad (4.2)$$

Note that the function arguments in Equation (4.2) were omitted to prevent convoluted notation. By sequentially applying these transformations across all links in the chain, the end-effector's position and orientation relative to the base frame can be calculated as a function of the angular displacement  $\boldsymbol{\vartheta}$  and linear displacement  $\mathbf{d}$ . Note that  $\mathbf{r}$  and  $\boldsymbol{\alpha}$  are constants derived from the robot's geometry and are therefore not explicitly written as function arguments. The overall transformation is then given by

$$\mathbf{T}_{EE}(\mathbf{d}, \boldsymbol{\vartheta}) = {}^1\mathbf{T}_n(\mathbf{d}, \boldsymbol{\vartheta}) = {}^1\mathbf{T}_2(d_1, \vartheta_1) {}^2\mathbf{T}_3(d_2, \vartheta_2) \dots {}^{n-1}\mathbf{T}_n(d_{n-1}, \vartheta_{n-1}).$$

For the rest of the work, it is assumed that the kinematic chain does not include prismatic joints, i.e.  $d_{k-1}$  is a constant. The end-effector's state can be expressed as

$$\mathbf{f}(\boldsymbol{\vartheta}) = {}^1\mathbf{T}_n(\boldsymbol{\vartheta}). \quad (4.3)$$

## 4.2 Inverse Kinematics

Inverse kinematics (IK) is the inverse operation of forward kinematics, aiming to determine a joint configuration  $\boldsymbol{\vartheta}^*$  that positions and orients the end-effector at a desired target  $\mathbf{p}^*$  while respecting physical constraints. For rigid end-effectors, analytical solutions exist only for 6-DoF manipulators. When the robot deviates from this configuration, the resulting system of equations might become under- or overdetermined. In overdetermined systems, additional constraints (e.g., energy efficiency, collision avoidance) must be considered, which often complicates finding a solution. Singularities – configurations where the robot loses a degree of freedom – must be exempt from the set of solutions. Jacobian-based numerical methods are widely used to address these challenges. These methods leverage the Jacobian matrix, which relates joint velocities to end-effector velocities. The roll  $\eta$ , pitch  $\xi$ , and yaw  $\psi$  can be extracted from the transformation matrix of the end-effector in (4.3):

$$\begin{aligned}\eta &= \text{atan2}(T_{EE32}, T_{EE33}) \\ \xi &= \text{atan2}\left(-\frac{T_{EE31}}{\sqrt{T_{EE32}^2 + T_{EE33}^2}}\right) \\ \psi &= \text{atan2}(T_{EE21}, T_{EE11}),\end{aligned}\tag{4.4}$$

where  $M_{ij}$  indicates the  $i$ -th row in column  $j$  of matrix  $\mathbf{M}$  and  $\text{atan2}$  is the piecewise defined inverse tan-function that is well-defined everywhere except at  $\text{atan2}(0, 0)$ . Then, the end-effector's position and orientation can be compactly expressed as

$$\mathbf{p}(\boldsymbol{\vartheta}) = \begin{bmatrix} x & y & z & \eta & \xi & \psi \end{bmatrix},$$

where the elements depend on  $\mathbf{f}(\boldsymbol{\vartheta})$  according to Equation (4.4). This representation simplifies the subsequent IK equations. The end-effector's rate of change as a function of the joint configuration can be expressed by the Jacobian

$$\mathbf{J}_p := \frac{\partial \mathbf{p}(\boldsymbol{\vartheta})}{\partial \boldsymbol{\vartheta}}.$$

The joint configuration is iteratively updated by

$$\boldsymbol{\vartheta}_{\ell+1} = \boldsymbol{\vartheta}_{\ell} + \Delta \boldsymbol{\vartheta}.\tag{4.5}$$

Using a modified Gauss-Newton approach and utilizing the Moore-Penrose pseudo-inverse of the Jacobian  $\mathbf{J}_p^+ = (\mathbf{J}_p^T \mathbf{J}_p)^{-1} \mathbf{J}_p^T$ , the update can be approximated as

$$\Delta \boldsymbol{\vartheta} \approx \mathbf{J}_p^+(\mathbf{p}^* - \mathbf{p}(\boldsymbol{\vartheta})).$$

## 4. MODEL-ENHANCED REINFORCEMENT LEARNING FOR ROBOTIC MANIPULATION WITH INVERSE KINEMATICS

---

The joint configuration is updated until the error defined by the difference between the current state and target state  $\mathbf{p}^*$  satisfies the convergence criterion  $|\mathbf{p}^* - \mathbf{p}| \leq \epsilon$ . This update arises from a linearization around  $\mathbf{p}(\boldsymbol{\theta})$ . It ensures minimal change in the joint configuration in the direction of a (local) minimum. The pseudo-inverse of the Jacobian is employed to address its potential non-invertibility, which arises due to the Jacobian often being non-square in robotic scenarios. However, this approach results in large joint updates in or near singular configurations, where the Jacobian loses rank. To mitigate numerical instabilities, the Levenberg method replaces the pseudo-inverse with a "dampened" inverse [97]:

$$\Delta\boldsymbol{\theta} = (\mathbf{J}_p^T \mathbf{J}_p + \nu \mathbf{I})^{-1} \mathbf{J}_p^T (\mathbf{p}^* - \mathbf{p}(\boldsymbol{\theta})), \quad (4.6)$$

where  $\mathbf{I}$  is the identity matrix and  $\nu \geq 0$  is the damping factor. Note that Equation (4.6) reduces to the modified Gauss-Newton if  $\nu = 0$ , whereas it approaches a gradient descent step for large  $\nu$ ,

$$\Delta\boldsymbol{\theta} \approx \frac{1}{\nu} \mathbf{J}_p^T (\mathbf{p}^* - \mathbf{p}(\boldsymbol{\theta})) \quad \text{if } \nu \gg 0.$$

The identity matrix induces an update that affects all parameter directions equally, meaning that differences in parameter contributions to the error is disregarded. To respect the error surface curvature, the identity matrix is replaced by the Gram matrix of the Jacobian, yielding the Levenberg-Marquardt update step [98]:

$$\Delta\boldsymbol{\theta} = (\mathbf{J}_p^T \mathbf{J}_p + \nu \text{diag}(\mathbf{J}_p^T \mathbf{J}_p))^{-1} \mathbf{J}_p^T (\mathbf{p}^* - \mathbf{p}(\boldsymbol{\theta})) \quad (4.7)$$

In the Levenberg-Marquardt algorithm (LMA), the current joint configuration is updated according to Equation (4.7) until the error between  $\mathbf{p}^*$  and  $\mathbf{p}$  is sufficiently small. LMA balances the effect of Gauss-Newton and gradient descent updates, with the former offering precision and the latter providing stability.

### 4.3 Separation of Planning and Execution of Movement: RL-IK

Motion planning is an abstraction of end-effector movements, where the manipulations required to achieve a task are represented as a sequence of end-effector states over discrete time steps. In the context of robot-specific reinforcement learning (RL), it is natural to distinguish between motion planning and the physical execution of the robot's movements. This separation allows for a more modular and interpretable control

### 4.3 Separation of Planning and Execution of Movement: RL-IK

---

architecture. "Reinforcement Learning - Inverse Kinematics" (RL-IK) operates on the basis of this task separation explicitly: At each decision step, the RL agent generates a target position and orientation for the end-effector based on its policy:

$$\mathbf{a}_t = \mathbf{p}^* \sim \pi(s_t).$$

The low-level controller, implemented using the Levenberg-Marquardt Algorithm (LMA), then assumes the responsibility of achieving the desired position and pose. No additional high-level actions are inferred from the RL policy until the end-effector reaches the target configuration. While the inverse kinematics phase is active, constraint violations – such as collisions, excessive force application, or breaches of time limits – must still be monitored. If such a violation occurs, the "IK phase" is prematurely terminated, the environment resets, and an appropriate reward is emitted to guide future policies. A detailed description of the RL-IK algorithm is described in Algorithm 4.1.

The adoption of LMA as a low-level controller, instead of allowing the RL agent to directly control individual joints, offers several advantages:

1. Reduction in action space: Regardless of the robots DoF, the action space will always consist of a fixed number of actions that describe the end-effector's position and orientation. This reduces the dimensionality of the RL problem and shifts computational complexity from the sample-inefficient RL component (see Section 2.1.3.3) to the computationally efficient Levenberg-Marquardt optimization unit [99].
2. Experimental simulations indicate that robot movements exhibit smoother trajectories with RL-IK. This smoothness can likely be attributed to joint updates that align with consistent spatial directions dictated by the LMA solver.
3. While LMA is generally robust to rank loss, a poorly chosen damping factor  $\nu$  can still result in singular configurations. In such cases, the agent may struggle to achieve its objectives. However, since the agent is designed to maximize a task-specific reward, it implicitly learns to avoid configurations that hinder efficient movement and manipulation capabilities.
4. Instead of evaluating and improving the policy at every intermediate step, it suffices to perform these operations after the end-effector transitions from its current pose  $\mathbf{p}$  to the target pose  $\mathbf{p}^*$ . This design generates fewer but qualitatively superior experience points. While this approach could theoretically exacerbate the credit assignment problem, experimental results in Section 4.5 indicate significant faster convergence to optimal rewards and improved training stability compared to direct joint-level control methods.

## 4. MODEL-ENHANCED REINFORCEMENT LEARNING FOR ROBOTIC MANIPULATION WITH INVERSE KINEMATICS

---

---

**Algorithm 4.1** RL-IK

---

```
1: Initialize  $\vartheta$ 
2:  $\mathcal{B} \leftarrow \emptyset$ 
3:  $\mathbf{p}_t \leftarrow \mathbf{p}(\vartheta)$  # from forward kinematics
4: repeat
5:   Receive initial observation state  $s_1$ 
6:   for each environment step do
7:      $\mathbf{p}_t^* \sim \pi(s_t)$ 
8:     repeat
9:        $\Delta\vartheta$  from Equation(4.7) using  $\mathbf{p}_t^*$  and  $\mathbf{p}_t$ 
10:       $\vartheta \leftarrow \vartheta + \Delta\vartheta$ 
11:       $\mathbf{p}_t \leftarrow \mathbf{p}(\vartheta)$  # from forward kinematics
12:      if constraint violation then
13:        Break
14:      until  $|p - p^*| \leq \epsilon$ 
15:      Observe transition  $s_{t+1} \sim P(s_t, a_t)$ 
16:      Observe reward  $r_t$ 
17:       $\mathcal{B} \leftarrow \mathcal{B} \cup \{(s_t, \mathbf{p}_t^*, r_t, s_{t+1})\}$ 
18:   Perform gradient steps
19: until Stopping criterion
20: return models
```

---

### 4.4 Reward Engineering for RL-IK

In RL-IK, designing an effective reward function is a critical factor for achieving successful motion planning and control. The reward function must carefully balance multiple objectives, including event-free motion trajectories (where no constraint violations occur), temporal efficiency, control cost, and partial task fulfilment. An improperly balanced reward structure can significantly hinder convergence or even result in undesirable "suicidal" behaviour, as discussed in Section 2.1.4. The reward function in RL-IK is composed of the following sub-rewards:

1. **Costs:** These include penalties for excessive time use and high control effort. Both must be minimized to ensure efficiency and smooth operation.
2. **Critical Rewards:** Collisions must be heavily penalized, but the temporal context is crucial. Collisions occurring earlier in the trajectory should be penalized more severely than those occurring later, as earlier collisions often indicate more

fundamental planning errors. However, severely punishing later collisions will hinder exploration and encourage the agent to terminate episodes early.

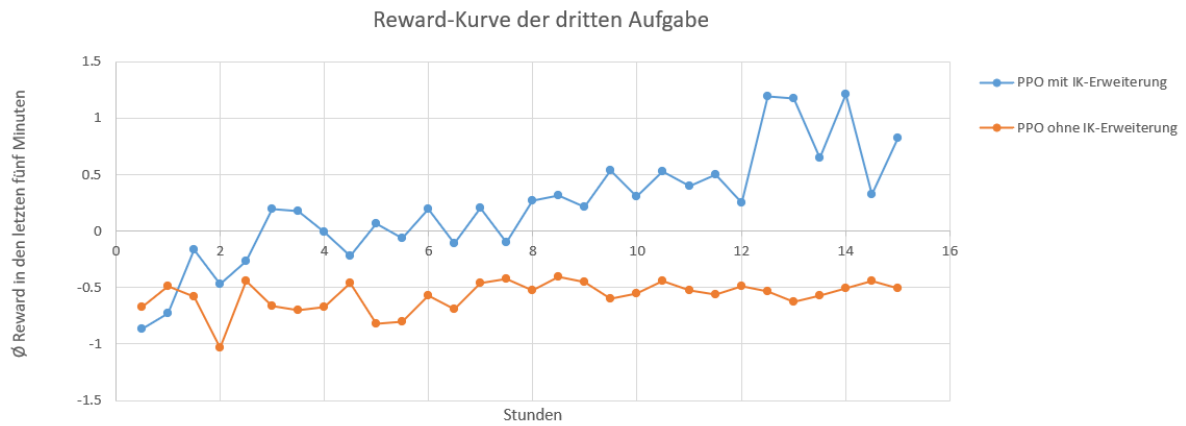
3. **Positive Reward:** Successfully gripping the target object and moving it towards the target location provides a positive reinforcement signal, incentivizing goal achievement.
4. **Ambivalent Rewards:** Factors such as the distance to the target object and alignment with the target position can be interpreted both positively and negatively, depending on the specific stage of the task and the overarching objective.

In a real-world scenario, image-based distance approximators and tactile registrators might be utilized to obtain spatial proximity data and detect collisions. This nuanced reward design ensures that the RL-IK can optimize both high-level motion planning and low-level control execution, ensuring to complete the robotic task efficiently.

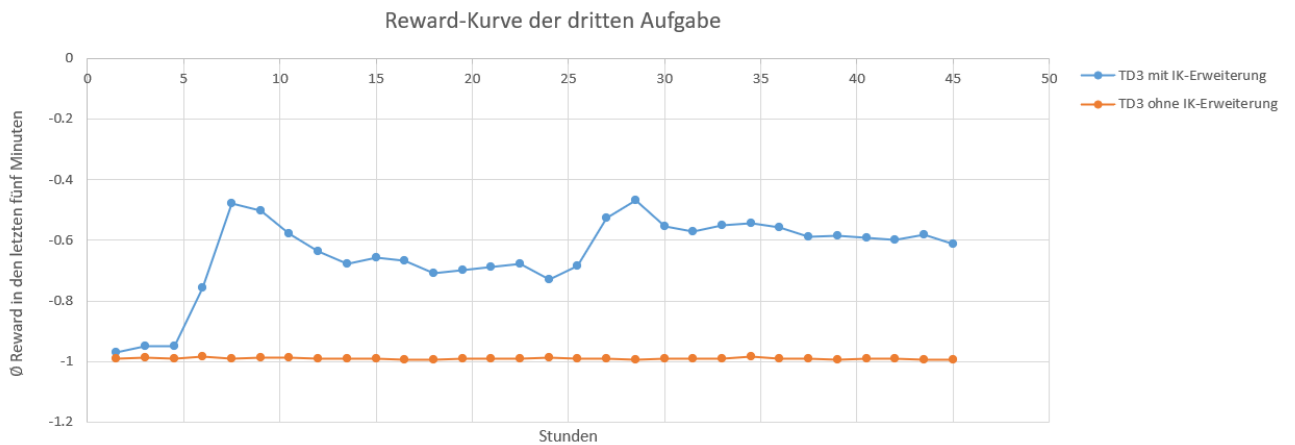
## 4.5 Experiments

The performance of RL-IK was assessed in a simulated robotic task illustrated in Figure 2.1. Successful task completion meant that the robot picked up a target object, transported it over a wall and placed it in a tray without colliding with the environment. PPO (Section 2.1.6.4) and a modified version of DDPG (Section 2.1.6.5) known as TD3 were employed to solve the task. IK-augmented variants of these algorithms, PPO-IK and TD3-IK, were then used in the same environment. Figure 4.3 shows the performances of PPO and PPO-IK over five runs. Because each step in the IK variant is a planned sub-trajectory rather than an incremental movement of the end-effector, PPO-IK needs considerably fewer "RL iterations" than PPO to converge to a task-solving reward. To allow comparison, the performance is reported as a function of wall-clock time in hours. All experiments were run on the same system and under the same circumstances. Figure 4.4 shows the performance for a similar set of experiments conducted with the TD3 algorithm. Although both variants do not manage to solve the task in a feasible time, TD3-IK converges to a significantly higher suboptimal reward. Because TD3 is highly hyperparameter sensitive, it is hypothesized that the parameters were not correctly configured for the task.

## 4. MODEL-ENHANCED REINFORCEMENT LEARNING FOR ROBOTIC MANIPULATION WITH INVERSE KINEMATICS



**Figure 4.3:** PPO used in the robotic task with and without inverse kinematics; graph taken from [93]. The difference in step count between the base and IK-augmented variant is drastic in favour of PPO-*IK*, therefore the graph illustrates the agents performance in hours.



**Figure 4.4:** TD3 used in the robotic task described in Figure 2.1 with and without inverse kinematics; graph taken from [93]. The difference in step count between the base and *IK*-augmented variant is drastic in favour of TD3-*IK*, therefore the graph illustrates the agents performance in hours.

## Introduction to Supervised Learning and Ensemble Strategies

As shown in Figure 1.1, machine learning encompasses a diverse range of paradigms aimed at enabling data-driven generation of models. While the primary focus of this thesis has been on reinforcement learning due to its applicability in developing advanced models for control, it is essential to acknowledge the complementary role of supervised learning (SL) in addressing specific challenges, particularly in machine perception and state identification. Supervised learning operates on labelled datasets, where the model learns to map input data to corresponding outputs based on provided examples. Through an iterative process, the model minimizes the discrepancy between predicted and actual outputs. In the SL paradigm, models adjust their parameters by comparing predicted outputs with ground-truth labels, reinforcing accurate predictions and correcting deviations. This chapter provides an introduction to the fundamental principles of SL and delves deeper into a methodology often employed in SL: Ensemble learning.

### 5.1 Supervised Learning

Supervised Learning is a foundational paradigm in machine learning, wherein models learn from labelled data (interpreted as supervisory signal, usually annotated by humans) to make predictions. Central to supervised learning is constructing a predictive model that can generalise well to unseen instances by learning patterns from labelled examples. In supervised learning, the overarching objective is to generate a statistical model reproducing a mapping that accurately predicts output labels from input features. A general SL task is to learn a function  $f : \mathbb{R}^\ell \rightarrow \mathbb{R}^m$  to approximate an unknown

## 5. INTRODUCTION TO SUPERVISED LEARNING AND ENSEMBLE STRATEGIES

---

mapping  $\Psi : \mathcal{X} \rightarrow \mathcal{Y}$  between the input space  $\mathcal{X}$  and output space  $\mathcal{Y}$ . Typically, the sets  $\mathcal{X}$  and  $\mathcal{Y}$  are not known and we only have a finite set of matching observations  $\{(\mathbf{x}_j, \mathbf{y}_j)\}_{j \in J} \subset \mathcal{X} \times \mathcal{Y}$  and  $J$  is its index set. The finite set of observations on which the model is trained is referred to as the training set. It must be stressed that in supervised learning, the potentially noisy labels  $\hat{\mathbf{y}} = \mathbf{y} + \boldsymbol{\epsilon}$ , are considered to be reasonably accurate  $\frac{1}{J} \sum_{j \in J} |\mathbf{y}_j - \hat{\mathbf{y}}_j| \approx 0$ . Neural network-based models play a pivotal role in approximating high-generalization-capable functions, i.e. functions that can perform well for data outside of the training set. With a hindsight to Section 1.2.3, gradients calculated from a given loss defined by the labels and the model's predictions, iteratively adjust network parameters. The choice of the neural network architecture often depends on the complexity of the problem and available computational capabilities. Image-based SL often require CNNs for efficient learning whereas transformer networks are well suited for problems in which long-term dependencies of the data dominate [12].

### 5.2 Ensemble Learning

Inspired by the "Wisdom of Crowds" principle [100], ensemble learning seeks to form a function  $\hat{f}$  by aggregating multiple weaker and diverse base learners (BLs) [101] to manage the bias-variance trade-off. When neural networks are employed as base learners, the approach is termed Deep Ensemble Learning (DEL) [102].

To formally define the convex ensemble model, consider a finite set of  $n$  base learners given by  $\mathcal{H} = \{h_i : \mathbb{R}^\ell \rightarrow \mathbb{R} \mid i \in [n]\}$ . The ensemble model prediction is given by

$$\hat{f}(\mathbf{x}) = \sum_{i \in [n]} \alpha_i h_i(\mathbf{x}), \quad (5.1)$$

where  $\alpha_i$  are weights of the individual BLs. In ensemble learning, the weights  $\alpha$  are usually restricted to the standard simplex:

$$\Delta_n = \left\{ \alpha \in \mathbb{R}^n \mid \sum_{i=1}^n \alpha_i = 1, \alpha_i \geq 0, i \in [n] \right\}. \quad (5.2)$$

An important property and motivation of ensemble learning is that the final ensemble model is not restricted to the model space of the base learners, *i.e.*, the combination of base learners spans a larger model space. Thus, the base learners can be restricted to a far simpler model space than the final ensemble model. This property is further described in the following proposition.

**Proposition 5.1.** Given a set of base learners  $h_1, \dots, h_n \in \mathcal{F}$ , the ensemble model  $\hat{f} := \sum_{i=1}^n \alpha_i h_i$  lies within the convex hull of  $\mathcal{F}$ , but it does not necessarily belong to the model space  $\mathcal{F}$ .

Proof. Let  $\mathcal{F}$  be defined as the model space of neural networks with a single ReLU node, *i.e.*, piecewise linear functions with one constant and one affine part. Consider a convex combination  $f$  of two base learners from  $\mathcal{F}$ . Then  $f$  is an ensemble model that is a piecewise linear function with one constant part and two affine parts and hence not an element of  $\mathcal{F}$ .  $\square$

### 5.2.1 Bootstrap Aggregating

Bootstrap aggregating, commonly referred to as bagging, is a widely used variance reduction technique in ensemble learning. It leverages multiple base learners trained on different subsets of the original training dataset to enhance predictive performance and robustness. By mitigating variance, bagging is particularly effective when applied to well-performing base learners or scenarios where data is sparse in specific regions of the training data manifold. The subsets used to train the base learners are generated through bootstrapping, a process that involves sampling with replacement from the original training set. These subsets maintain the same size as the original dataset but typically contain duplicate samples. As a result, each subset represents approximately 63.2% of the unique training data, as shown in [103]. A predefined number of  $n$  base learners are then independently trained on these bootstrapped subsets. The final ensemble prediction is obtained by averaging their outputs, assigning an equal weight to each learner: The weight of  $h_i(\mathbf{x})$  in Equation (5.1) becomes  $\alpha_i = \frac{1}{n} \quad \forall i \in [n]$ . Bagging can be computationally expensive, as the base learners must be reasonably performative, requiring relative complex base learners. In summary, bagging is a powerful technique for reducing variance and improving the generalization capabilities of predictive models, particularly in cases where base learners are sufficiently performant and computational constraints are manageable.

### 5.2.2 Boosting

Boosting is an ensemble learning technique designed to reduce bias by sequentially training base learners. The key idea is that each subsequent model "boosts" the ensemble's overall performance by focusing on and correcting the mistakes made by its predecessors. In each iteration, the training data is re-weighted to reflect the model's

## 5. INTRODUCTION TO SUPERVISED LEARNING AND ENSEMBLE STRATEGIES

---

accuracy on individual samples. Misclassified or poorly predicted examples are assigned higher weights, incentivizing subsequent models to focus on these challenging data points. This iterative re-weighting mechanism enables boosting algorithms to progressively refine their predictive performance. Once a sufficient number of base learners have been trained, their predictions are aggregated using Equation (5.1). Base learners that demonstrate higher accuracy, particularly on difficult examples, are assigned greater model weights  $\alpha_i$ . The base learners employed in boosting are generally weak and computationally simple. Some boosting approaches use decision stumps as base learners. Prominent examples of boosting algorithms include AdaBoost [104, 105], Gradient Boosting Machines [106] and Integer Programming Boost [107]. Prior to neural network-based machine learning, these algorithms have been used across diverse domains due to their ability to generate a complex predictor from simple base learners.

### 5.2.3 Voting

Voting in ensemble learning refers to an approach in which independently trained base learners collaboratively contribute to the final prediction. This methodology is predominantly employed in classification tasks, where each base learner casts a "vote" for a specific class label. Since base learners are trained on the same training data set, the effectiveness of the voting mechanism hinges on the diversity of these learners. Diversity in prediction can be achieved through techniques like random initialization in stochastic models, see Section 2.1.5.4. Without sufficient diversity, the ensemble risks become biased towards shared weaknesses among the base learners. Two principal voting schemes are commonly used: hard voting and soft voting. In the hard voting scheme, each base learner outputs a discrete class prediction, and the final ensemble prediction corresponds to the majority class label among the votes. This approach assumes that each learner contributes equally to the final decision. In contrast, the soft voting scheme aggregates the predicted probability distributions from all base learners. The class with the highest aggregated probability is selected as the final prediction. Soft voting often leverages the confidence levels of individual learners, leading to more nuanced and potentially more accurate ensemble predictions.

## 5.3 Relevance for Robotics

Perception is a cornerstone of effective robotic control, enabling robots to interpret and respond to highly dynamic environments. Perception involves reliably identifying the robot's state space (see Chapter 2), including detecting target objects, obstacles,

agent states, poses, and potential collisions. This capability is essential for decision-making processes and fundamental to control-generating algorithms described in earlier chapters.

Robotic perception relies on computer vision and "artificial tactile systems". Computer vision allows robots to process and understand visual data. Vision modules can be independently trained with on annotated datasets to identify key elements of the environment, such as targets, obstacles, or paths, and even perform tasks like depth estimation from standard visual sensors. Similarly, tactile systems enable robots to sense and interpret physical interactions in the environment. Supervised learning approaches can help these systems differentiate between desirable physical contact, such as pushing a swing door, and unintended collisions. It will be shown in Chapter 6 that ensemble models are generally more performative than their simpler counterparts. This improvement, however, comes with a significant computational overhead, as ensembles are complex structures. In real-world robotic applications, state-space identification and environmental interpretation must occur in real-time to ensure short response time. Balancing the accuracy benefits of ensemble learning with the computational constraints of real-time inference remains a critical challenge in deploying these models effectively on robotic hardware.

Beyond perception, supervised learning can play a role in robotic control, particularly through imitation learning, commonly referred to as behavioural cloning [108]. In this paradigm, the control problem is initially framed as a supervised learning problem, where the robot learns to imitate expert demonstrations. This approach provides a "warm start" for reinforcement learning algorithms, allowing them to begin training from an informed initial policy rather than random exploration.

The integration of supervised and ensemble learning strategies into robotic perception and control pipelines represents a step toward more autonomous and reliable systems. Vision and tactile systems, trained with supervised learning, can deliver robust sensory information representing the observation space in reinforcement learning algorithms. The ongoing challenge lies in bridging the gap between computational complexity and real-time deployment to enable robot operations in highly dynamic environments.



## Optimizing Ensemble Models via Column Generation

This chapter presents a novel boosting algorithm that integrates principles from column generation, a well-established technique in classical optimization, with neural-network-based supervised learning. The proposed approach, named data-reduction ensemble learning algorithm (DDA), optimizes an ensemble of models by solving a linear programming (LP) problem defined by the predictions on the training dataset.

The algorithm constructs an optimal convex combination from a set of base learners, leveraging the concept of a soft margin commonly employed in support vector machines. Through this optimization, the ensemble aims to maximize predictive performance on the training data while establishing and maintaining generalization capabilities. The set of base learners is iteratively extended. The generation of new base learners is guided by the ensemble's performance on specific training data points. Neural network approximators and supervised learning techniques are employed to fit new and diverse base learners to correct the ensemble's previous mistakes. This work attempts to bridge a methodological gap between traditional optimization techniques and modern neural network-based learning frameworks.

### 6.1 Introduction

In recent years, there has been a growing interest among researchers in bridging the fields of mathematical optimization and machine learning by applying optimization methodologies to machine learning problems. A particularly promising avenue of exploration lies in ensemble learning—a widely adopted approach in the machine learning community at the time of writing [109, 102, 110, 111]. Motivated by the structural parallels between ensemble learning and Column Generation (CG), the authors have

## 6. OPTIMIZING ENSEMBLE MODELS VIA COLUMN GENERATION

---

recently focused on investigating ensemble methods in the context of "classical" optimization. Ensemble learning methods combine multiple models to enhance predictive performance. There are plenty of real-world applications of ensemble learning such as object identification [112, 113, 114, 115], fault detection and diagnosis [116, 117, 118], semantic segmentation, edge detection [119, 120, 116], and for black box optimization [121].

It is known that deep learning networks may require vast computational effort to perform well in complex tasks. For example, state-of-the-art large language models can have more than 100 billion parameters [122]. From that point of view, combining easier or smaller learners may be an attractive approach. Ensemble learning shares similarities with the concept of Column Generation (CG), see [123, 111]. The latter computes solution estimates by combining partial solution candidates (options) in master optimization problems, which provide information to generate more options (columns). Efficient CG algorithms for solving nonlinear optimization problems, like in [124], are based on fast generation of initial columns in a starting phase and solving low-dimensional (pricing) sub-problems for generating columns. Ensemble models typically combine simpler base learners to obtain a more accurate and complex model. [125] show that the generalization error of an ensemble model strongly depends on the so-called diversity, see Section 6.2.4. Our research focuses on investigating whether base learners can be generated in a manner analogous to Column Generation (CG), utilizing subsets of the data rather than the entire training set. This approach aims to enhance both the diversity and generalizability of the resulting ensemble model while simultaneously reducing the computational overhead associated with processing the full training dataset for each new base learner. The main research question in this chapter is how to efficiently decompose the optimization problem of training an ensemble model for classification problems and understanding the effect of dynamic data selection. One of the goals of this work has been to develop a decomposition framework that splits the full optimization problem into smaller sub-problems with optimality and convergence guarantees. This chapter demonstrates that such a decomposition approach exists under these assumptions, attempting to improve understanding in using global optimization techniques for determining globally optimal ensemble models.

To show this, a specific ensemble learning algorithm is designed which improves an initial model of unspecified complexity by generating and combining new base learners into an ensemble. The base learners and the initial model are combined into an ensemble by solving a linear programming (LP) master problem, which maximizes the soft margin of the ensemble classifier. The resulting dual solutions of the master problem are further

used in sub-problems, which train new base learners to boost the previous ensemble. Several types of sub-problems are considered for training base learners in the algorithm. To illustrate the efficiency of the ensemble algorithm, we measure the diversity and the empirical generalization error on several test instances. The main contribution of this chapter is summarized as follows:

- A novel method is designed for (binary) classification tasks, characterized by the following key properties:
  - Enabling a dual decomposition that splits the classification problem into sub-problems considering small base learners and a sub-set of the data points. In this way, the learning task is decomposed into simpler sub-problems with fewer variables.
  - The method is a framework for decomposing certain classification problems, that is globally convergent if sub-problems are solved to global optimality.
  - The data sub-sets defining the sub-problems are computed using sparse dual solutions of the LPBoost master problem, see [123]. The choice of the number of data points in the set defines a trade-off between accuracy and computational complexity.
  - Sub-problems are learning problems with either linear or nonlinear loss.
  - The base learners generated by the method are diverse.
- We show the method can reliably improve an existing (ensemble) classifier.
- By analysing the validation data we found that the generalization error decreases as diversity increases, thus supporting the use of model diversity as a performance measure.
- The proposed algorithm computes an optimal ensemble model within the convex hull of a given model space. Our results indicate that the ensemble model has the potential to achieve accuracy surpassing that of individual base learners. Notably, the ensemble convex hull constitutes a richer model space, as the base learner model space forms a proper subset of it. Additionally, we present a general theoretical result characterizing the ensemble model space.

This chapter is structured as follows. Section 6.2 focuses on LPBoost-based ensemble learning and the mechanisms involved in the algorithm, namely data selection, diversity and ensemble learning. Section 6.3 describes the algorithm in a modular fashion; first it provides an overview over the general scheme and then describes the core features of the algorithm. Section 6.4 presents experiments with standard LPBoost and the new variant for two data sets with varying difficulties. Section 6.5 summarizes our findings.

## 6.2 Diverse Ensemble Learning and LPBoost

We provide background on the problem formulation, LPBoost, and some fundamental theory. We first sketch the problem formulation and the ensemble learning concept in Section 6.2.1. Then, we describe the LP master problem of LPBoost and the pricing problem in Section 6.2.2. We give an intuition for the potential of only considering a sub-set of data points in the generation of base learners by exploiting the sparsity in Section 6.2.3. Section 6.2.4 describes the role of model diversity based on considerations from literature.

### 6.2.1 Problem formulation and ensemble learning (EL)

In this chapter, a binary classification problem is considered with  $\mathcal{Y} = \{-1, 1\}$ . For one single arbitrary observation  $(x, y)$ ,  $y$  denotes the target of this observation with  $y \in \{-1, 1\}$ . The binary classifier  $f : \mathbb{R}^\ell \rightarrow \{-1, 1\}$  as  $f(x) = \text{sign}(\hat{f}(x))$  is defined by the score function  $\hat{f} : \mathbb{R}^\ell \rightarrow \mathbb{R}$ ,  $\hat{f} \in \mathcal{F}$ . Given a model space  $\mathcal{F}$ , the learning task aims at obtaining a predictive score model  $\hat{f}^* \in \mathcal{F}$  which minimizes the expectation of a loss function  $\mathcal{L}$

$$\hat{f}^* = \operatorname{argmin} \left\{ \mathbb{E}[\mathcal{L}(\hat{f}(x), y)] : \hat{f} \in \mathcal{F} \right\}, \quad (6.1)$$

where  $\mathbb{E}$  denotes the expectation regarding a joint probability distribution  $P(x, y)$  over  $\mathcal{X} \times \mathcal{Y}$ . As we are typically restricted to a finite training set  $T = \{(\mathbf{x}_j, \mathbf{y}_j)\}_{j \in J}$ , the learning task is usually approximated by

$$\hat{f}^* = \operatorname{argmin} \left\{ \sum_{j \in J} \mathcal{L}(\hat{f}(x_j), y_j) : \hat{f} \in \mathcal{F} \right\}. \quad (6.2)$$

Throughout the chapter, we assume that the model space  $\mathcal{F}$  is limited to a predefined ML architecture, *e.g.*, a ANN with a specific number of nodes and layers. In binary classification, a loss function could be for instance a linear classification error (6.18) or a non-linear binary cross entropy loss (6.17). Considering a training set  $T := \{(\mathbf{x}_j, y_j)\}_{j \in J}$ , the learning task in ensemble learning is two-fold: learn the ensemble weights  $\alpha$  and construct the set of base learners  $\mathcal{H}$ . Consider for instance deep neural networks with a given architecture. The base learners  $h_i$  are thus parameterized by the weights and biases of the ANN. To simplify the notation, we define

$$h_i(\mathbf{x}) := \varphi(\mathbf{x}; \mathbf{W}_i), \quad (6.3)$$

where  $\varphi$  denotes the input-output mapping of a ANN with weights and biases given by a parameter vector  $\mathbf{W}_i \in \mathbb{R}^k$ . We are assuming that all the weights and biases are

bounded, *i.e.*,  $\|\mathbf{W}_i\|_1 \leq M < \infty$ . The space of base learners is then reduced to the Hilbert space  $\mathcal{F} = \{\varphi(\mathbf{x}; \mathbf{W}) \mid \mathbf{W} \in [-M, M]^k\}$ . Furthermore, by limiting the number of BLs in the ensemble model to  $n$ , the ensemble learning problem can be written as

$$\begin{aligned} \min_{\alpha, h_1, \dots, h_n} \quad & \sum_{j \in J} \mathcal{L} \left( \sum_{i \in [n]} \alpha_i h_i(\mathbf{x}_j), y_j \right) \\ \text{s.t.} \quad & \alpha \in \Delta_n, \\ & h_i \in \mathcal{F}, \quad i \in [n], \end{aligned} \tag{6.4}$$

where  $\Delta_n$  denotes the standard  $n$ -dimensional simplex. The minimizer of (6.4) is an optimal ensemble model for the given training set  $T$ , loss function  $\mathcal{L}$ , base learner architecture, and number of base learners. By Proposition 2.1, we know that the ensemble model given by the minimizer of problem (6.4) is within the convex hull of  $\mathcal{F}$ , denoted  $\text{conv}(\mathcal{F})$ . With an additional assumption on the number of base learners we get a more general optimality property for the ensemble model given by problem (6.4), as described in the following proposition.

**Proposition 6.1.** Given that the number of base learners in problem (6.4) is sufficiently large, then the resulting ensemble model will also be optimal within  $\text{conv}(\mathcal{F})$  for the given objective function.

Proof. Follows from Carathéodory's theorem [126]. □

Propositions 5.1 and 6.1 are significant because they establish that the ensemble model, derived by solving (6.4), is optimal within a broader and more complex model space compared to the space of the base learners. This insight serves as a primary motivation for designing an algorithm that iteratively optimizes over the comparatively simpler base learners to effectively address (6.4) and construct a robust ensemble model.

Each base learner  $h_i$  in problem (6.4) is parameterized by  $k$  parameters. The complete problem contains  $k \times n + n$  variables. In practice, however, a substantial number of auxiliary variables may be required to represent the activation states of the base learners for each data point, significantly increasing the dimensionality of the optimization problem. Even when disregarding this large variable count, problem (6.4) remains inherently challenging due to the non-convexity, as formally described in the following proposition.

**Proposition 6.2.** Unless  $\mathcal{F}$  is a singleton set, problem (6.4) is in general non-convex.

## 6. OPTIMIZING ENSEMBLE MODELS VIA COLUMN GENERATION

Proof. With  $\mathcal{F}$  not restricted to a singleton set, the terms  $\alpha_i h_i(\mathbf{x}_j)$  in the objective results in bilinear terms as function  $h_i$  is parameterized by weights and biases.  $\square$

Considering that non-trivial classification problems typically require DNNs with thousands to millions of parameters, problem (6.4) still remains well out of reach of today’s general purpose deterministic global optimization solvers. To solve problem (6.4), we present an algorithm inspired by the LPBoost algorithm [123] that follows a column generation approach. By this decomposition approach, we are able to determine one base learner, *i.e.*, column, at a time and thus limit each sub-problem to a smaller sub-set of variables. The presented algorithm also allows to only consider a sub-set of the training data when determining a new base learner, further simplifying each sub-problem.

### 6.2.2 LPBoost and its LP master problem

To determine an optimal convex combination of a given finite set of fixed base learners, the LPBoost algorithm of [123] considers an LP optimization problem which maximizes the soft margin of the combined base learners in a boosting approach. They consider a given set of base learners  $\mathcal{H} = \{h_i | i \in [n]\}$ . The LP model determines an ensemble model  $f(x) = \sum_{i \in I} a_i h_i(x)$  that maximizes the soft margin of training data set  $\{X, Y\}$ :

$$\begin{aligned}
 & \max_{\alpha, \xi, \rho} \quad \rho - \lambda \sum_{j \in J} \xi_j \\
 \text{s.t.} \quad & y_j \sum_{i \in [n]} \alpha_i h_i(\mathbf{x}_j) + \xi_j \geq \rho, \quad j \in J \\
 & \alpha \in \Delta_n \\
 & \xi_j \geq 0, \quad j \in J \\
 & \rho \in \mathbb{R}
 \end{aligned} \tag{6.5}$$

where  $y_j \sum_{i \in [n]} \alpha_i h_i(\mathbf{x}_j)$  is the margin of data point  $(x_j, y_j)$  for a finite set of base learners. Variable  $\xi_j$  is the margin slack variable of  $(x_j, y_j) \in \mathcal{X} \times \{-1, 1\}$  and  $\lambda$  a positive penalty parameter, *e.g.*,  $\lambda = 0.1$  which influences the size of the non-zero dual

variables. A dual problem of (6.5) is given by

$$\begin{aligned}
 \min_{\mathbf{u}, \beta} \quad & \beta \\
 \text{s.t.} \quad & \sum_{j \in J} u_j y_j h(\mathbf{x}_j) \leq \beta, \quad h \in \mathcal{H} \\
 & \sum_{j \in J} u_j = 1 \\
 & 0 \leq u_j \leq \lambda, \quad j \in J
 \end{aligned} \tag{6.6}$$

In this formulation,  $u$  represents a dual solution of master problem (6.5), where  $u_j$  represents the miss-classification cost assigned to data point  $(x_j, y_j)$ . Following the column generation approach in [123], a new base learner  $\hat{h}$  can be generated or selected by solving the following pricing problem, which in the machine learning context represents the BL generation

$$\hat{h} = \arg \max_{h \in \mathcal{N}} \sum_{j \in J} u_j y_j h(x_j), \tag{6.7}$$

where  $\mathcal{N} \subset \mathcal{F}$  is a large finite set of base learners and  $u$  is the information of the master problem in terms of a dual solution of (6.5). Leveraging column generation to solve problem (6.5) provides a computationally efficient framework for managing a potentially vast set of base learners  $\mathcal{N}$  by iteratively operating on a smaller, more tractable subset  $\mathcal{H}$ . The objective (6.7) can be used as a termination criterion of the boosting method, see [123]. Note that the method in [123] is limited to a large set of predefined base learners. Our approach extends this framework by also considering the generation of new base learners. This is done by solving the less restrictive pricing problem, similar to the approach in [127]

$$\hat{h} = \arg \max_{h \in \mathcal{F}} \sum_{j \in J} u_j y_j h(x_j). \tag{6.8}$$

Since this method does not rely on a limited set of initial base learners, it is more flexible and potentially yields better solutions.

### 6.2.3 Sparse dual solution, dynamic data selection and BL generation

Numerical experiments show that the dual master problem (6.6) solution  $u$  is typically a sparse vector, with most of its components taking values of either 0 or  $\lambda$ . In linear programming, a sparse vector means that most of its entries are equal to 0. Let

$$R := \{(x_j, y_j) \in T : u_j > 0\} \tag{6.9}$$

## 6. OPTIMIZING ENSEMBLE MODELS VIA COLUMN GENERATION

---

be the set of data points with a non-zero dual value, which we refer to as active data points. Note, that a dual value of  $u_j = 0$  indicates that the associated data point is well-classified with a sufficient margin, as shown in [123]. Since  $u \in \Delta_{|J|}$  and  $u_j = \lambda$  if  $u_j > 0$ , the number of non-zero components of  $u$  is  $|R| = \frac{1}{\lambda}$  for  $\frac{1}{\lambda} \in \mathbb{N}$ . The relative number of active data points is  $\nu := \frac{1}{\lambda|T|}$ , such that  $|R| = \nu|T|$ . Using  $\nu$  as a parameter to control the number of active points, we can set  $\lambda = \frac{1}{\nu T}$  in (6.6) to influence the sparsity of the dual solution that allows us to identify critical data points and discard the irrelevant data points.

For the non-active data points in  $T \setminus R$ , the objective function contribution coefficient  $u_j = 0$  in BL generation sub-problem (6.8) is zero, such that the pricing problem (6.8) is equivalent to

$$\hat{h} = \arg \max_{h \in \mathcal{F}} \sum_{(x_j, y_j) \in R} u_j y_j h(x_j). \quad (6.10)$$

This means that we can limit the generation of new learners to the active data points in  $R$ , which tend to be miss-classified or well-classified with insufficient margin. Numerical experiments showed that, in practice, it is more effective to use a nonlinear objective function in the BL generation problem

$$\hat{h} = \arg \max_{h \in \mathcal{F}} \sum_{(x, y) \in R} \mathcal{L}(h(x), y) \quad (6.11)$$

rather than the traditional linear objective in (6.10). It is noteworthy that sub-problems in learning tasks are generally not solved to global optimality, attributable to their complexity. Note that since data set  $R$  is changing in each iteration, the related base learners tend to be more diverse which is an additional advantage of working with a smaller set of active points. Section 6.2.4 summarizes some results of [125], showing that the generalization error of an ensemble model is related to the diversity of its base learners.

### 6.2.4 Diversity measurement for ensemble

An interesting approach for analysing an ensemble model is to compute the diversity of its base learners. [125] show that the generalization error of an ensemble depends mainly on the diversity and accuracy of the base learners. We will also use the diversity concept for analysing the investigated ensemble models. The following describes this approach. Consider an error function for binary classification based on  $\text{margin}(f(x), y) = yf(x)$

### 6.3 A Data-Reduction Ensemble Learning Algorithm (DDA)

---

defined as

$$\text{Err}(f(x), y) := \begin{cases} 1 & \text{if } \text{margin}(f(x), y) = -1 \\ 0.5 & \text{if } \text{margin}(f(x), y) = 0 \\ 0 & \text{if } \text{margin}(f(x), y) = 1. \end{cases}$$

Hence,  $\text{Err}(f(x), y) = -\frac{1}{2}(yf(x) - 1)$ . The generalization error of  $f$  is defined by

$$G := \mathbb{E}[\text{Err}(f(x), y)].$$

The diversity of  $f$  at  $x \in \mathcal{X}$  is defined by the difference between ensemble error and the weighted average error of individual classifiers (base learners):

$$\begin{aligned} \text{div}(f(x), y) &:= \text{Err}(f(x), y) - \sum_{i \in I} \alpha_i \text{Err}(f_i(x), y) \\ &= \frac{1}{2} \text{margin}(f(x), y) - \frac{1}{2} \sum_{i \in I} \alpha_i \text{margin}(f_i(x), y). \end{aligned} \tag{6.12}$$

The relation between generalization error and diversity according to [125] is

$$G = \frac{1}{2} (1 - \mathbb{E}[\text{sign}(\gamma(x))]),$$

where  $\gamma(x) := \text{sign}(\text{div}(f(x), x)) - 2\text{div}(f(x), x)$ . This shows that increasing the diversity at points where  $\gamma(x)$  is small, will decrease  $G$ . In particular, the defined diversity measure and  $G$  varies in different ranges of diversity

- When  $\text{div}(f(x), y)$  is negative, the ensemble classifier has miss-classified instance  $x$ . Increasing diversity will improve the ensemble performance.
- When  $\text{div}(f(x), y)$  is positive, diversity needs to be decreased to increase the generalization performance of the ensemble.

More details of derivations for diversity and generalization errors can be found in [125]. Note that the ensemble model in [125] combines the classifier output of base learners  $f_i(x) \in \{-1, 1\}$  instead of score function  $h_i(x)$ , which is different from the ensemble described in Section 6.2. We show in Section 6.4 experimentally that for base learners obtained from solving sub-problem (6.11), the generalization error is decreasing and the average diversity is increasing iteratively.

### 6.3 A Data-Reduction Ensemble Learning Algorithm (DDA)

The initialization procedure, designed to generate an initial set of base learners, accelerates the computation process. The core of the algorithm involves solving a sequence

## 6. OPTIMIZING ENSEMBLE MODELS VIA COLUMN GENERATION

---

of master problems (6.6) to determine convex combinations of base learners, alongside solving sub-problems (6.10) or (6.11) to generate new base learners.

We present an ensemble learning algorithm, which iteratively generates base learners  $h_i(x) \in [-1, 1]$ . The algorithm defines low-dimensional BL generation sub-problems using dynamic data selection. The sub-problems can be solved using an arbitrary method. The initialization procedure, designed to generate an initial set of base learners, accelerates the computation process. The core of the algorithm involves solving a sequence of master problems (6.6) to determine convex combinations of base learners, alongside solving sub-problems (6.10) or (6.11) to generate new base learners. We show in Section 6.3.2 that if the sub-problems, *i.e.*, pricing problems, are solved exactly, then the algorithm converges to an optimal solution of problem (6.4), which is equivalent to problem (6.14).

### 6.3.1 Algorithm

Algorithm 6.1 describes the new ensemble learning algorithm. It consists of three phases, named initial, generation and refinement phase. The initial phase generates a base-line model  $h$  using procedure `SOLVEINITIAL` by fitting the training set, such that  $h$  initiates the base learner set  $\mathcal{H}$ . In the generation and refinement phases, diverse base learners are generated and added to base learner set  $\mathcal{H}$ . After adding a new base learner, the primal and dual solution  $(\alpha, u)$  of LP master problem (6.5) are used to update ensemble model  $\hat{f}$  by calling function `SOLVEPANDUPDATEENS`. Employing the optional generation phase in conjunction with the refinement phase typically yields superior results compared to using the refinement phase alone. The refinement phase also extends  $\mathcal{H}$ , but the base learners generated in this phase are fitted on a sub-set of active training data determined by  $u$ .

In order to accelerate the CG procedure, in the beginning, columns are generated using a start heuristic disregarding the master problem. This is motivated by a potentially poor search direction if the master problem is defined by few columns, see [124]. Our ensemble algorithm employs a similar strategy, which we refer to as the generation phase. This phase is based on error-driven bootstrap aggregation [128] and is outlined in Algorithm 6.2. Error-based bootstrapping generates a new data set  $S$  of size  $|S| = |T|$ . Unlike conventional bootstrapping that forms data set  $S$  by random sampling from the original data set with replacement, error-based bootstrapping adds to the set of miss-classified data points  $M \subset T$  a uniform sample of  $T \setminus M$  with replacement, such that it focuses on miss-classified instances from a base learner  $h$ . This procedure is denoted by function `BOOTSTRAPDATA(h, T)` in Line 6 of Algorithm 6.1.

### 6.3 A Data-Reduction Ensemble Learning Algorithm (DDA)

---

**Algorithm 6.1** Main Algorithm, Dynamic Data-Reduction Ensemble Learning

---

**Input:** Training data  $T = X \times Y \subset \mathcal{X} \times \{-1, 1\}$

**Parameters:** BL configuration and data reduction value  $\nu \in (0, 1)$  defining  $\lambda = \frac{1}{\nu|T|}$

**Output:**  $f(\cdot)$

```

1: function DDA
2:    $h \leftarrow \text{SOLVEINITIAL}(T)$                                      # initial phase
3:    $\mathcal{H} \leftarrow \{h\}, R \leftarrow \emptyset$ 
4:   repeat                                                         # generation phase (optional)
5:      $\hat{R} = R$ 
6:      $S \leftarrow \text{BOOTSTRAPDATA}(h, T)$ 
7:      $h \leftarrow \text{SOLVEBTSUBPROBLEM}(S)$ 
8:      $\mathcal{H} \leftarrow \mathcal{H} \cup \{h\}$ 
9:      $(\hat{f}, u) \leftarrow \text{SOLVEMPANDUPDATEENS}(\mathcal{H}, T)$ 
10:     $R \leftarrow \text{ACTIVEDATA}(T, u)$ 
11:  until  $|R \cap \hat{R}| > \sigma|R|$ 
12:  repeat                                                         # refinement phase
13:     $R \leftarrow \text{ACTIVEDATA}(T, u)$ 
14:     $h \leftarrow \text{SOLVEREFSUBPROBLEM}(R)$ 
15:     $\mathcal{H} \leftarrow \mathcal{H} \cup \{h\}$ 
16:     $(\hat{f}, u) \leftarrow \text{SOLVEMPANDUPDATEENS}(\mathcal{H}, T)$ 
17:  until stopping criterion
18:  return  $\text{sign}(\hat{f})$ 

```

---

The miss-classified set is obtained according to the predictions of  $h$  instead of preliminary ensemble  $\hat{f}$ . If instead the predictions of  $\hat{f}$  are used, *i.e.*,  $\text{BOOTSTRAPDATA}(\hat{f}, T)$  is called (see Algorithm 6.2), the generation phase might not terminate quickly and transition to the subsequent refinement phase takes longer. The set of miss-classified data points is updated in each iteration. In Lines 7 and 8 of Algorithm 6.1, a new base learner  $h_i(x)$  is fitted on  $S$  by  $\text{SOLVEBTSUBPROBLEM}(S)$  and added to set  $\mathcal{H}$ .  $\text{SOLVEMPANDUPDATEENS}(\mathcal{H}, T)$  in Line 9 of Algorithm 6.1 constructs a new master problem (6.6) given the updated base learner inferences and labels. The master problem is solved by an LP solver, and it returns the primal and dual solution. The ensemble score function  $\hat{f}$  is obtained via the primal solution which contains the base learner weight vector.

The details of  $\text{SOLVEMPANDUPDATEENS}$  can be found in Algorithm 6.3. The

## 6. OPTIMIZING ENSEMBLE MODELS VIA COLUMN GENERATION

---

**Algorithm 6.2** Generates error-based bootstrapping data set  $S$  of size  $|S| = |T|$

---

```

1: function BOOTSTRAPDATA( $h, T$ )
2:    $M \leftarrow$  MISSCLASSIFIEDDATA( $h, T$ )
3:    $S \leftarrow$  Uniform sample of  $|T| - |M|$  data points from  $T$ 
4:    $S \leftarrow S \cup M$ 
5:   return  $S$ 

```

---

sparse dual solution  $u$  is used to determine the set of active data (6.9). Function ACTIVE DATA( $T, u$ ) in Line 10 of Algorithm 6.1 selects active data from  $T$  fulfilling (6.9) and returns set  $R$ .

---

**Algorithm 6.3** Primal/dual solution of master-problem and update ensemble

---

```

1: function SOLVEMPANDUPDATEENS( $\mathcal{H}, T$ )
2:    $(\alpha, u) \leftarrow$  PRIMALDUALSOLVEMP( $\mathcal{H}, T, \lambda$ )
3:    $\hat{f} = \sum_{h_i \in \mathcal{H}} \alpha_i h_i$ 
4:   return  $(\hat{f}, u)$ 

```

---

The generation phase stops if the difference between former active data set  $\hat{R}$  and the current active set  $R$  is not significant. The stopping criterion parameter  $\sigma$  determines the minimum change in the active data set to continue the iterations in the generation phase. It is recommended to set its value to at least 0.2, indicating a change of 20%. Note that a small value of  $\sigma$  can result in the generation of many base learners that do not contribute to the ensemble. The refinement phase of Algorithm 6.1 augments the set  $\langle$  by adding base learners  $h_i(x)$ . It achieves this by solving the LP master problem, assigning weights to data points and training base learners on sufficiently weighted points. The function ACTIVE DATA() is called to determine the training set  $R$ , which is then used by SOLVEREF SUBPROBLEM() to solve either (6.11) or (6.10), depending on whether a linear or nonlinear objective is used. Since data points with large weights are selected, Parameter  $\nu$  is used in Algorithm 6.3 to calculate  $\lambda$  in (6.5) adjusting the sparsity of the dual solution to provide a limit on the number of considered data points. We discuss the choice of its value in Section 6.4. The termination criterion of the refinement phase in Algorithm 6.1 can be based on the dual feasibility criterion from the original LPBoost algorithm, see [123]. If the following condition is met, the base learner can not further improve the ensemble:

$$\sum_{j \in J} u_j y_j h_i(x_j) \leq \beta + \epsilon, \quad (6.13)$$

### 6.3 A Data-Reduction Ensemble Learning Algorithm (DDA)

---

where  $\epsilon$  is a small positive value. There are two main advantages of the column generation approach in Algorithm 6.3 for solving (6.4). First, it decomposes the problem into simpler sub-problems compared to simultaneously optimizing all base learners and their combination. Secondly, it allows us to only consider a small sub-set of active data points which results in easier sub-problems for determining new base learners, *i.e.*, easier pricing problems.

#### 6.3.2 Convergence

We analyse the convergence properties of Algorithm 6.1 and focus on the soft margin loss function, although the algorithm is not limited to this specific loss function. We show that Algorithm 6.1 computes an optimal solution of the extended master problem

$$\begin{aligned}
 \max_{f, \xi, \rho} \quad & \rho - \lambda \sum_{j \in J} \xi_j \\
 \text{s.t.} \quad & y_j f(\mathbf{x}_j) + \xi_j \geq \rho, \quad j \in J \\
 & \xi_j \geq 0, \quad j \in J \\
 & \rho \in \mathbb{R}, \quad f \in \text{conv}(\mathcal{F}).
 \end{aligned} \tag{6.14}$$

Note that problem (6.14) is equivalent to problem (6.4) when using the soft margin loss function and allowing sufficiently many base learners. The proof requires a few intermediate results that are presented in Propositions 6.3, 6.4 and 6.5.

**Proposition 6.3.** Problem (6.14) is convex.

Proof. Consider the feasible set of (6.14) as

$$\Omega := \{(f, \xi, \rho) : y_j f(\mathbf{x}_j) + \xi_j \geq \rho, j \in J, \rho \in \mathbb{R}, f \in \text{conv}(\mathcal{F})\}.$$

Let  $g_1 = (f_1, \xi_1, \rho_1)$ ,  $g_2 = (f_2, \xi_2, \rho_2)$  and  $g = (f, \xi, \rho) := tg_1 + (1-t)g_2$ ,  $t \in [0, 1]$ . From  $g_1, g_2 \in \Omega$  we have

$$y_j f(\mathbf{x}_j) + \xi_j = y_j (tf_1(\mathbf{x}_j) + (1-t)f_2(\mathbf{x}_j)) + t\xi_{1j} + (1-t)\xi_{2j} \geq t\rho_1 + (1-t)\rho_2 = \rho.$$

Since  $f = tf_1 + (1-t)f_2 \in \text{conv}(\mathcal{F})$ ,  $g \in \Omega$ , which proves  $\Omega$  is convex.  $\square$

**Proposition 6.4.** Problem (6.14) is equivalent to the semi-infinite dual-LP:

$$\begin{aligned}
 \min_{\mathbf{u}, \beta} \quad & \beta \\
 \text{s.t.} \quad & \sum_{j \in J} u_j y_j h(\mathbf{x}_j) \leq \beta, \quad h \in \mathcal{F} \\
 & \sum_{j \in J} u_j = 1 \\
 & 0 \leq u_j \leq \lambda, \quad j \in J.
 \end{aligned} \tag{6.15}$$

## 6. OPTIMIZING ENSEMBLE MODELS VIA COLUMN GENERATION

---

Proof. Consider the Lagrangian of (6.5)

$$c(f, \rho, \xi, \mu) := \rho - \lambda \sum_{j \in J} \xi_j + \sum_{j \in J} u_j (y_j f(\mathbf{x}_j) + \xi_j - \rho).$$

Then

$$\begin{aligned} c^*(u) &:= \max \left\{ c(f, \rho, \xi, u) : (f, \xi, \rho) \in \text{conv}(\mathcal{F}) \times \mathbb{R}_+^{|J|} \times \mathbb{R} \right\} \\ &= \max_{\rho \in \mathbb{R}} \left( 1 - \sum_{j \in J} u_j \right) \rho + \max_{\xi \in \mathbb{R}_+^{|J|}} \sum_{j \in J} (u_j - \lambda) \xi_j + \max_{f \in \text{conv}(\mathcal{F})} \sum_{j \in J} u_j y_j f(\mathbf{x}_j) \\ &= \max_{f \in \text{conv}(\mathcal{F})} \sum_{j \in J} u_j y_j f(\mathbf{x}_j) = \max_{h \in \mathcal{F}} \sum_{j \in J} u_j y_j h(\mathbf{x}_j). \end{aligned} \quad (6.16)$$

Since problem (6.14) is convex, strong duality applies, *i.e.*, the optimal value of (6.14) is equal to

$$c^* := \min \left\{ c^*(u) : u \in \mathbb{R}_+^{|J|} \right\} = \min \left\{ \beta : \beta \geq c^*(u), u \in \mathbb{R}_+^{|J|} \right\}.$$

For dual points with  $c^*(u) < \infty$ , we have  $\sum_{j \in J} u_j = 1$  and  $0 \leq u_j \leq \lambda$ ,  $j \in J$ . From (6.16) follows

$$c^* = \min \left\{ \beta : \beta \geq \max_{h \in \mathcal{F}} \sum_{j \in J} u_j y_j h(\mathbf{x}_j), \sum_{j \in J} u_j = 1, 0 \leq u_j \leq \lambda, j \in J \right\}.$$

Hence, the optimal value of (6.15) is  $c^*$ .  $\square$

**Proposition 6.5.** If the loop in the refinement phase of Algorithm 6.1 is not stopped, and pricing problem (6.8) is solved exactly, the primal-dual solution  $(\hat{f}, u)$  converges towards a primal-dual solution of extended master problem (6.14).

Proof. Let  $U$  denote the set of dual points  $u$  computed by Algorithm 6.1. Since  $U$  is bounded, there exists a subsequence  $\{u\} \subset U$  converging to a point  $\hat{u}$ , where  $u$  is a dual solution of (6.6). Let  $\mathcal{H}$  be the set of base learners related to the subsequence. Then  $\hat{u}$  solves (6.6) regarding  $\mathcal{H}$ . Let  $\hat{h}$  denote the solution of pricing problem (6.8) at  $\hat{u}$ . Since  $\hat{u}$  is a limit point and by continuity of the dual function, we have

$$\max_{h \in \mathcal{H}} \sum_{j \in J} \hat{u}_j y_j h(\mathbf{x}_j) = \sum_{j \in J} \hat{u}_j y_j \hat{h}(\mathbf{x}_j) = \lim_{u \rightarrow \hat{u}} \max_{h \in \mathcal{F}} \sum_{j \in J} u_j y_j h(\mathbf{x}_j).$$

Hence, we can replace  $\mathcal{H}$  by  $\mathcal{F}$  in (6.6), which shows that  $\hat{u}$  is a solution of (6.15). Since (6.15) is the dual of (6.14), the proposition follows.  $\square$

The theoretical results show that Algorithm 6.1 can generate an optimal solution of problem (6.4). Keep in mind that this requires the BL generation sub-problem to be solved exactly. Even if the dimension of a sub-problem is much smaller than of the original problem, it may still not be feasible to solve the problem exactly. The results show that the algorithm follows a solid decomposition technique and if the sub-problems are solved approximately then the algorithm becomes more heuristic.

## 6.4 Numerical Experiments

We study the characteristics of Algorithm 6.1 on a set of binary classification problems (6.1). Algorithm 6.1 was implemented in Python, as a part of the general ensemble framework DECOLEARN<sup>1</sup> which contains learning methods to generate four types of initial models and base learners with any configuration in the generation and refinement phases: MLP, CNNs, CART-based decision trees and XGBoost models [129, 130, 131]. The experiment focuses on convolutional neural networks (CNN) base learners. The pricing sub-problems (6.10) or (6.11) which yields new base learners, is solved using TensorFlow 2.3. This approach does not guarantee global optimality for the pricing sub-problems, but allows to efficiently work with non-trivial base learners. Solving the pricing sub-problems to guaranteed global optimality in a reasonable time is still not feasible for CNN base learners of relevant size. DECOLEARN also uses GUROBI 9.5.1 for solving the LP master problem. All computational experiments were carried out on a computer with i7-1165G7 4-Core 2.80 GHz CPU and 16GB RAM.

The test instances for binary classification were obtained from the CIFAR-10 data set [132]. The data set contains 60,000  $32 \times 32$  color images in 10 classes, with 6,000 images per class. There are 50,000 training images and 10,000 test images, divided into five training batches and one test batch, each with 10,000 images. The test batch contains 1,000 randomly-selected images from each class, while the training batches contain the remaining images in random order. However, some training batches may contain more images from one class than another, with exactly 5,000 images from each class in total. For our experiments on binary classification using CIFAR-10, we chose two pairs of labels with varying degrees of similarity. The labels of categories 3 and 8 are considered easy and labels 3 and 2 hard to distinguish.

For base learner generation with a non-linear loss, we used CNNs with VGG architectures [133] in both the generation and refinement phase based on the binary cross

---

<sup>1</sup>see <https://github.com/ouyang-w-19/decolearn>

## 6. OPTIMIZING ENSEMBLE MODELS VIA COLUMN GENERATION

entropy:

$$\mathcal{L}(h(x), y) = y \log(Q(h(x)) + \epsilon) + (1 - y) \log(1 - Q(h(x)) + \epsilon), \quad (6.17)$$

where the clipping is defined as  $Q(y) := \max(\min(y, 1 - \epsilon), \epsilon)$ . The tolerance  $\epsilon$  is added for stability reasons and is a small positive number, e.g.  $10^{-7}$ . If the base learners have an output range of  $[0, 1]$ , the output is re-scaled to  $[-1, 1]$  or alternatively to  $\mathbb{R}$ . Note that (6.17) requires a different binary label type for base learner training, namely  $y \in \{0, 1\}$ . The complexity of the models were adjusted to the degree of similarity between two labels. The details of the algorithm parameters, model types and hyper-parameters are documented in Appendix A, Tables 6.8 and 6.9.

Evaluating the algorithm involves several key performance indicators, including ensemble binary accuracy, diversity, and the solution to the master problem at each iteration. The diversity is defined per data point. To describe the overall diversity of an ensemble, we will use the average diversity over the training set defined by

$$\overline{\text{div}}(f) := \frac{1}{|T|} \sum_{j \in T} \text{div}(f(x_j), y_j).$$

We compare different variants of the algorithm, which differ in the types of sub-problems and parameter configurations. One of the variants solves sub-problems with active data only and another variant encompasses all training data when solving the "sub-problem". For the data selection based ensemble method, we use two variants. The linear variant (LL) applies the standard LPBoost with a sub-problem that minimizes the reduced cost (pricing problem) in the refinement phase with a linear expected loss function

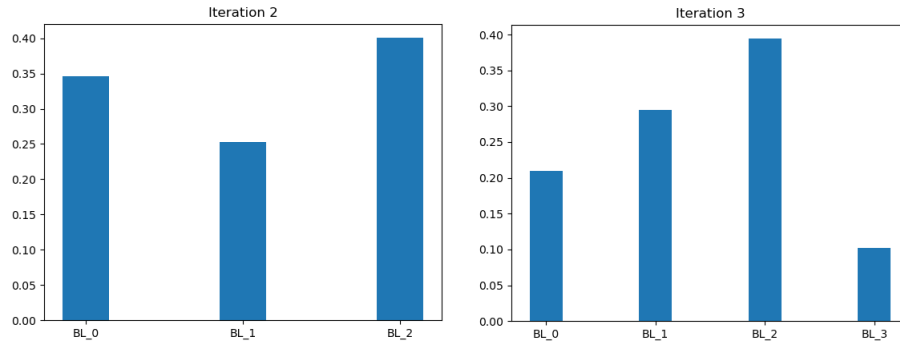
$$\max_{h \in \mathcal{F}} \sum_{(x_j, y_j) \in T} u_j y_j h(x_j). \quad (6.18)$$

We compare that to a nonlinear variant called NL which considers sub-problems on the active set using the nonlinear loss function (6.17). We will first focus on the characteristics of the variants of DDA in the generation phase and then on their behaviours in the refinement phase.

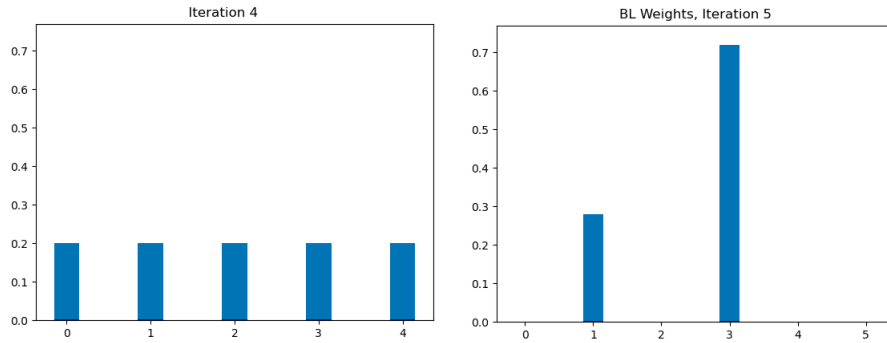
### 6.4.1 Generation Phase

The generation phase is optional and serves the purpose of generating base learners through an error-based bootstrapping approach. The master problem is only solved to calculate an ensemble in each iteration and its corresponding active set  $R$ , providing a stopping criterion for this phase. Active columns in the master problem refer to the

columns with non-zero weights of the corresponding base learner. In the best case, the generated columns remain active in the refinement phase. This is illustrated when running the DDA generation phase on an instance of CIFAR-10, where after 2 iterations the phase ends and all generated base learners stay active, see Figure 6.1. The case in Fig. 6.2 illustrates a run over MNIST instances where four out of six BLs generated in the generation phase have zero weights in the next phase.



**Figure 6.1:** Generated base learners with their master problem weight. The generation phase ends after 2 iterations. In the next iteration, the old base learners remain active together with a new active column.



**Figure 6.2:** Six base learners (corresponding weights on the y-axis) are generated in the generation phase. As shown by the second figure, most of the columns become inactive in the next iteration.

We conducted experiments by varying the stopping criterion of the generation phase. Specifically, we compare a variant with a fixed number of 5 iterations to one that uses the criterion based on the relative progress described in Line 11 of Algorithm 6.1. The results, shown in Table 6.1, present the average performance over 20 runs of DDA,

## 6. OPTIMIZING ENSEMBLE MODELS VIA COLUMN GENERATION

one iteration following the generation phase, for both stopping criterion variants. The number of inactivate columns corresponds to the ensemble obtained from the generation phase. A relative criterion has a relatively smaller number of inactive columns than a fixed number of 5 iterations. The number of iterations for the relative criterion is 2.6 on average. The last row gives the number of miss-classified instances in the validation set of the ensemble after the generation. This illustrates that with more iterations, the result generalizes better. Specifically, continuing bootstrap aggregation for an average of twice as many iterations results in an ensemble that generalizes more effectively, with an improvement factor of approximately two.

**Table 6.1:** Effect of fixed versus relative stopping criterion on number of columns becoming inactive, number of iterations and miss-classified data in the validation set

Stopping criterion:	Fix 5 iterations	relative progress
Inactivated col. (frac.)	$0.49 \pm 0.15$	$0.29 \pm 0.23$
Iterations	5	$2.6 \pm 0.66$
nr. Miss-classified train	$2084 \pm 949$	$2271 \pm 1064$
nr. Miss-classified val.	$335 \pm 184$	$649 \pm 461$

To study the contribution of the generation phase to DDA, we run the algorithm with and without the phase for an instance of CIFAR-10 with binary labels (categories 3 and 2). Ten runs of the algorithm are performed and Table 6.2 presents the obtained statistical results. We measured the classification quality for the final ensemble as the number of miss-classified data points of the training set and of the validation set. The results show the average values over ten runs with the standard deviation. Table 6.2 indicates that the generation phase can improve the quality of the final ensemble, that is, the number of miss-classified data points (testing/validation set) decrease from 723 to 576 on average. On the other hand, the average computational time of DDA is doubled when adding the generation phase.

The next section will illustrate that a big improvement (lower) in miss-classification score may occur immediately after the generation phase at a relatively low computational cost.

### 6.4.2 Refinement Phase

The refinement phase adds base learners to set  $\mathcal{H}$ . We used two data sets that vary in difficulty to compare the behavior for the variants LL and NL. Experiment A uses

**Table 6.2:** DDA with and without generation phase on categories 3 and 2 of the CIFAR-10 data set

DDA	without gen. phase	with gen. phase
Initial nr. Miss-classified train	4882 $\pm$ 303	
Initial nr. Miss-classified val	978 $\pm$ 54	
nr. Miss-classified train	3024 $\pm$ 1511	1607 $\pm$ 772
nr. Miss-classified val	723 $\pm$ 200	576 $\pm$ 34
Total time (s)	205.84 $\pm$ 177.15	412.22 $\pm$ 234.26

categories 3 and 8 from the CIFAR-10 dataset and is considered easy. Experiment B run on categories 3 and 2 from the same dataset and is considered more challenging. We measure the miss-classification counts over the validation and training set and the algorithm tun time.

**Table 6.3:** Comparing LL and NL variant on category 3 and 8 instances of CIFAR-10, Experiment A

	LL variant	NL variant
nr. Runs	10	10
Initial nr. Miss-classified train	3283 $\pm$ 1342	
Initial nr. Miss-classified val	654 $\pm$ 267	
nr. Miss-classified train	1924 $\pm$ 767	840 $\pm$ 274
nr. Miss-classified val	389 $\pm$ 146	194 $\pm$ 49
Total time (s)	65.21 $\pm$ 40.61	51.53 $\pm$ 20.85

Tables 6.3 and 6.4 show the results of the two variants over the two experiments. In both cases, the nonlinear loss function variant appears to lead to less miss-classifications. The performance on the validation set approximates the generalization capability of the method. The total computational time varies significantly, but seems reasonable for both variants.

Figure 6.3 contains various graphs depicting different measures for the experiment. The x-axis gives the number of iterations (or number of BLs) and the y-axis the miss-classification score of training set and validation set for the LL and NL variant of the DDA algorithm. Experiment A refers to the run on the CIFAR-10 instance with categories 3 and 8. Note that the miss-classification goes down fast when using the nonlinear loss function.

## 6. OPTIMIZING ENSEMBLE MODELS VIA COLUMN GENERATION

**Table 6.4:** Comparing LL and NL variant results on categories 3 and 2 instances of CIFAR-10, Experiment B

	LL variant	NL variant
nr. Runs	10	10
Initial nr. Miss-classified train	$4661 \pm 312$	
Initial nr. Miss-classified val	$933 \pm 60$	
nr. Miss-classified train	$4248 \pm 642$	$1607 \pm 772$
nr. Miss-classified val	$861 \pm 124$	$576 \pm 34$
Total time (s)	$717.19 \pm 717$	$412.22 \pm 234.26$

We now study the behaviour of both variants of DDA using the data of Experiment B. Instead of using a stopping criterion in the refinement phase, we fix the number of iterations. By doing so, we can illustrate that the miss-classification does not improve anymore for the LL variant after the sixth iteration. Using the nonlinear loss function, the miss-classification keeps an improving trend when generating more base learners, *i.e.*, with increasing iterations. This is illustrated in Figure 6.4. Interestingly, the average diversity has the same tendency as illustrated in Figures 6.4e and 6.4f.

In our experiment, we observe the behaviour of the algorithm when we vary the data reduction parameter  $\nu$ , which influences the number of data points that are considered active, setting  $\lambda = \frac{1}{\nu|T|}$  in (6.6). We use the data of Experiment B and use five runs. We measure the miss-classification of the data and the total running time. The resulting values are reported in Table 6.5.

**Table 6.5:** Results for Experiment B data varying data reduction  $\nu \in (0, 1)$

$\nu$	Miss-classified train	Miss-classified val.	Total time
0.1	$2380 \pm 908$	$625 \pm 72$	$254.27 \pm 227.60$
0.4	$2520 \pm 862$	$592 \pm 114$	$134.47 \pm 77.12$
0.5	$2061 \pm 486$	$511 \pm 52$	$138.88 \pm 93.65$
0.6	$2324 \pm 163$	$510 \pm 26$	$112.93 \pm 63.15$
0.7	$2273 \pm 132$	$502 \pm 24$	$135.25 \pm 49.40$
0.8	$3828 \pm 670$	$787 \pm 112$	$236.41 \pm 227.02$
1.0	$4365 \pm 442$	$874 \pm 78$	$94.87 \pm 29.49$

**Table 6.6:** Highest base learner and ensemble accuracy with strong base learners (s. 6.10), data reduction  $\nu = 0.5$ 

	Highest BL Accuracy		Highest Ensemble Accuracy	
	Train	Validation	Train	Validation
Run 1	94.16	78.80	<b>96.80</b>	<b>82.00</b>
Run 2	94.33	78.20	<b>94.91</b>	<b>78.65</b>
Run 3	95.08	79.25	<b>95.47</b>	<b>79.45</b>
Run 4	95.25	78.65	<b>95.48</b>	78.65
Run 5	74.68	73.60	<b>75.05</b>	<b>77.40</b>

### 6.4.3 Comparison of base learner and ensemble accuracy

In our last series of experiments we compare the highest base learner and ensemble accuracy, to measure the predictive prowess of the model. The first experiment is run with base learners that have an average validation accuracy of 77.7% over five runs. They are considered to be strong base learners. The second experiment is run with weaker base learners that have an average validation accuracy of 72.47% over five runs. The hyper-parameter configuration for the strong and weak base learners is listed in Table 6.10. They are considered to be weak base learners.

**Table 6.7:** Highest base learner and ensemble accuracy with weaker base learners (s. 6.10), data reduction  $\nu = 0.5$ 

	Highest BL Accuracy		Highest Ensemble Accuracy	
	Train	Validation	Train	Validation
Run 1	71.44	71.6	<b>81.8</b>	<b>79.3</b>
Run 2	74.41	73.2	<b>83.37</b>	<b>78.65</b>
Run 3	73.38	73.35	<b>83.7</b>	<b>78.4</b>
Run 4	73.15	70.65	<b>84.08</b>	<b>79.85</b>
Run 5	75.25	73.55	<b>82.54</b>	<b>78.45</b>

## 6.5 Discussion and Conclusion

This chapter introduces and analyses an algorithm that iteratively enhances an ensemble of base learners using a column generation approach. Following an initial improvement through error-based bootstrapping, the algorithm employs a data selection

## 6. OPTIMIZING ENSEMBLE MODELS VIA COLUMN GENERATION

---

strategy that prioritizes active data points, i.e., those with non-zero dual values, which are particularly sensitive to the margin (objective). The implementation of the algorithm is part of the ensemble framework DECOLEARN and it allows various experimental setups to reliably improve an initial solution by adding base learners to an ensemble. All experiments were executed in a binary classification setting. The focus of the algorithm is on the refinement phase that performs linear programming on an iteratively growing set of score functions to maximize a soft margin. The LP constructs an ensemble from weighted score functions, populates the set of active data until an adjustable limit is reached and sets a stopping criterion for the algorithm. In numerical experiments, we investigated a variant with the traditional linear loss function and a nonlinear loss function based on cross entropy.

Experiments demonstrate that the nonlinear variant, which focuses on solving nonlinear sub-problems defined by active data only, is potentially more effective in reducing generalization error compared to the variant that employs a linear loss function. Additionally, experiments with this variant reveal a correlation between generalization error and a broader definition of ensemble diversity, based on base learners  $h_i(x) \in [-1, 1]$ , thus relaxing the original assumption set by [125]. We hypothesize that, because the underlying distribution of the active data set can deviate significantly from the original distribution, the proposed algorithm enhances both accuracy and diversity, particularly when the active data set –by definition, containing the most challenging data points–undergoes significant changes. We summarize our findings by the following bullet points:

- We have presented a column generation based decomposition algorithm DDA for optimizing ensemble models, optimizing the ensemble one base learner at a time and only using a sub-set of the training data for each base learner.
- We showed that algorithm DDA computes an optimal ensemble model in the convex hull of a given ANN model space.
- The results in Table (6.6) and (6.7) show that the generated ensemble models have the potential to reach an accuracy which is significantly better than the accuracy of the individual base learners, demonstrating that the ANN model space used in the experiments is a proper sub-set of its convex hull.
- Using reduced data and non-linear sub-problems seems beneficial to the LPBoost approach.
- We pose the following hypothesis: Diversity emerges from fitting the model with active data set  $R$ . The size of  $R$  can be chosen deliberately small. However, a small active data set may deteriorate accuracy.

- Model diversity as defined by [125], seems to be correlated to the generalization error.
- We found that the generation phase can improve the final ensemble results at higher computational cost.
- The implementation of the generation phase accelerates convergence. It is also recommended not to use a fixed iteration number as the stopping criterion during the generation phase because it significantly increases the likelihood of generating ineffective base learners.”
- Ultimately, the additional computational cost to generate an ensemble model does not justify the minor accuracy improvements of the model.

## Appendix: Configuration of Algorithm Meta-parameters and Model Hyper-parameters

The settings of the algorithm meta-parameters in the numerical experiments are summarized in Table 6.8. The model type of the base learners obtained in both the generation and refinement phases is set as CNN. The threshold of active data is set as zero, that is, all instances with non-zero dual values are considered active. The meta-parameter  $\nu$  of Algorithm 6.3 is taken as  $\nu = 0.2$ , while the meta-parameter  $\sigma$  of the termination condition for the generation phase in Algorithm 6.1 is set to  $\sigma = 0.2$ .

**Table 6.8:** Decolearn Metaparameters

Parameter Name	Value
Nu	2e-1
Generation model type	CNN
Refinement model type	CNN
Combining Mechanism	Averaging
Active Data Threshold $\phi$	0.0
Min. Generation Phase Change $\sigma$	2e-1

**Base Learner Hyper-parameters for DDA** The hyper-parameters and configurations for the base learners in the DDA algorithm are presented in Table 6.9. The configurations of base learners in the generation phase are slightly different from the ones in the refinement phase, *i.e.*, the number of episodes in the refinement phase is increased from five to eight. The NL variant uses binary cross entropy, while the LL

## 6. OPTIMIZING ENSEMBLE MODELS VIA COLUMN GENERATION

variant is based on a linear loss. Table 6.10 shows the base learner configuration in the DDA algorithm of base learners that significantly differ in their accuracy.

**Table 6.9:** Configuration of base learners for NL and LL comparison.

BCE: BinaryCross-entropy, L: Linear

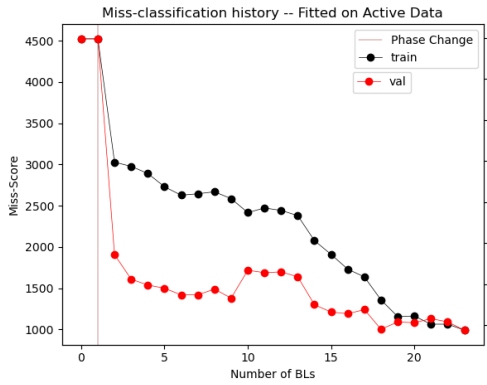
Parameter Name	Generation	Refinement
Architecture	VGG	VGG
Learning Rate	1e-3	1e-3
Loss	BCE/LL	BCE/LL
Activation Function	tanh	tanh
Last Layer Activation	tanh	tanh
Batch Size	50	50
Episodes	5	8
Layers	(32, 32, 16)	(32, 32, 16)
Filter Size	(2, 2)	(2, 2)
Max. Pooling Filter	(2, 2)	(2, 2)
FC-Layers	(8, 1)	(8, 1)

The investigation of this chapter has been published in [134].

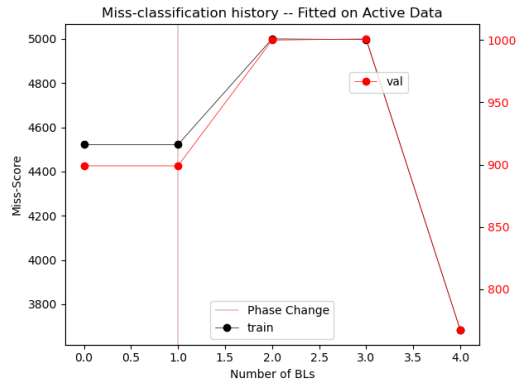
**Table 6.10:** Configuration of base learners. BCE: Binary Cross-entropy.

Parameter Name	Strong BLs		Weak BLs	
	Generation	Refinement	Generation	Refinement
Architecture	VGG	VGG	VGG	VGG
Learning Rate	1e-3	1e-3	1e-3	1e-3
Loss	BCE	BCE	BCE	BCE
Activation Function	tanh	tanh	tanh	tanh
Last Layer Activation	sigmoid	sigmoid	relu	relu
Batch Size	25	25	25	25
Episodes	7	7	7	7
Layers	(64, 32, 32)	(64, 32, 32)	(64, 32, 32)	(64, 32, 32)
Filter Size	(2, 2)	(2, 2)	(2, 2)	(2, 2)
Max. Pooling Filter	(2, 2)	(2, 2)	(2, 2)	(2, 2)
FC-Layers	(15, 1)	(15, 1)	(15, 1)	(15, 1)

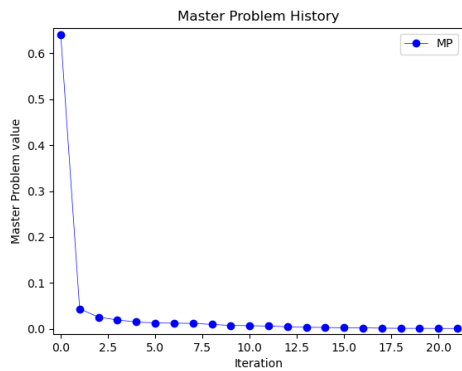
## 6. OPTIMIZING ENSEMBLE MODELS VIA COLUMN GENERATION



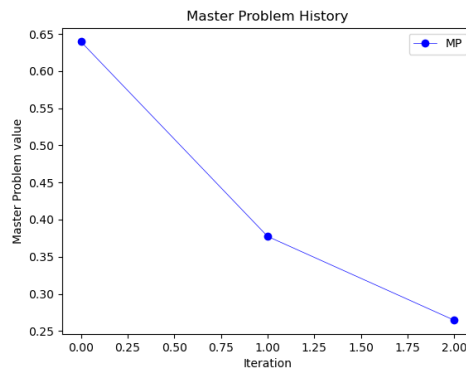
(a) Miss-classification increasing iterations NL variant DDA



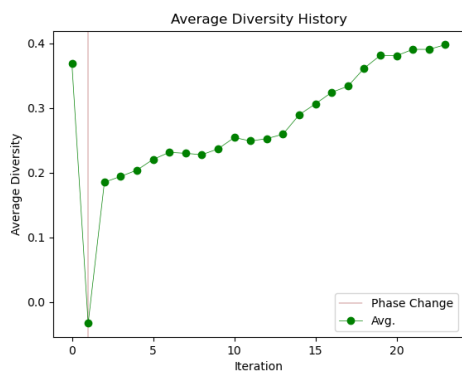
(b) Miss-classification increasing iterations LL variant DDA



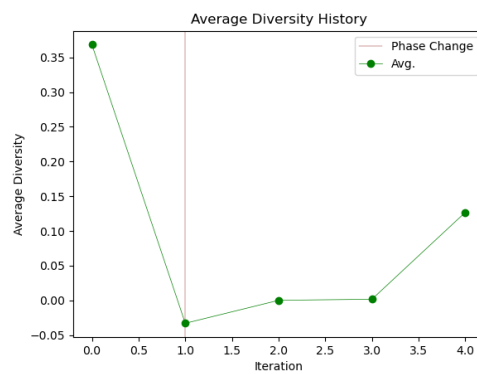
(c) Master objective function value over refinement iterations for NL variant



(d) Master objective function value over refinement iteration for LL variant

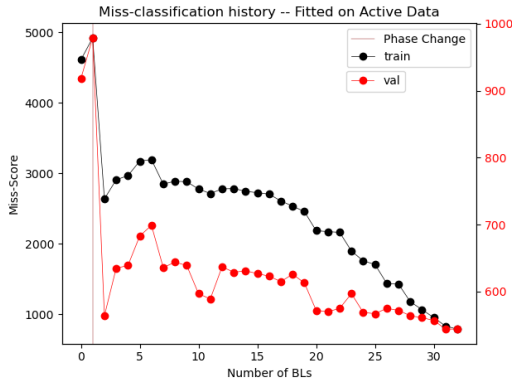


(e) Average diversity development for NL variant DDA

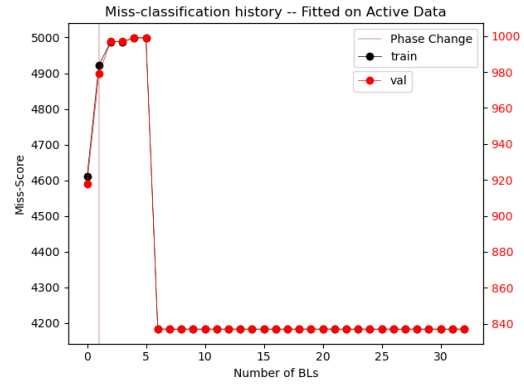


(f) Average diversity development LL variant DDA

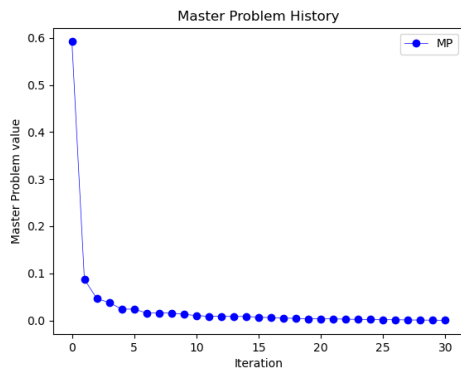
**Figure 6.3: Experiment B.** Miss-classification for increasing iterations (generated BLs) for the LL and NL variants



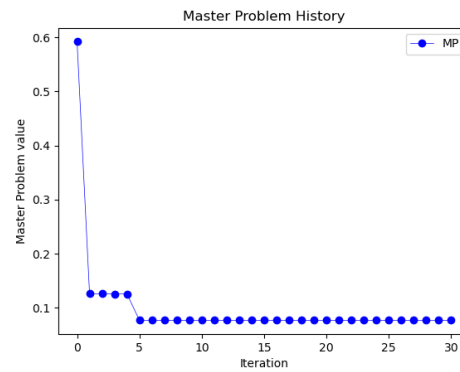
(a) Miss-classification increasing iterations  
NL variant DDA



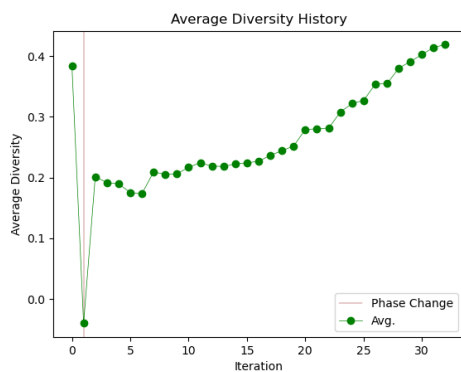
(b) Miss-classification increasing iterations  
LL variant DDA



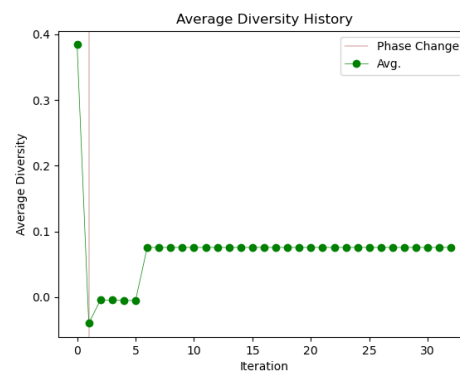
(c) Master objective function value trend of  
the NL variant for refinement phase



(d) Master objective function value trend of  
the LL variant for refinement phase



(e) Average diversity trend of NL variant



(f) Average diversity trend of LL variant

**Figure 6.4:** Experiment fixing 30 refinement phase iterations. Miss-classification of data points for NL and LL variant



## Latent Space Boosting

The previous chapter introduced a linear methodology for constructing ensembles. This chapter presents a non-linear ensemble generation technique that leverages the properties of neural networks. Concretely, the approach exploits the fact that neural networks, as composite functions, map input data through multiple layers of abstraction. For classification tasks, the last layer encodes the distribution over class predictions, representing the highest level of abstraction. The hidden layer activations (outputs) – often referred to as latent variables – encode rich representations at lower levels of abstraction. This chapter introduces an innovative approach that constructs ensembles not from the final output layer but instead from these latent variables. The approach operates on the assumption that latent variables, when taken as a whole from one layer, offer a much richer representation of the data distribution. By forming the ensemble from latent variables and separately training a discriminator on top of them, the proposed method enhances the generalization capacity of the ensemble. To further improve the accuracy of the ensemble, the training loss function is augmented with an additional diversity term. This term explicitly measures the difference between the parameters of each base learner and the aggregate parameters of the ensemble. The objective is to encourage base learners to make accurate predictions while establishing parameter configurations that are sufficiently distinct from those of other learners. Ideally, each base learner should occupy an orthogonal direction in the function space, ensuring maximum independence and minimal redundancy among base learners.

## 7.1 Introduction

Currently minor improvements in image classification tasks are achieved by significantly increasing model complexity. This trend has been ongoing for the last years and high performance models now usually have millions of parameters [15, 135, 136]. Inspired by the results of [137], we investigate the potential of Ensemble Learning (EL) to iteratively extend an ensemble of feature extractors and fit a classifier (combination layer) on the features as latent variables. A new EL algorithm is implemented as part of the EL-framework DECOLEARN, containing the LPboost EL-method of [134]. The main differences of the new EL approach compared to other work of the authors [134] are

1. Formulating a non-linear network-based ensemble learning problem to solve classification tasks, replacing linear ensemble construction methods.
2. The base learners are used to generate non-linear feature transformations. These transformations are integrated into the neural network through a combination layer.

Hence, our contribution to the work of [137] is to generate a novel ensemble architecture through a boosting approach in an efficient manner. Efficiency is defined to be the ratio between accuracy and model complexity. An approximate measure for model complexity is the total number of base and classifier model parameters. The base learners, serving as feature transformers, are intentionally overfitted on a subset of the training data, aligning with the approach of [137]. We sketch the idea of boosting in the latent space in Section 7.2 and a description of the algorithm in Section 7.3. We reflect on the research question in Section 7.4.

## 7.2 Latent Space Boosting

A general learning task is to learn a function  $f : \mathbb{R}^l \rightarrow \mathbb{R}^m$  to approximate an unknown mapping  $\Psi : \mathcal{X} \rightarrow \mathcal{Y}$  between the input space  $\mathcal{X}$  and output space  $\mathcal{Y}$ . Typically the sets  $\mathcal{X}$  and  $\mathcal{Y}$  are not known and we only have a finite set of matching observations  $\{(\mathbf{x}_j, \mathbf{y}_j)\}_{j \in J} \subset \mathcal{X} \times \mathcal{Y}$ . We consider multi-class classification problems, that is  $\mathcal{Y} = [M] := \{1, 2, \dots, M\}$ . Given a training set  $T = X \times Y$  and a model space  $\mathcal{F}$ , the problem aims to obtain a predictive model  $f^* \in \mathcal{F}$  which minimizes an expected loss function  $\mathcal{L}$ , i.e. a risk, regarding hypothesis  $f \in \mathcal{F}$

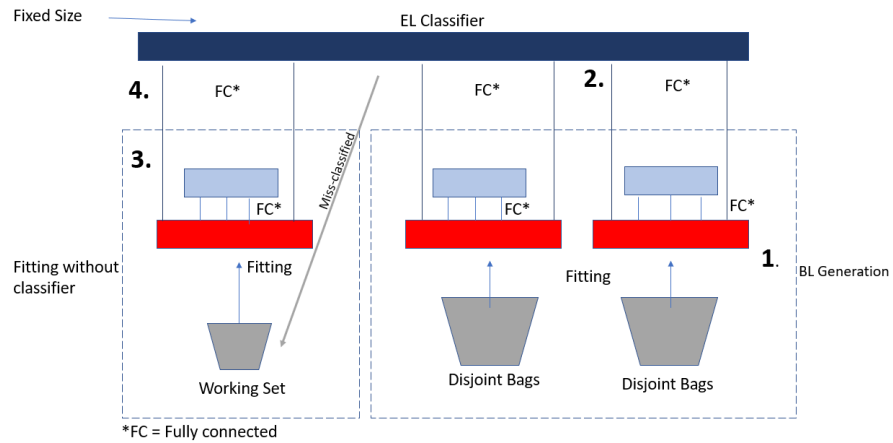
$$f^* = \operatorname{argmin}\{\mathbb{E}[\mathcal{L}(f(x), y)] : f \in \mathcal{F}\}, \quad (7.1)$$

where  $\mathbb{E}$  denotes the expectation regarding a joint probability distribution  $P(x, y)$  over  $\mathcal{X} \times \mathcal{Y}$ . The training set  $T$  is thought of as drawn i.i.d. from  $P(x, y)$ .

Let  $\phi$  denote a final layer and  $g$  a transformer (transformation layer) and  $f = \phi \circ g$  be an approximated solution of (7.1). The image space of  $g$  is called *latent space*. A latent space Base Learner (BL) is any weak transformer that maps the features from the original space to points in the latent space. The classifier takes the transformations of the base learners as input. Together they compromise the new ensemble architecture.

Let  $\mathcal{G} = \{g_1, \dots, g_Q\}$  be a set of latent space BLs. Consider a combination layer MLP, represented as  $\hat{\phi}(z)$  and a transformation layer  $G = (g_q)_{q \in [Q]}$  mapping from  $X$  to  $Z^Q$ . Here,  $Q = |\mathcal{G}|$  denotes the index of transformers in the ensemble and  $Z^Q$  represents the latent space considering all transformers. Each base learner is fitted on a subset of data  $S \subset X \times Y$ , called working set. "Fitting" is a standard learning problem and there is a large variety of efficient tools for solving such problems to at least local optimality.

An ensemble learning (EL) classifier  $\hat{f}(x, w) = \hat{\phi}(G(x, w_1); w_2)$  is a parameterized machine learning model, where  $(w_{1q})_{q \in [Q]}$  is the parameter vector of the transformation layer and  $w_2$  is the parameter vector of the combination layer. Figure 7.1 depicts this concept. The complexity of the transformers are arbitrary. Experiments with both low and high model complexities will be carried out to investigate performance.



**Figure 7.1:** 1. Base learners are fitted on disjoint bags, Point 1. The classifier is fitted on  $\mathcal{G}$  while  $w_1$  is being fixed (Point 3). Generation of working set from miss-classified data points and the new base learner are fitted on the working set. parameters  $w_2$  is fixed (Point 4). Retrain the classifier on the new  $\mathcal{G}$

## 7. LATENT SPACE BOOSTING

---

### 7.3 Algorithm

---

#### Algorithm 7.1 Latent Space Boosting

---

1. Partition training dataset  $T$  into  $p$  subsets  $S_1, \dots, S_p$
  2. For all  $j \in [p]$ : solve `LATENTSPACESUBPROBLEM`( $S_j$ ) and add the transformation  $\tilde{g}$  to  $\mathcal{G}$
  3. Compute  $\hat{f}$  by solving `COMBINATIONLEARNINGPROBLEM`( $\mathcal{G}$ )
  4. Determine the set  $S$  of miss-classified data points of  $\hat{f}$
  5. Solve `DIVERSELATENTSPACESUBPROBLEM`( $S$ ) and add the transformation  $\tilde{g}$  to  $\mathcal{G}$
  6. Compute  $\hat{f}$  by solving `COMBINATIONLEARNINGPROBLEM`( $\mathcal{G}$ )
  7. Repeat steps 4-6 until stopping criterion is reached
  8. Return  $f(x) = \operatorname{argmax}_{m \in [M]} \hat{f}_m(x)$
- 

Algorithm 7.1 describes the new boosting method to be investigated. It generates initial base learners by solving the sub-problem `LATENTSPACESUBPROBLEM`( $S$ ):

$$\tilde{\phi} \circ \tilde{g} = \operatorname{argmin} \left\{ \sum_{(x,y) \in S} \mathcal{L}(f(x, w), y) : w \in \mathcal{W}_{\mathcal{CNN}} \right\}, \quad (7.2)$$

where  $\mathcal{L}$  denotes the loss of point  $(x, y)$  and  $\mathcal{W}_{\mathcal{CNN}}$  denotes the parameter space of a convolutional neural network, although the concept is also applicable to other learners. Each iteration of the main loop (steps 4-6) generates a new (diverse) base learner fitted on the miss-classified data by solving sub-problem

`DIVERSELATENTSPACESUBPROBLEM`( $S$ ):

$$\tilde{\phi} \circ \tilde{g} = \operatorname{argmin} \left\{ \sum_{(x,y) \in S} \mathcal{L}(f(x, w), y) - \delta \|\bar{w}(\mathcal{G}) - w\|^2 : w \in \mathcal{W}_{\mathcal{CNN}} \right\}, \quad (7.3)$$

where  $\bar{w}(\mathcal{G})$  denotes the average network parameters of the previously generated base learners of  $\mathcal{G}$ , as suggested in [138]. An ensemble learning classifier is obtained by solving the sub-problem `COMBINATIONLEARNINGPROBLEM`( $\mathcal{G}$ ):

$$\hat{f} = \operatorname{argmin} \left\{ \sum_{(x,y) \in T} \mathcal{L}(f(x, w), y) : w \in \mathcal{W}_{\mathcal{MLP}}(\mathcal{G}) \right\}, \quad (7.4)$$

where  $\mathcal{W}_{\mathcal{MLP}}(\mathcal{G})$  denotes the parameter space of a MLP regarding base learners  $\mathcal{G}$ .

## 7.4 Results and Conclusion

This chapter presented an exploratory investigation into a novel ensemble learning strategy, termed Latent Space Boosting. The approach deviates from conventional ensemble methodologies by aggregating the latent representations of base learners, rather than their final-layer outputs. While this concept introduces an innovative perspective to ensemble learning, its practical implications remain limited.

Numerical experiments indicate that the proposed method results in a slight reduction in empirical generalization loss and are therefore not discussed in detail. However, the improvement is marginal at best, and the performance gains are insufficient to justify the additional computational complexity introduced by each base learner. Therefore, while the strategy demonstrates theoretical interest, it does not offer tangible benefits for real-world applications.

In conclusion, Latent Space Boosting should be regarded primarily as an academic exercise that enriches the understanding of ensemble strategies, but it lacks immediate applicability in practical scenarios.



## 8.1 Conclusion

The main objective of this thesis was to develop a novel DRL algorithm informed by an extensive literature review, analysis of existing RL and DRL approaches, and exploratory experiments with statistical models. To this end, a maximum-entropy, state-of-the-art, Cramér-based DRL algorithm, termed Cramér-based Distributional Soft Actor-Critic, C-DSAC, was proposed and rigorously investigated. In Chapter 2, foundational concepts of RL, DRL, and entropy-augmented RL were reviewed. Section 2.3 provided a detailed examination of moment-agnostic and fixed-moment distributional approaches. Gaussian and Gaussian Mixture Models were analysed as viable modelling choices for the latter. Experimental results indicated that Gaussian models are both expressive and numerically stable, making them suitable for distributional modelling. However, the restrictive assumption of a specific return distribution introduces a parametric approximation bias. Empirical evidence in Chapter 3 showed that the benefits of using Gaussian models outweigh the additional bias.

The main contribution of this thesis, the C-DSAC algorithm, was introduced in Chapter 3. C-DSAC is a novel off-policy, maximum-entropy distributional RL Algorithm within the Actor-Critic framework, demonstrating superior performance compared to existing RL algorithms. The C-DSAC convergence properties were thoroughly analysed and it was demonstrated that C-DSAC converges in the policy evaluation setting. Moreover, under mild assumptions described both in Chapter 3 and in Section 2.1.2, C-DSAC was shown to converge in the control setting. Subsequent analysis revealed that policy iteration with C-DSAC leads to convergence to the optimal value distribution and optimal policy. A noteworthy discovery during this investigation is

## 8. CONCLUSIONS AND FUTURE WORK

---

the phenomenon of confidence-driven model updates: state-action values with high variance, induced either by approximation noise or inherent stochasticity in the environment, are interpreted as having low confidence in its accuracy, resulting in updates that have a less effect on the approximator parameters. This mechanism was proven to mitigate overestimation bias effectively. Empirical evaluations in complex benchmark environments demonstrated that C-DSAC achieves higher returns in fewer iterations compared to established RL algorithms, albeit with some trade-off in stability. **These findings support the conclusion that C-DSAC surpasses current state-of-the-art RL/DRL algorithms.**

Chapter 4 explored a novel model-augmented RL meta-approach tailored for robotic applications, named Reinforcement Learning-Inverse Kinematics (RL-IK). This approach decouples the planning and execution tasks, assigning an RL algorithm for motion planning and the Levenberg-Marquardt algorithm for movement execution. Experimental results demonstrated the superiority of RL-IK over vanilla RL variants, indicating high potential for addressing challenging robotic control tasks.

Chapters 6 and 7 shifted focus to the supervised learning (SL) aspects of the thesis. Chapter 6 examined an ensemble learning approach that combines strong base learners by solving a linear program defined from their predictions on the training set. Chapter 7 investigated a non-linear ensemble approach where each base learner is defined to be a sub-graph of the entire ensemble graph. After adding each base learner graph to the ensemble, a multi-layer perceptron (MLP) classifier is retrained on the resulting extended latent representation of the input. While both methods reduce empirical generalization error, the non-linear graph-based approach performs slightly better. However, the marginal improvement in generalization error does not justify the substantial computational overhead of training additional complex base learners. These results suggest that future research in SL should prioritize the development of innovative network architectures rather than ensemble methods.

In conclusion, this thesis achieves the defined objectives. The development, analysis, and theoretical validation of the novel state-of-the-art DRL algorithm C-DSAC represent a significant contribution to the field. Parallel research efforts in SL provided insights into ensemble learning, ultimately recommending a shift toward innovative network designs over computationally intensive ensemble methods. This shift should mirror the transformative progression from recurrent neural networks to transformers for sequential prediction problems. Future research directions are discussed in the following section.

## 8.2 Future Work

Several promising research avenues can be explored to further advance the field of (distributional) reinforcement learning and – more generally – machine learning. Investigating fundamental aspects of machine learning can even contribute to a broader understanding of human intelligence and consciousness. In the immediate future, improvements to C-DSAC may include enhancing value representations through multi-modal distributions. Leveraging the meta-learning techniques discussed in Section 2.1.5 could further optimize performance. Additionally, incorporating an adaptive temperature adjustment mechanism, as proposed in [139], may improve the exploration-exploitation trade-off, leading to more robust learning dynamics.

In the broader context, reinforcement learning is poised to play a dominant role in the pursuit of Artificial General Intelligence (AGI). For instance, the Alberta Plan [53] outlines a research strategy that emphasizes continuous RL in real-world environments while considering computational constraints and multi-agent interactions. This plan introduces a generalized agent concept that extends traditional RL architectures by integrating additional perception, transition, policy, and critic modules, each optimizing distinct objectives. The transition module, in particular, holds significant potential for implementing temporal abstraction, akin to RL-IK (Chapter 4). However, unlike RL-IK, which relied on a predefined low-level controller, the Alberta Plan envisions a more flexible, model-free learning paradigm where models are learned and continuously refined through experience. Future research could explore the adaptation of these principles within the context of distributional RL.

Another promising research direction is the integration of RL with foundation models, such as large language models (LLMs). Foundation models encode vast amounts of knowledge about the world from a human perspective. This knowledge can be leveraged to effectuate a warm-start mechanism for decision-making problems. Conversely, RL algorithms are increasingly used in fine-tuning LLMs, as evidenced by the recent application of a pure RL flywheel to optimize model behaviour [140]. While RL-driven training loops hold promise for other domains, current approaches to most problems remain computationally prohibitive, necessitating the innovation of more efficient RL training paradigms.

Beyond algorithmic advancements, a paradigm shift in the architecture of function approximators may drive the next major leap in performance. This paradigm shift would be analogous to the transition from RNNs to transformers in sequence modelling tasks. Notably, approaches based on the transformer architecture have already

## 8. CONCLUSIONS AND FUTURE WORK

---

been explored in offline RL settings with "Decision Transformers", yielding promising yet debated results [141, 142]. Another critical consideration is the role of hardware acceleration in RL and ML in general. Just as Tensor Processing Units exhibit large performance advantages in deep learning, tailored hardware could significantly enhance RL scalability. Quantum Reinforcement Learning, grounded in fundamental quantum principles, represents a promising frontier. The inherent properties of quantum mechanics – superposition, interference, and entanglement – could significantly enhance state-space exploration, action sampling, and model optimization. However, defining quantum RL metrics remains a challenging task. An alternative approach to novel hardware involves the integration of biological neural networks into machine learning. The human brain demonstrates exceptional efficiency in both spatial and energetic terms, suggesting that biologically inspired architectures may hold the key to more advanced approaches in this field with profound implications to RL.

Ultimately, continued advancements in machine learning will not only drive technological progress but may also provide deeper insights into the fundamental nature of intelligence itself. Given the striking parallels between biological and artificial learning systems – both in their ability to interact with the world and form abstractions – it is the author's firm belief that future developments in this field will contribute significantly to the understanding of cognition, creativity, and intelligence.



## Publications arising from this Thesis

The investigation for the present thesis resulted in a number of publications.

### A.1 Journal Publication

Vanya Aziz, Ivo Nowak, Ouyang Wu, Eligius M. T. Hendrix, and Jan Kronqvist. On optimizing ensemble models using column generation. *J Optim Theory Appl*, 203(1):1794–1819, 2024

Impact factor JCR 2023: 1.6. Subject categories: Applied Mathematics - Q2 (89/332); Operations Research & Management Science - Q3 (65/106);

### A.2 Submitted Journal Publication

Vanya Aziz, Ivo Nowak, and Eligius M. T. Hendrix. Distributional RL Using the Cramér distance. *J. Mach. Learn. Res.*, *submitted*, 2024

Impact factor JCR 2023: 4.1. Subject categories: Automation & Control Systems - Q1 (21/84); Computer Science, Artificial Intelligence - Q2 (45/197).

### A.3 Publication in International Conference Proceedings

Vanya Aziz, Ivo Nowak, Ouyang Wu, Jan Kronqvist, and Eligius M. T. Hendrix. Boosting via Latent Representation. In *Proceedings OLA'2023 International Conference on Optimization and Learning*, pages 150–153, 2023



### 8.4 Resumen

El término "Industria 4.0" describe la cuarta revolución industrial y se caracteriza por la integración de la tecnología digital en los procesos de manufactura. Los conceptos transformadores de la Industria 4.0 permiten una producción económica con tamaños de lote radicalmente pequeños, lo que exige niveles sin precedentes de automatización y adaptabilidad. Estos requisitos en las instalaciones de producción requieren sistemas que actúen de manera autónoma y se optimicen a sí mismos. El campo del aprendizaje automático ofrece soluciones prometedoras para lograr los objetivos de la Industria 4.0, particularmente al permitir la toma de decisiones basada en datos y mecanismos de control adaptativos. Esta disertación se centra en el Deep Reinforcement Learning, (DRL), un enfoque basado en redes neuronales para resolver Procesos de Decisión de Markov en espacios de alta dimensionalidad con dinámicas de transición desconocidas. La principal contribución de esta tesis es el desarrollo de un novedoso algoritmo de aprendizaje por refuerzo distribuido de vanguardia dentro del marco Actor-Crítico de máxima entropía. Este algoritmo, denominado "Cramér-based Soft Distributional Soft Actor-Critic" (C-DSAC), demuestra un rendimiento superior en comparación con otros algoritmos de aprendizaje por refuerzo, especialmente en entornos con espacios de alta dimensionalidad y dinámicas complejas. Su rendimiento se debe en parte a un fenómeno que surge en el Aprendizaje por Refuerzo Distribucional basado en la métrica de Cramér, denominado "actualizaciones del modelo impulsadas por la confianza". Este mecanismo garantiza que el aproximador de la función de valor se actualice de manera más conservadora cuando la confianza en sus estimaciones es baja. Se presentan

justificaciones teóricas para el algoritmo, demostrando su convergencia en el escenario de evaluación de políticas y, bajo suposiciones ampliamente aceptadas, también en el escenario de control. Más allá de la investigación algorítmica fundamental, esta tesis contribuye a la aplicación práctica del aprendizaje por refuerzo en la robótica. Dada la importancia crucial de los sistemas robóticos multiarticulados en la tecnología de producción moderna, se diseña un meta-algoritmo de aprendizaje por refuerzo denominado "Reinforcement Learning - Inverse Kinematics" (RL-IK). Este enfoque mejora la aplicabilidad del aprendizaje por refuerzo en tareas de control robótico al acelerar significativamente la convergencia hacia políticas cercanas a óptimas en comparación con los métodos estándar de aprendizaje por refuerzo. Un requisito esencial para las aplicaciones del aprendizaje por refuerzo en sistemas de control del mundo real es la percepción de la máquina para la identificación de estados. Para abordar los desafíos en este campo, esta tesis explora nuevos enfoques de Aprendizaje Supervisado (Supervised Learning, SL), validados en tareas de clasificación de imágenes, con un enfoque en estrategias de aprendizaje por conjuntos (ensemble learning). Para proporcionar una comprensión teórica más profunda, los métodos se analizan desde la perspectiva de la optimización clásica y se comparan con enfoques heurísticos. Esta investigación ofrece conocimientos valiosos sobre las dificultades de encontrar conjuntos de clasificadores óptimos a nivel global.

## Capítulo 1

El primer capítulo establece el marco de referencia para esta tesis. Introduce el concepto de aprendizaje automático junto con varios términos esenciales y conceptos clave necesarios para comprender los temas más avanzados que se exploran en los capítulos siguientes. El capítulo incluye una breve digresión técnica sobre uno de los principales impulsores del éxito de las herramientas contemporáneas de aprendizaje automático: los poderosos aproximadores de funciones, en particular las redes neuronales artificiales, inspiradas en los principios de funcionamiento de las neuronas biológicas. Finalmente, este capítulo expone las principales motivaciones y objetivos de la investigación presentada en esta tesis, a saber, contribuir a la comprensión del aprendizaje por refuerzo distribuido y desarrollar un nuevo algoritmo de aprendizaje por refuerzo distribuido que idealmente supere a los algoritmos de vanguardia actuales. Además, se investigan métodos de aprendizaje por conjuntos como un enfoque complementario para tareas de aprendizaje supervisado. Esta sección detalla los objetivos específicos que guían la investigación desarrollada en los capítulos siguientes. La principal métrica para la

mejora con respecto a los algoritmos existentes de aprendizaje por refuerzo se define como la convergencia a niveles deseados de desempeño en menos iteraciones que los algoritmos de vanguardia actuales.

## Capítulo 2

Este capítulo introduce al lector a la importancia del "Reinforcement Learning" (RL) y contextualiza la investigación presentada en la tesis. Comienza con una breve presentación de la historia de la investigación en RL. Se presentan los Procesos de Decisión de Markov con dinámicas de transición desconocidas y grandes espacios de solución como el marco en el que se pueden describir los problemas de RL. Tras una descripción intuitiva del problema, se define formalmente y se establecen algunas suposiciones ampliamente utilizadas y poco restrictivas sobre el mismo. En este capítulo se revisan diversas estructuras dentro del RL: enfoques basados únicamente en críticos, aprendizaje por diferencia temporal (TD Learning), búsqueda directa de políticas, arquitecturas actor-crítico, función de ventaja y estimador de ventaja generalizada, algoritmos on-policy y off-policy, así como el aprendizaje por refuerzo multiagente. A continuación, se detallan diversos desafíos inherentes al campo del RL, incluyendo el dilema de exploración-explotación, el sesgo de sobreestimación en métodos basados en Q, la complejidad temporal en entornos de gran dimensión y la brecha sim-to-real. Dado su papel crucial, la ingeniería de recompensas se aborda por separado en la Sección 2.1.4 y se describe en detalle. En particular, se advierte al lector sobre cómo señales de recompensa aparentemente intuitivas pueden inducir comportamientos no deseados en los agentes, y se analiza este problema en el contexto de agentes hipotéticamente superinteligentes. Además, en la Sección 2.1.4.1, se expone el desafío paradigmático en RL y se introduce la noción de esparsidad en la función de recompensa. La densidad de la recompensa influye de manera significativa en la eficiencia del aprendizaje: aunque aprender a partir de recompensas escasas es el escenario ideal, a menudo conduce a una convergencia lenta o fallida. El diseño de recompensas densas puede acelerar el entrenamiento, pero requiere una considerable intervención humana, asemejándolo al aprendizaje supervisado y potencialmente entrando en conflicto con los principios de aprendizaje autónomo. Asimismo, se explica el concepto de alineación de valores, un aspecto clave para generar modelos de control que respeten principios éticos humanos, en el contexto del modelado de recompensas. En la Sección 2.1.5, se presentan algunos de los principales meta-algoritmos de RL que han demostrado reducir el sesgo de sobreestimación, mejorar la estabilidad numérica en el entrenamiento de modelos y priorizar

experiencias más valiosas para acelerar la convergencia. También se introducen diversas técnicas de inicialización de modelos para establecer parámetros iniciales bien definidos, y se detalla la atenuación progresiva de la tasa de aprendizaje (learning rate annealing), que contribuye a la estabilidad en las actualizaciones del modelo. Para ilustrar los conceptos expuestos en secciones anteriores, la Sección 2.1.6 introduce varios algoritmos de RL, cada uno destacando distintos aspectos fundamentales de este campo. Algunos de estos algoritmos también sirven como base para métodos más avanzados. La Sección 2.1.6.1 explica el algoritmo clásico de Q-Learning en su versión tabular, introduciendo el concepto clave de funciones de valor Q, sobre el cual se basa la principal contribución de esta tesis. En la Sección 2.1.6.3, se presenta REINFORCE, un algoritmo de RL basado en políticas directas que incorpora la función de ventaja. Las siguientes dos secciones introducen PPO y DDPG, dos algoritmos representativos del paradigma actor-crítico en configuraciones on-policy y off-policy, respectivamente. Posteriormente, el capítulo explora el concepto de aprendizaje por refuerzo con máxima entropía. Los conceptos novedosos presentados en esta tesis tienen su fundamento en este marco, reconocido por su robustez y alto rendimiento en tareas complejas y de alta dimensionalidad. La notación utilizada en esta sección refleja las bases teóricas del esquema actor-crítico. La formulación del objetivo con entropía aumentada se detalla en la Sección 2.2.1, seguida de la presentación del marco soft actor-critic y sus justificaciones teóricas. El capítulo concluye con una exposición sobre el "Distributional Reinforcement Learning" (DistRL). A diferencia de los métodos tradicionales basados en el valor esperado, DistRL se centra en modelar la distribución completa de los retornos. Esta aproximación, relativamente reciente, ha demostrado en múltiples ocasiones un rendimiento superior en comparación con sus variantes base y, por lo tanto, es de gran interés en la comunidad de investigación en RL. En la sección dedicada a DistRL, se presentan demostraciones rigurosas sobre la convergencia en los escenarios de evaluación y control bajo la métrica de Wasserstein. La elección del modelo de distribución en DistRL tiene un impacto significativo en el rendimiento, por lo que se exploran diversas estrategias. Se introduce el concepto de distribuciones de momentos libres, que no imponen restricciones sobre sus momentos y permiten una adaptación arbitraria de los mismos a costa de un alto costo computacional. En contraste, las distribuciones de momentos fijos están completamente caracterizadas por sus dos primeros momentos, lo que las hace computacionalmente más eficientes. Se llevan a cabo experimentos exhaustivos utilizando modelos Gaussianos y Gaussianos Mixtos, lo que lleva a la selección de núcleos Gaussianos individuales como el modelo de distribución preferido para las principales contribuciones de esta tesis, las cuales se detallan en el siguiente capítulo.

## Capítulo 3

El Capítulo 3 describe el trabajo principal de esta tesis, el desarrollo de un novedoso algoritmo de "Reinforcement Learning" (RL) de última generación en el marco de Actor-Crítico de máxima entropía denominado "Cramér-based Distributional Soft Actor-Critic" (C-DSAC).

Basado en la evaluación de políticas utilizando la métrica de Cramér, se deriva una función de pérdida regularizada que mejora la eficiencia del aprendizaje y la robustez del algoritmo. Este algoritmo representa un avance respecto al estado del arte en aprendizaje por refuerzo y cuenta con un marco teórico riguroso que garantiza la convergencia y la estabilidad de las actualizaciones, incluso en presencia de un alto nivel de ruido ambiental.

En la Sección 3.2, se explora un nuevo enfoque distribucional para RL y se presenta un método para actualizar adaptativamente los modelos de valores-Q distribucionales utilizando la distancia de Cramér regularizada. Sobre estas bases, y en combinación con el aprendizaje por refuerzo de máxima entropía, se propone el marco Cramér-based Distributional Soft Actor-Critic (C-DSAC). Además, se derivan las fórmulas correspondientes para su implementación con redes neuronales y se presentan datos experimentales que demuestran un rendimiento de última generación en los entornos evaluados.

Para explicar la eficiencia del nuevo enfoque, se realiza un análisis exhaustivo de la dinámica del algoritmo. Se pone énfasis en el comportamiento del algoritmo en presencia de errores de aproximación y ruido del sistema, y se derivan ecuaciones para los gradientes del valor. Se demuestra que la adaptación a los valores-Q es más lenta en pares de estado-acción donde la varianza del modelo distribucional es grande, es decir, cuando la confianza en su estimación es baja. Además, se enfatiza cómo estas ecuaciones revelan el mecanismo inherente del algoritmo para mitigar el sesgo de sobreestimación.

Los teoremas en esta sección proporcionan una base teórica para el enfoque distribucional basado en Cramér para SAC. La Sección 3.2.1 define la distancia de Cramér regularizada e introduce la noción de una métrica ideal. Se concluye que, debido a tres propiedades— invarianza ante sumas, sensibilidad a la escala y muestras de gradiente no sesgadas—, la distancia de Cramér regularizada constituye una métrica ideal.

En la Sección 3.2.2, se define el objetivo distribucional de máxima entropía. La recompensa se modela explícitamente como una variable aleatoria, lo que permite el uso de una función de valor distribucional. Similar al algoritmo SAC (Sección 2.2.3), la entropía de la política se mide mediante la probabilidad logarítmica negativa.

La evaluación de políticas en C-DSAC se demuestra tener un punto fijo único en la Sección 3.2.3. La demostración utiliza las propiedades de Cramér para mostrar que la función de valor distribucional predice eventualmente el retorno de una política arbitraria. La mejora de políticas se describe en la Sección 3.2.4. Se sigue un razonamiento similar al de la mejora de políticas en SAC, con la adición de que la media de la función de valor distribucional es considerada en el paso de mejora. Posteriormente, se demuestra que la iteración de políticas para C-DSAC converge a la distribución de valores y política óptimos cuando los pasos de evaluación y mejora se alternan.

En la Sección 3.3, los resultados teóricos de las secciones anteriores se utilizan para derivar funciones objetivo parametrizadas adecuadas para la implementación. Para manejar los grandes espacios de estados y acciones en tareas complejas, se emplean redes neuronales como aproximadores de funciones paramétricas con el objetivo de optimizar sus parámetros. Se asume que la variable de valor sigue una distribución gaussiana, es decir, C-DSAC se implementa como un algoritmo de momento fijo, como se describe en la Sección 2.3.5.1. Los detalles de implementación incluyen estrategias para la evaluación y mejora de políticas. Basado en estos resultados, el algoritmo completo de C-DSAC se presenta en el Algoritmo 3.1.

El análisis en la Sección 3.4 muestra el efecto de la distancia de Cramér regularizada en la dinámica del algoritmo y tiene como objetivo proporcionar una posible explicación de su rendimiento superior observado. Además, se demuestra que C-DSAC posee un mecanismo inherente para contrarrestar el sesgo de sobreestimación.

El análisis concluye que la actualización de la red neuronal de distribución de valores mediante gradientes de la función de pérdida de Cramér regularizada  $J_Z(\theta)$  mejora el modelo especialmente en aquellas regiones donde la varianza—medida de la incertidumbre—es pequeña. Se argumenta que la actualización del modelo en presencia de alto ruido (baja confianza) es más propensa a la sobreestimación. Sin embargo, dado que el gradiente tiene un impacto reducido en los parámetros del modelo en estos casos, este efecto mitiga naturalmente la sobreestimación.

La sección experimental (3.5) evalúa C-DSAC en entornos de referencia de diversa complejidad, cuantificando su rendimiento para sustentar las afirmaciones de la tesis respecto a su superioridad sobre enfoques alternativos.

Los experimentos se realizan con los mismos valores de hiperparámetros que en [30], excepto por los específicos de C-DSAC (por ejemplo, límites en la desviación estándar, etc.). Se reconoce que detalles auxiliares al algoritmo principal y la "optimización a nivel de código" pueden impactar significativamente el rendimiento de un algoritmo [90]. Se asume que la implementación original utilizada para los resultados base en

[30] no aplicó técnicas de optimización adicionales como recorte de la función de valor, inicialización ortogonal, escalado de capas, ajuste de la tasa de aprendizaje, recorte de observaciones y recompensas, o normalización de observaciones y recompensas. Para ser coherente con estos supuestos, en los experimentos no se utilizan estas técnicas de optimización a nivel de código.

El agente es entrenado en los entornos de OpenAI: Hopper-v4 ( $\mathcal{S} \times \mathcal{A}$ )  $\in \mathbb{R}^{11} \times \mathbb{R}^3$ , Ant-v4 ( $\mathcal{S} \times \mathcal{A}$ )  $\in \mathbb{R}^{111} \times \mathbb{R}^8$ , Humanoid-v4 ( $\mathcal{S} \times \mathcal{A}$ )  $\in \mathbb{R}^{376} \times \mathbb{R}^{17}$ , HalfCheetah-v4 ( $\mathcal{S} \times \mathcal{A}$ )  $\in \mathbb{R}^{17} \times \mathbb{R}^6$  y Walker2d-v4 ( $\mathcal{S} \times \mathcal{A}$ )  $\in \mathbb{R}^{17} \times \mathbb{R}^6$ . Siguiendo las recomendaciones de [92], se reporta la acción media de la política objetivo cada 1.000 iteraciones, con una duración de episodio de 1.000 pasos. El agente es entrenado por un millón de iteraciones en cada entorno y los resultados se promedian sobre cinco semillas aleatorias. Se muestra que C-DSAC supera a SAC en 4 de los 5 entornos evaluados. En particular, en el entorno más complejo, Humanoid-v4, C-DSAC supera claramente a SAC y alcanza una recompensa media de aproximadamente 5730, valor que SAC solo logra después de aproximadamente cuatro millones de iteraciones. Se concluye que C-DSAC es superior a SAC.

## Capítulo 4

Este capítulo explora la aplicación de algoritmos de aprendizaje por refuerzo (RL) en el campo de la robótica y presenta un enfoque novedoso diseñado específicamente para tareas de manipulación robótica. El algoritmo propuesto, denominado Reinforcement Learning - Inverse Kinematics (RL-IK) [93], integra el conocimiento de la cinemática del robot y puede interpretarse como un algoritmo de RL basado en modelos (model-based RL). Esta integración permite un aprendizaje más eficiente al aprovechar el modelo físico subyacente del robot, particularmente porque el movimiento del robot y la planificación de la tarea se separan y asignan a módulos distintos. El enfoque de este capítulo está centrado en robots de múltiples articulaciones en un espacio tridimensional (3D), controlados mediante políticas aprendidas dentro del marco actor-crítico (actor-critic, AC). El objetivo del algoritmo es mejorar el rendimiento de los robots en tareas que implican manipulación y posicionamiento, combinando cinemática inversa con aprendizaje por refuerzo para la planificación y toma de decisiones en un nivel de abstracción superior. El capítulo comienza con una introducción detallada a la cinemática directa, estableciendo el estándar Denavit-Hartenberg para la descripción de matrices de transformación. A continuación, se definen las transformaciones rotacionales y traslacionales, las cuales se utilizan para especificar la posición y orientación

del efector final (end-effector) del robot en el espacio. Se introduce el problema de la cinemática inversa y se exponen los diversos métodos numéricos empleados para calcular la configuración articular del robot a partir de un estado objetivo. Tras analizar diferentes enfoques, se establece que el método de Levenberg-Marquardt es una solución numérica suficientemente eficiente para abordar este problema. En la Sección 4.3, se describe la idea central de RL-IK: en el contexto del aprendizaje por refuerzo aplicado a la robótica, resulta natural distinguir entre la planificación del movimiento y la ejecución física de los desplazamientos del robot. Esta separación permite una arquitectura de control más modular e interpretable. RL-IK se basa explícitamente en esta segmentación de tareas: en cada paso de decisión, el agente de RL genera una posición y orientación objetivo para el efector final en función de su política. Las ventajas de RL-IK se exponen en la misma sección. En primer lugar, reduce el espacio de acciones. Independientemente del número de grados de libertad (Degrees of Freedom, DoF) del robot, el espacio de acción se mantiene fijo y definido únicamente por las posiciones y orientaciones del efector final. En segundo lugar, los movimientos del robot resultan más suaves. En tercer lugar, se evitan singularidades en la configuración del robot. Por último, la eficiencia en la utilización de muestras (sample efficiency) mejora de manera drástica en la tarea gracias a RL-IK. Para el algoritmo propuesto, el diseño de la función de recompensa desempeña un papel crucial. Se plantea un esquema de recompensas que incentiva al agente a resolver la tarea de manera eficiente sin infringir restricciones de seguridad ni condiciones operativas. Finalmente, en la Sección 4.5, RL-IK se compara con variantes básicas de RL en la resolución de una tarea robótica. Los resultados obtenidos evidencian que RL-IK acelera significativamente el proceso de aprendizaje.

## Capítulo 5

Este capítulo marca un cambio en el enfoque de esta tesis, pasando del Aprendizaje por Refuerzo (RL) al Aprendizaje Supervisado (SL), y sirve como introducción a los capítulos posteriores. Aunque el núcleo principal de este trabajo ha sido el aprendizaje por refuerzo, debido a su aplicabilidad en el desarrollo de modelos avanzados para el control, es fundamental reconocer el papel complementario del aprendizaje supervisado en la resolución de ciertos desafíos específicos, particularmente en la percepción automática y la identificación de estados dentro de un entorno dinámico. El capítulo comienza explicando que el aprendizaje supervisado es un paradigma fundamental dentro del aprendizaje automático. En este enfoque, los modelos aprenden

a partir de datos etiquetados—entendidos como señales de supervisión, generalmente anotadas por humanos—para realizar predicciones precisas. El objetivo central del aprendizaje supervisado es la construcción de un modelo predictivo que pueda generalizar adecuadamente a instancias no vistas, identificando patrones relevantes a partir de ejemplos etiquetados. En términos formales, la meta del aprendizaje supervisado es desarrollar un modelo estadístico capaz de aproximar con precisión una función de mapeo que asocie las características de entrada con sus etiquetas de salida correspondientes. Dentro de este paradigma, los modelos basados en redes neuronales desempeñan un papel fundamental debido a su capacidad para aproximar funciones altamente generalizables, es decir, funciones que pueden desempeñarse eficazmente en datos distintos de los utilizados durante el entrenamiento. Por esta razón, la parte de esta tesis dedicada al aprendizaje supervisado se centra específicamente en el contexto de las redes neuronales. En la Sección 5.2, se introduce el Aprendizaje en Conjunto (Ensemble Learning) como el eje central de este trabajo en aprendizaje supervisado. Inspirado en el principio de la "Sabiduría de las Multitudes" (Wisdom of Crowds), el aprendizaje en conjunto busca construir modelos robustos mediante la combinación de múltiples modelos base (Base Learners, BLs) débiles pero diversos, con el objetivo de gestionar el equilibrio entre sesgo y varianza. Posteriormente, en las Secciones 5.2.1, 5.2.2 y 5.2.3, se profundiza en las estrategias fundamentales del aprendizaje en conjunto: "Bagging" (Bootstrap Aggregating), "Boosting" y "Majority Voting", analizando sus principios, ventajas y aplicaciones en distintos contextos de aprendizaje automático

## Capítulo 6

Este capítulo presenta un novedoso algoritmo de boosting que integra principios de generación de columnas, una técnica bien establecida en optimización clásica, con el aprendizaje supervisado basado en redes neuronales. El enfoque propuesto, denominado algoritmo de aprendizaje en conjunto con reducción de datos (DDA, por sus siglas en inglés), optimiza un conjunto de modelos resolviendo un problema de programación lineal (LP) definido por las predicciones en el conjunto de entrenamiento. El algoritmo construye una combinación convexa óptima a partir de un conjunto de modelos base, aprovechando el concepto de margen suave, generalmente utilizado en máquinas de soporte vectorial. A través de esta optimización, el conjunto busca maximizar el rendimiento predictivo en los datos de entrenamiento mientras mantiene su capacidad de generalización. El conjunto de modelos base se amplía de manera iterativa, y la generación de nuevos modelos base está guiada por el desempeño del conjunto en puntos de

datos específicos del conjunto de entrenamiento. Para ello, se emplean aproximadores basados en redes neuronales y técnicas de aprendizaje supervisado con el fin de ajustar nuevos modelos base diversos que corrijan los errores previos del conjunto. Este trabajo busca tender un puente metodológico entre las técnicas tradicionales de optimización y los marcos modernos de aprendizaje basados en redes neuronales. La Sección 6.2 introduce LPBoost [123] y la Generación de Columnas como los principales fundamentos del algoritmo propuesto. Se centra en tareas de clasificación binaria de imágenes y en la construcción lineal de un modelo en conjunto. El problema de aprendizaje supervisado se define formalmente en la Sección 6.2.1, donde se proporciona una justificación teórica para un enfoque de generación de columnas basado en conjuntos de redes neuronales. En la Sección 6.2.2 se define el problema de combinación convexa lineal y se introduce una función de pérdida sustituta, el margen suave, para abordar el error de generalización. Con base en este problema, se formula un problema dual para determinar si un nuevo modelo base contribuye a aumentar el margen suave. La Sección 6.2.3 introduce el concepto de dispersión en un problema de programación lineal y reporta que la solución del problema maestro dual es típicamente un vector disperso. Esta observación resulta clave para implementar el algoritmo DDA de manera eficiente. Un enfoque interesante para estimar el rendimiento del modelo en conjunto es analizar la diversidad de sus modelos base en combinación con su precisión. La Sección 6.2.4 define la métrica de diversidad y la monitorea como un indicador adicional de la capacidad de generalización del conjunto. Basándose en la teoría establecida, la Sección 6.3 describe el algoritmo DDA, el cual se define mediante dos fases: una fase de generación basada en el principio de Bagging y una fase de refinamiento que opera sobre los principios del algoritmo LPBoost. La fase de generación proporciona las columnas iniciales necesarias para la formulación del programa lineal. Bajo supuestos relativamente fuertes, se demuestra que el algoritmo converge hacia el margen óptimo suave. El capítulo continúa con la presentación de experimentos numéricos en la Sección 6.4. Se evalúa la mejora en la precisión tras la ejecución del algoritmo y se presentan resultados intermedios de la fase de generación para analizar la contribución de cada una de las fases. Finalmente, el capítulo concluye con la observación de que, si bien DDA genera modelos en conjunto que superan en precisión a los modelos base individuales, su costo computacional no justifica las leves mejoras de precisión obtenidas. DDA permanece como un ejercicio académico que arroja luz sobre la intersección entre el aprendizaje automático y la optimización clásica y, aunque no resulta útil en la práctica, explora conceptos teóricos relevantes para el análisis y la comprensión de los principios de funcionamiento de los métodos de aprendizaje supervisado.

## Capítulo 7

Este capítulo presenta una técnica innovadora para la generación de conjuntos no lineales que aprovecha las propiedades de las redes neuronales. Concretamente, el enfoque propuesto explota el hecho de que las redes neuronales, como funciones compuestas, transforman los datos de entrada a través de múltiples capas de abstracción.

En tareas de clasificación, la última capa de la red codifica la distribución sobre las predicciones de clases, representando el nivel más alto de abstracción. Sin embargo, las activaciones de las capas ocultas —frecuentemente denominadas variables latentes— encapsulan representaciones ricas en niveles inferiores de abstracción. En este contexto, este capítulo introduce un enfoque innovador que construye conjuntos no a partir de la capa de salida final, sino a partir de estas variables latentes.

La metodología propuesta parte de la premisa de que las variables latentes, cuando se consideran en su conjunto dentro de una misma capa, proporcionan una representación significativamente más amplia de la distribución de los datos. Así, en lugar de fusionar las predicciones finales de múltiples modelos, la técnica construye el conjunto directamente sobre estas representaciones intermedias. Para potenciar aún más la capacidad de generalización del conjunto, se incorpora un discriminador entrenado sobre las variables latentes.

Además, con el fin de mejorar la precisión del conjunto, la función de pérdida del entrenamiento se complementa con un término adicional de diversidad. Este término cuantifica explícitamente la diferencia entre los parámetros de cada aprendiz base y los parámetros agregados del conjunto. El objetivo de este mecanismo es doble: por un lado, incentivar que los aprendices base realicen predicciones precisas y, por otro, garantizar que sus configuraciones de parámetros sean lo suficientemente distintas entre sí. En términos ideales, cada aprendiz base debería ocupar una dirección ortogonal en el espacio funcional, maximizando así la independencia y minimizando la redundancia entre los modelos que conforman el conjunto.

Después de una breve introducción al problema del aprendizaje supervisado, la Sección 7.2 formaliza los conceptos relevantes. Se definen las representaciones latentes de los datos de entrada, así como los módulos de clasificación y abstracción. Se introduce el concepto de transformador, cuyo espacio imagen se denomina espacio latente. Dentro de esta arquitectura, un clasificador toma como entrada las transformaciones generadas por los aprendices base y, en conjunto, configuran la nueva estructura del conjunto. Adicionalmente, la función de pérdida se enriquece con un término que mide la distancia en norma  $L2$  entre los parámetros promedio de los aprendices base en la

iteración previa y los parámetros del aprendizaje base en la iteración actual del boosting.

Las conclusiones de este capítulo son similares a las del capítulo anterior. Aunque la leve mejora en la precisión del conjunto no justifica el costo computacional adicional del enfoque, los resultados empíricos obtenidos ofrecen hallazgos interesantes respecto a la diversidad de representaciones y las estrategias de conjunto basadas en grafos. Estos aspectos proporcionan nuevas perspectivas para el desarrollo de métodos más eficientes y diversificados en el aprendizaje supervisado.

## Capítulo 8

Este capítulo concluye la tesis. Se reitera en él los objetivos fundamentales del trabajo, se describen las metodologías empleadas en el desarrollo de la investigación y se presentan los hallazgos más relevantes en los ámbitos del aprendizaje por refuerzo (RL) y del aprendizaje supervisado (SL). La Sección 8.2 expone diversas líneas de investigación prometedoras que pueden ser exploradas para avanzar en el campo del aprendizaje por refuerzo, tanto en su forma convencional como en su variante distribucional (DistRL), y, de manera más general, en el aprendizaje automático. Se postula que el estudio de los aspectos fundamentales del aprendizaje automático no solo impulsa el progreso en este campo, sino que también amplía nuestra comprensión de la inteligencia humana. Entre los objetivos a corto y mediano plazo, se plantea la mejora algorítmica del método C-DSAC. En particular, se considera la integración de heurísticas comunes en RL descritas en el capítulo introductorio, así como la incorporación de la técnica de auto-entropía, cuya omisión en los experimentos realizados tuvo el propósito de garantizar una comparación justa con el algoritmo SAC. Asimismo, se espera que la representación más sofisticada de distribuciones a través de Modelos de Mezcla Gaussiana (GMM) permita un incremento significativo en el rendimiento de C-DSAC. Desde una perspectiva más amplia, el aprendizaje por refuerzo está llamado a desempeñar un papel central en la búsqueda de la Inteligencia Artificial General (AGI). En este contexto, se menciona el "Plan Alberta" de Richard Sutton, una estrategia de investigación que enfatiza la necesidad de desarrollar sistemas de RL continuos que operen en entornos del mundo real, teniendo en cuenta las limitaciones computacionales y las interacciones multiagente. Este plan introduce el concepto de un agente generalizado que amplía las arquitecturas tradicionales de RL al integrar módulos adicionales de percepción, transición, política y crítica, cada uno de los cuales optimiza objetivos distintos. Otra línea de investigación de gran interés es la integración del aprendizaje por refuerzo con los modelos fundacionales, tales como los grandes modelos de lenguaje (LLMs). Los mod-

elos fundacionales encapsulan vastos volúmenes de conocimiento sobre el mundo desde una perspectiva humana. Dicho conocimiento puede ser aprovechado para implementar mecanismos de arranque eficiente en problemas de toma de decisiones. A su vez, los algoritmos de RL están siendo empleados con creciente frecuencia en la adaptación y ajuste fino de estos modelos. Un ejemplo reciente de ello es la aplicación de un ciclo de RL puro en el modelo DeepSeek R1 para optimizar su comportamiento. Más allá de los avances algorítmicos, se argumenta que un cambio de paradigma en la arquitectura de los aproximadores de funciones podría impulsar el próximo gran salto en el rendimiento de los modelos. En particular, este capítulo enfatiza la importancia de la aceleración mediante hardware innovador, posiblemente inspirado en la biología neuronal o basado en principios cuánticos. Finalmente, la tesis concluye con una afirmación ambiciosa: se predice que el aprendizaje automático no solo impulsará el progreso tecnológico, sino que también contribuirá a una comprensión más profunda de la naturaleza fundamental de la inteligencia. Dadas las sorprendentes similitudes entre los sistemas de aprendizaje biológicos y artificiales—tanto en su capacidad de interactuar con el entorno como en su habilidad para formar abstracciones—se sostiene que los desarrollos futuros en este campo jugarán un papel clave en la exploración de los mecanismos subyacentes a la cognición, la creatividad y la inteligencia misma.



## Bibliography

- [1] Oxford English Dictionary. Oxford english dictionary, 2023.
- [2] Sarker Iqbal H. Machine learning: Algorithms, real-world applications and research directions. *SN Computer Science*, 2021.
- [3] Aureilien Gelron. *Hands-on machine learning with Scikit-Learn, Keras and TensorFlow: concepts, tools, and techniques to build intelligent systems*. O'Reilly, California, 2 edition, 2019.
- [4] Vladimir Nasteski. An overview of the supervised machine learning methods. *Horizons*, 4:51–62, 2017.
- [5] Richard S. Sutton and Andrew G. Barto. *Reinforcement Learning: An Introduction*. The MIT Press, second edition, 2018.
- [6] Hebb Donald O. *The organization of behavior: A neuropsychological theory*. Wiley, New York, 1 edition, 1949.
- [7] Warren Mcculloch and Walter Pitts. A logical calculus of ideas immanent in nervous activity. *Bulletin of Mathematical Biophysics*, 5:127–147, 1943.
- [8] F. Rosenblatt. The perceptron - a perceiving and recognizing automaton. Technical Report 85-460-1, Cornell Aeronautical Laboratory, Ithaca, New York, January 1957.
- [9] Marvin Minsky and Seymour A. Papert. *Perceptrons: An Introduction to Computational Geometry*. The MIT Press, 1969.
- [10] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Alex Graves, Ioannis Antonoglou, Daan Wierstra, and Martin Riedmiller. Playing atari with deep reinforcement learning, 2013.
- [11] Yongshuai Liu, Avishai Halev, and Xin Liu. Policy learning with constraints in model-free reinforcement learning: A survey. In Zhi-Hua Zhou, editor, *Proceedings of the Thirtieth International Joint Conference on Artificial Intelligence, IJCAI-21*, pages 4508–4515. International Joint Conferences on Artificial Intelligence Organization, 8 2021. Survey Track.
- [12] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017.

## BIBLIOGRAPHY

---

- [13] T.P. Nguyen, S. Choi, and S.J. et al. Park. Inspecting method for defective casting products with convolutional neural network (cmn). *International Journal of Precision Engineering and Manufacturing-Green Technology*, 8:583–594, 2021.
- [14] A. Azadeh, S.F. Ghaderi, and S. Sohrabkhani. Annual electricity consumption forecasting by neural network in high energy consuming industrial sectors. *Energy Conversion and Management*, 49(8):2272–2278, 2008.
- [15] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale, 2021.
- [16] OpenAI. Dall-e: Creating images from text. Blog post, 2021. <https://openai.com/blog/dall-e/>, Accessed: 2022-02-10.
- [17] Sercan Arik, Jitong Chen, Kainan Peng, Wei Ping, and Yanqi Zhou. Neural voice cloning with a few samples. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 31. Curran Associates, Inc., 2018.
- [18] Rainer Drath and Alexander Horch. Industrie 4.0: Hit or hype?[industry forum]. *IEEE industrial electronics magazine*, 8(2):56–58, 2014.
- [19] Klaus Schwab. *The fourth industrial revolution*. Crown Currency, 2017.
- [20] Marc G. Bellemare, Will Dabney, and Rémi Munos. A distributional perspective on reinforcement learning. In *Proceedings of the 34th International Conference on Machine Learning - Volume 70*, ICML'17, page 449–458. JMLR.org, 2017.
- [21] Alhussein Fawzi, Matej Balog, Aja Huang, Thomas Hubert, Bernardino Romera-Paredes, Mohammadamin Barekatin, Alexander Novikov, Francisco J. R. Ruiz, Julian Schrittwieser, Grzegorz Swirszcz, David Silver, Demis Hassabis, and Pushmeet Kohli. Discovering faster matrix multiplication algorithms with reinforcement learning. *Nature*, 610:47 – 53, 2022.
- [22] Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano, Jan Leike, and Ryan Lowe. Training language models to follow instructions with human feedback, 2022.
- [23] Shixiang Shane Gu, Ethan Holly, Timothy P. Lillicrap, and Sergey Levine. Deep reinforcement learning for robotic manipulation with asynchronous off-policy updates. In *Proceedings of the IEEE International Conference on Robotics and Automation (ICRA)*, pages 3389–3396, Piscataway, NJ, USA, May 2017. IEEE.
- [24] Cosmin Paduraru, Daniel Jaymin Mankowitz, Gabriel Dulac-Arnold, Jerry Li, Nir Levine, Sven Goyal, and Todd Hester. Challenges of real-world reinforcement learning: definitions, benchmarks and analysis. *Machine Learning*, 110:2419 – 2468, 2021.
- [25] Hado van Hasselt, Arthur Guez, and David Silver. Deep reinforcement learning with double Q-learning. In *Proceedings of the Thirtieth AAAI Conference on Artificial Intelligence*, AAAI'16, page 2094–2100. AAAI Press, 2016.
- [26] Timothy P. Lillicrap, Jonathan J. Hunt, Alexander Pritzel, Nicolas Heess, Tom Erez, Yuval Tassa, David Silver, and Daan Wierstra. Continuous control with deep reinforcement learning, 2019.

## BIBLIOGRAPHY

---

- [27] John Schulman, Sergey Levine, Philipp Moritz, Michael I. Jordan, and Pieter Abbeel. Trust region policy optimization, 2017.
- [28] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms, 2017.
- [29] Scott Fujimoto, Herke van Hoof, and David Meger. Addressing function approximation error in actor-critic methods. In Jennifer Dy and Andreas Krause, editors, *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pages 1587–1596. PMLR, 2018.
- [30] Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. In Jennifer G. Dy and Andreas Krause, editors, *Proceedings of the 35th International Conference on Machine Learning, ICML 2018, Stockholm, Sweden, July 10-15, 2018*, volume 80 of *Proceedings of Machine Learning Research*, pages 1856–1865. PMLR, 2018.
- [31] Marc G. Bellemare, Will Dabney, and Mark Rowland. *Distributional Reinforcement Learning*. MIT Press, 2023. <http://www.distributional-rl.org>, Accessed: 2023-09-15.
- [32] Timothy H. Muller, James L. Butler, Sebastijan Veselic, Bruno Miranda, Timothy Edward John Behrens, Zeb Kurth-Nelson, and Steven Wayne Kennerley. Distributional reinforcement learning in prefrontal cortex. *Nature Neuroscience*, 27:403 – 408, 2021.
- [33] Martin L. Puterman. *Markov Decision Processes: Discrete Stochastic Dynamic Programming*. John Wiley & Sons, Inc., New York, NY, USA, 1st edition, 1994.
- [34] Bob Givan and Ron Parr. An introduction to markov decision processes.
- [35] Venkat N. Gudivada, Dhana Rao, and Vijay V. Raghavan. Chapter 9 - big data driven natural language processing research and applications. In Venu Govindaraju, Vijay V. Raghavan, and C.R. Rao, editors, *Big Data Analytics*, volume 33 of *Handbook of Statistics*, pages 203–238. Elsevier, 2015.
- [36] Theresa Eimer, Marius Lindauer, and Roberta Raileanu. Hyperparameters in reinforcement learning and how to tune them. In Andreas Krause, Emma Brunskill, Kyunghyun Cho, Barbara Engelhardt, Sivan Sabato, and Jonathan Scarlett, editors, *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pages 9104–9149. PMLR, 23–29 Jul 2023.
- [37] Gerald Tesauro. Td-gammon: A self-teaching backgammon program. In Alan F. Murray, editor, *Applications of Neural Networks*, pages 267–285. Springer US, Boston, MA, 1995.
- [38] Ivo Grondman, Lucian Busoniu, Gabriel A. D. Lopes, and Robert Babuska. A survey of actor-critic reinforcement learning: Standard and natural policy gradients. *IEEE Transactions on Systems, Man, and Cybernetics, Part C (Applications and Reviews)*, 42(6):1291–1307, 2012.
- [39] Dimitri P. Bertsekas and John N. Tsitsiklis. Gradient convergence in gradient methods with errors. *SIAM J. Optim.*, 10:627–642, 1999.
- [40] Thomas Degris, Martha White, and Richard S. Sutton. Off-policy actor-critic, 2013.
- [41] John Schulman, Philipp Moritz, Sergey Levine, Michael Jordan, and Pieter Abbeel. High-dimensional continuous control using generalized advantage estimation, 2018.
- [42] Evan Greensmith, Peter L. Bartlett, and Jonathan Baxter. Variance reduction techniques for gradient estimates in reinforcement learning. *Journal of machine learning research*, 5:1471–1530, 2004.

## BIBLIOGRAPHY

---

- [43] Huy Xuan Pham, Hung Manh La, David Feil-Seifer, and Aria Nefian. Cooperative and distributed reinforcement learning of drones for field coverage, 2018.
- [44] Martin A. Riedmiller, Andrew W. Moore, and Jeff G. Schneider. Reinforcement learning for cooperating and communicating reactive agents in electrical power grids. In Markus Hannebauer, Jan Wendler, and Enrico Pagello, editors, *Balancing Reactivity and Social Deliberation in Multi-Agent Systems, From RoboCup to Real-World Applications (selected papers from the ECAI 2000 Workshop and additional contributions)*, volume 2103 of *Lecture Notes in Computer Science*, pages 137–149. Springer, 2000.
- [45] Pawel Ladosz, Lilian Weng, Minwoo Kim, and Hyondong Oh. Exploration in deep reinforcement learning: A survey. *Information Fusion*, 85:1–22, 2022.
- [46] Sebastian Thrun and Anton Schwartz. Issues in using function approximation for reinforcement learning. In Michael Mozer, Paul Smolensky, David Touretzky, Jeffrey Elman, and Andreas Weigend, editors, *Proceedings of the 1993 Connectionist Models Summer School*, pages 255–263. Lawrence Erlbaum, 1993.
- [47] Mohammad Gheshlaghi Azar, Rémi Munos, and Hilbert J Kappen. Minimax pac bounds on the sample complexity of reinforcement learning with a generative model. *Machine Learning*, 91(3):325–349, 2013.
- [48] R. A. Howard. *Dynamic Programming and Markov Processes*. MIT Press, Cambridge, MA, 1960.
- [49] Bruno Scherrer. Improved and generalized upper bounds on the complexity of policy iteration. In C.J. Burges, L. Bottou, M. Welling, Z. Ghahramani, and K.Q. Weinberger, editors, *Advances in Neural Information Processing Systems*, volume 26. Curran Associates, Inc., 2013.
- [50] Yashraj Narang, Kier Storey, Iretoiyo Akinola, Miles Macklin, Philipp Reist, Lukasz Wawrzyniak, Yunrong Guo, Adam Moravanszky, Gavriel State, Michelle Lu, Ankur Handa, and Dieter Fox. Factory: Fast contact for robotic assembly, 2022.
- [51] Josh Tobin, Rachel Fong, Alex Ray, Jonas Schneider, Wojciech Zaremba, and Pieter Abbeel. Domain randomization for transferring deep neural networks from simulation to the real world, 2017.
- [52] Arnab Nilim and Laurent El Ghaoui. Robust control of markov decision processes with uncertain transition matrices. *Operations Research*, 53(5):780–798, 2005.
- [53] Richard S. Sutton, Michael Bowling, and Patrick M. Pilarski. The Alberta plan for AI research, 2023.
- [54] Aravind Singh, Abhishek Gupta, Aravind Rajeswaran, Vikash Kumar, and Sergey Levine. End-to-end robotic reinforcement learning without reward engineering. In *Robotics: Science and Systems*, 2019.
- [55] Nuttapon Chentanez, Andrew Barto, and Satinder Singh. Intrinsically motivated reinforcement learning. In L. Saul, Y. Weiss, and L. Bottou, editors, *Advances in Neural Information Processing Systems*, volume 17. MIT Press, 2004.
- [56] Jason Gabriel and Vafa Ghazavi. 336the challenge of value alignment: From fairer algorithms to AI safety. In *Oxford Handbook of Digital Ethics*. Oxford University Press, 12 2023.
- [57] Thomas Arnold, Daniel Kasenberg, and Matthias Scheutz. Value alignment or misalignment—what will keep systems accountable? In *Workshops at the thirty-first AAAI conference on artificial intelligence*, 2017.
- [58] Jan Leike and Ilya Sutskever. Introducing superalignment, 2023. Accessed: 2024-10-28.
- [59] Nazri Mohd Nawi, Walid Hasen Atomi, and M.Z. Rehman. The effect of data pre-processing on optimized training of artificial neural networks. *Procedia Technology*, 11:32–39, 2013. 4th International Conference on Electrical Engineering and Informatics, ICEEI 2013.

- 
- [60] Ping Luo, Xinjiang Wang, Wenqi Shao, and Zhanglin Peng. Towards understanding regularization in batch normalization, 2019.
- [61] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift, 2015.
- [62] Tom Schaul, John Quan, Ioannis Antonoglou, and David Silver. Prioritized experience replay, 2016.
- [63] Xavier Glorot and Yoshua Bengio. Understanding the difficulty of training deep feedforward neural networks. In Yee Whye Teh and Mike Titterton, editors, *Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics*, volume 9 of *Proceedings of Machine Learning Research*, pages 249–256, Chia Laguna Resort, Sardinia, Italy, 13–15 May 2010. PMLR.
- [64] Eyal Even-Dar and Y. Mansour. Learning rates for q-learning. *J. Mach. Learn. Res.*, 5:1–25, 2004.
- [65] Diederik Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *International Conference on Learning Representations (ICLR)*, San Diego, CA, USA, 2015.
- [66] Christopher Watkins and Peter Dayan. Q-learning. *Machine Learning*, 8(3):279–292, 1992.
- [67] M. G. Bellemare, Y. Naddaf, J. Veness, and M. Bowling. The arcade learning environment: An evaluation platform for general agents. *Journal of Artificial Intelligence Research*, 47:253–279, June 2013.
- [68] Ronald J. Williams. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine Learning*, 8:229–256, 2004.
- [69] David Silver. Lecture 7: Policy gradient methods, 2015. Slides for lecture on reinforcement learning, Accessed: 2023-06-05.
- [70] David Silver, Guy Lever, Nicolas Manfred Otto Heess, Thomas Degris, Daan Wierstra, and Martin A. Riedmiller. Deterministic policy gradient algorithms. In *International Conference on Machine Learning*, 2014.
- [71] Antonin Raffin, Ashley Hill, Adam Gleave, Anssi Kanervisto, Maximilian Ernestus, and Noah Dormann. Stable-baselines3: Reliable reinforcement learning implementations. *Journal of Machine Learning Research*, 22(268):1–8, 2021.
- [72] Diederik P Kingma and Max Welling. Auto-Encoding Variational Bayes, 2022.
- [73] Kaiwen Wang, Kevin Zhou, Runzhe Wu, Nathan Kallus, and Wen Sun. The benefits of being distributional: Small-loss bounds for reinforcement learning, 2023.
- [74] Kaiwen Wang, Owen Oertell, Alekh Agarwal, Nathan Kallus, and Wen Sun. More benefits of being distributional: Second-order bounds for reinforcement learning. In *Proceedings of the 41st International Conference on Machine Learning*, volume 235 of *Proceedings of Machine Learning Research*, pages 51192–51213. PMLR, 2024.
- [75] Marc G. Bellemare, Ivo Danihelka, Will Dabney, Shakir Mohamed, Balaji Lakshminarayanan, Stephan Hoyer, and Rémi Munos. The Cramer distance as a solution to biased Wasserstein gradients, 2017.
- [76] Will Dabney, Mark Rowland, Marc G. Bellemare, and Rémi Munos. Distributional reinforcement learning with quantile regression. In *AAAI Conference on Artificial Intelligence*, 2017.
- [77] Will Dabney, Georg Ostrovski, David Silver, and Rémi Munos. Implicit quantile networks for distributional reinforcement learning, 2018.

## BIBLIOGRAPHY

---

- [78] Mark Rowland, Marc G. Bellemare, Will Dabney, Rémi Munos, and Yee Whye Teh. An analysis of categorical distributional reinforcement learning. In *International Conference on Artificial Intelligence and Statistics*, 2018.
- [79] Karl Pearson. Contributions to the mathematical theory of evolution. *Philosophical Transactions of the Royal Society of London*, 185:71–110, 1894.
- [80] Karl Cobbe, Oleg Klimov, Chris Hesse, Taehoon Kim, and John Schulman. Quantifying generalization in reinforcement learning. In Kamalika Chaudhuri and Ruslan Salakhutdinov, editors, *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pages 1282–1289. PMLR, 09–15 Jun 2019.
- [81] Aviral Kumar, Rishabh Agarwal, Tengyu Ma, Aaron C. Courville, G. Tucker, and Sergey Levine. Dr3: Value-based deep reinforcement learning requires explicit regularization. *ICML*, abs/2112.04716, 2021.
- [82] Jingliang Duan, Yang Guan, Shengbo Eben Li, Yangang Ren, Qi Sun, and BO CHENG. Distributional soft actor-critic: Off-policy reinforcement learning for addressing value estimation errors. *IEEE Transactions on Neural Networks and Learning Systems*, 33(11):6584–6598, nov 2022.
- [83] Xiaoteng Ma, Li Xia, Zhengyuan Zhou, Jun Yang, and Qianchuan Zhao. DSAC: Distributional soft actor critic for risk-sensitive reinforcement learning, 2020.
- [84] Derek Yang, Li Zhao, Zichuan Lin, Tao Qin, Jiang Bian, and Tieyan Liu. *Fully parameterized quantile function for distributional reinforcement learning*. Curran Associates Inc., Red Hook, NY, USA, 2019.
- [85] Alix Lhéritier and Nicolas Bondoux. A Cramér distance perspective on quantile regression based distributional reinforcement learning, 2022.
- [86] Daniel Wontae Nam, Younghoon Kim, and Chan Y. Park. GMAC: A distributional perspective on actor-critic framework. In Marina Meila and Tong Zhang, editors, *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pages 7930–7939. PMLR, 18–24 Jul 2021.
- [87] V M Zolotarev. Metric distances in spaces of random variables and their distributions. *Mathematics of the USSR-Sbornik*, 30(3):373, apr 1976.
- [88] Dino Sejdinovic, Bharath Sriperumbudur, Arthur Gretton, and Kenji Fukumizu. Equivalence of distance-based and rkhs-based statistics in hypothesis testing. *The Annals of Statistics*, 41(5):2263–2291, 2013.
- [89] Qingfeng Lan, Yangchen Pan, Alona Fyshe, and Martha White. Maxmin q-learning: Controlling the estimation bias of q-learning. In *International Conference on Learning Representations*, 2020.
- [90] Logan Engstrom, Andrew Ilyas, Shibani Santurkar, Dimitris Tsipras, Firdaus Janoos, Larry Rudolph, and Aleksander Madry. Implementation matters in deep policy gradients: A case study on PPO and TRPO, 2020.
- [91] Greg Brockman, Vicki Cheung, Ludwig Pettersson, Jonas Schneider, John Schulman, Jie Tang, and Wojciech Zaremba. Openai gym. *arXiv preprint arXiv:1606.01540*, 2016.
- [92] Peter Henderson, Riashat Islam, Philip Bachman, Joelle Pineau, Doina Precup, and David Meger. Deep reinforcement learning that matters. In *AAAI Conference on Artificial Intelligence*, 2017.
- [93] Vanya Aziz. Robotersteuerung mit modellgestütztem Reinforcement Learning. Master’s thesis, Technische Hochschule Köln, Germany, 2021.
- [94] International Federation of Robotics. Standardisation: Definition of a robot, 2024. Accessed: 2024-12-12.

## BIBLIOGRAPHY

---

- [95] Sami Haddadin, Sven Parusel, Lars Johannsmeier, Saskia Golz, Simon Gabl, Florian Walch, Mohamadreza Sabaghian, Christoph Jähne, Lukas Hausperger, and Simon Haddadin. The franka emika robot: A reference platform for robotics research and education. *IEEE Robotics & Automation Magazine*, 29(2):46–64, 2022.
- [96] Francesco Cursi, Weibang Bai, Eric Yeatman, and Petar Kormushev. Globdesopt: A global optimization framework for optimal robot manipulator design. *IEEE Access*, PP:1–1, 01 2022.
- [97] Kenneth Levenberg. A method for the solution of certain non – linear problems in least squares. *Quarterly of Applied Mathematics*, 2:164–168, 1944.
- [98] Donald W. Marquardt. An algorithm for least-squares estimation of nonlinear parameters. *Journal of the Society for Industrial and Applied Mathematics*, 11(2):431–441, 1963.
- [99] Kungurtsev Bergou, Diouane. Convergence and complexity analysis of a levenberg–marquardt algorithm for inverse problems. *Journal of Optimization Theory and Applications*, 185:927–944, 2020.
- [100] James Surowiecki. *The Wisdom of Crowds*. Anchor, 2005.
- [101] Lior Rokach. Ensemble-based classifiers. *Artificial intelligence review*, 33:1–39, 2010.
- [102] Balaji Lakshminarayanan, Alexander Pritzel, and Charles Blundell. Simple and scalable predictive uncertainty estimation using deep ensembles. *Advances in Neural Information Processing Systems*, 30, 2017.
- [103] Kohavi Bauer. An empirical comparison of voting classification algorithms: Bagging, boosting, and variants. *Machine Learning*, 36:105–139, 1999.
- [104] Yoav Freund and Robert E. Schapire. A decision-theoretic generalization of on-line learning and an application to boosting. In Paul Vitányi, editor, *Computational Learning Theory*, pages 23–37, Berlin, Heidelberg, 1995. Springer Berlin Heidelberg.
- [105] Robert E. Schapire. Explaining adaboost. In *Empirical Inference*, 2013.
- [106] Jerome H Friedman. Greedy function approximation: a gradient boosting machine. *Annals of statistics*, pages 1189–1232, 2001.
- [107] Marc Pfetsch and Sebastian Pokutta. Ipboost–non-convex boosting via integer programming. In *International Conference on Machine Learning*, pages 7663–7672. PMLR, 2020.
- [108] Faraz Torabi, Garrett Warnell, and Peter Stone. Behavioral cloning from observation. In *Proceedings of the 27th International Joint Conference on Artificial Intelligence, IJCAI’18*, page 4950–4957. AAAI Press, 2018.
- [109] Thomas G. Dietterich. Ensemble methods in machine learning. In *Multiple Classifier Systems*, pages 1–15, Berlin, Heidelberg, 2000. Springer Berlin Heidelberg.
- [110] Omer Sagi and Lior Rokach. Ensemble learning: A survey. *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery*, 8(4):e1249, 2018.
- [111] Zhi-Hua Zhou. *Ensemble Methods: Foundations and Algorithms*. CRC Press, 2012.
- [112] Andreas Aakerberg, Kamal Nasrollahi, and Thomas Heder. Improving a deep learning based rgb-d object recognition model by ensemble learning. In *2017 Seventh International Conference on Image Processing Theory, Tools and Applications (IPTA)*, pages 1–6. IEEE, 2017.

## BIBLIOGRAPHY

---

- [113] Yen-Yu Lin, Tyng-Luh Liu, and Chiou-Shann Fuh. Local ensemble kernel learning for object category recognition. In *2007 IEEE Conference on Computer Vision and Pattern Recognition*, pages 1–8. IEEE, 2007.
- [114] Tomasz Malisiewicz, Abhinav Gupta, and Alexei A Efros. Ensemble of exemplar-svms for object detection and beyond. In *2011 International Conference on Computer Vision*, pages 89–96. IEEE, 2011.
- [115] Jie Xu, Wei Wang, Hanyuan Wang, and Jinhong Guo. Multi-model ensemble with rich spatial information for object detection. *Pattern Recognition*, 99:107098, 2020.
- [116] Hao Jiang, Xiaojie Qiu, Jing Chen, Xinyu Liu, Xiren Miao, and Shengbin Zhuang. Insulator fault detection in aerial images based on ensemble learning with multi-level perception. *IEEE Access*, 7:61797–61810, 2019.
- [117] Yibin Li, Yan Song, Lei Jia, Shengyao Gao, Qiqiang Li, and Meikang Qiu. Intelligent fault diagnosis by fusing domain adversarial training and maximum mean discrepancy via ensemble learning. *IEEE Transactions on Industrial Informatics*, 17(4):2833–2841, 2020.
- [118] Yueping Wang, Kun Zhu, Mengyun Sun, and Yueyu Deng. An ensemble learning approach for fault diagnosis in self-organizing heterogeneous networks. *IEEE Access*, 7:125662–125675, 2019.
- [119] Abolfazl Abdollahi, Biswajeet Pradhan, and Abdullah M Alamri. An ensemble architecture of deep convolutional segnet and unet networks for building semantic segmentation from high-resolution aerial images. *Geocarto International*, pages 1–16, 2020.
- [120] Mohd Hanafi Ahmad Hijazi, Stefanus Kieu Tao Hwa, Abdullah Bade, Razali Yaakob, and Mohammad Saf-free Jeffree. Ensemble deep learning for tuberculosis detection using chest x-ray and canny edge detected images. *IAES International Journal of Artificial Intelligence*, 8(4):429, 2019.
- [121] Alexander Thebelt, Jan Kronqvist, Miten Mistry, Robert M Lee, Nathan Sudermann-Merx, and Ruth Misener. Entmoot: A framework for optimization over ensemble tree models. *Computers & Chemical Engineering*, 151:107343, 2021.
- [122] Jan Kocoń, Igor Cichecki, Oliwier Kaszyca, Mateusz Kochanek, Dominika Szydło, Joanna Baran, Julita Bielaniec, Marcin Gruza, Arkadiusz Janz, Kamil Kanclerz, Anna Kocoń, Bartłomiej Koptyra, Wiktoria Mieleśczenko-Kowszewicz, Piotr Milkowski, Marcin Oleksy, Maciej Piasecki, Łukasz Radliński, Konrad Wojtasik, Stanisław Woźniak, and Przemysław Kazienko. Chatgpt: Jack of all trades, master of none. *Information Fusion*, 99:101861, 2023.
- [123] Ayhan Demiriz, Kristin P Bennett, and John Shawe-Taylor. Linear programming boosting via column generation. *Machine Learning*, 46(1):225–254, 2002.
- [124] Pavlo Muts, Stefan Bruche, Ivo Nowak, Ouyang Wu, Eligius MT Hendrix, and George Tsatsaronis. A column generation algorithm for solving energy system planning problems. *Optimization and Engineering*, 24:317–351, 2023.
- [125] Yijun Bian and Huanhuan Chen. When does diversity help generalization in classification ensembles? *IEEE Transactions on Cybernetics*, 52(9):9059–9075, 2022.
- [126] V. M. Tikhomirov. *Convex Analysis*, pages 1–92. Springer Berlin Heidelberg, 1990.
- [127] Sandra Benítez-Peña, Emilio Carrizosa, Vanesa Guerrero, M. Dolores Jiménez-Gamero, Belén Martín-Barragán, Cristina Molero-Río, Pepa Ramírez-Cobo, Dolores Romero Morales, and M. Remedios Sillero-Denamiel. On sparse ensemble methods: An application to short-term predictions of the evolution of covid-19. *European Journal of Operational Research*, 295(2):648–663, 2021.

## BIBLIOGRAPHY

---

- [128] Göksu Tüysüzöğlü and Dokuz Birant. Enhanced bagging (ebagging): A novel approach for ensemble learning. *International Arab Journal of Information Technology*, 17(4):515–528, 2020.
- [129] D.W. Ruck, S.K. Rogers, M. Kabrisky, and J.P. Mills. Target recognition: Conventional and neural network approaches. In *International 1989 Joint Conference on Neural Networks*, pages 608 vol.2–, 1989.
- [130] K. Fukushima. Neocognitron: A Self-organizing Neural Network Model for a Mechanism of Pattern Recognition Unaffected by Shift in Position. *Biological Cybernetics*, 36(4):193–202, 1980.
- [131] Leo Breiman, Jerome Friedman, Charles J. Stone, and R.A. Olshen. *Classification and Regression Trees*. Chapman and Hall/CRC, 1984.
- [132] Alex Krizhevsky. Learning multiple layers of features from tiny images. *University of Toronto*, 05 2012.
- [133] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition, 2015.
- [134] Vanya Aziz, Ivo Nowak, Ouyang Wu, Eligius M. T. Hendrix, and Jan Kronqvist. On optimizing ensemble models using column generation. *J Optim Theory Appl*, 203(1):1794–1819, 2024.
- [135] Alexander Kolesnikov, Lucas Beyer, Xiaohua Zhai, Joan Puigcerver, Jessica Yung, Sylvain Gelly, and Neil Houlsby. Big transfer (bit): General visual representation learning. In Andrea Vedaldi, Horst Bischof, Thomas Brox, and Jan-Michael Frahm, editors, *Computer Vision – ECCV 2020*, pages 491–507, Cham, 2020. Springer International Publishing.
- [136] Linnan Wang, Saining Xie, Teng Li, Rodrigo Fonseca, and Yuandong Tian. Sample-efficient neural architecture search by learning actions for monte carlo tree search. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 44(9):5503–5515, 2022.
- [137] Bruno Antonio, Davide Moroni, and Massimo Martinelli. Efficient adaptive ensembling for image classification. *Expert Systems*, n/a(n/a), 2023.
- [138] Hassam Sheikh, Mariano Phielipp, and Ladislau Boloni. Maximizing ensemble diversity in deep reinforcement learning. In *International Conference on Learning Representations*, 2022.
- [139] Tuomas Haarnoja, Aurick Zhou, Kristian Hartikainen, George Tucker, Sehoon Ha, Jie Tan, Vikash Kumar, Henry Zhu, Abhishek Gupta, Pieter Abbeel, and Sergey Levine. Soft actor-critic algorithms and applications, 2019.
- [140] DeepSeek-AI. Deepseek-R1: Incentivizing Reasoning Capability in LLMs via Reinforcement Learning. Technical report, deepseek, 2025.
- [141] Lili Chen, Kevin Lu, Aravind Rajeswaran, Kimin Lee, Aditya Grover, Michael Laskin, Pieter Abbeel, Aravind Srinivas, and Igor Mordatch. Decision transformer: Reinforcement learning via sequence modeling. In A. Beygelzimer, Y. Dauphin, P. Liang, and J. Wortman Vaughan, editors, *Advances in Neural Information Processing Systems*, 2021.
- [142] Scott Emmons, Benjamin Eysenbach, Ilya Kostrikov, and Sergey Levine. Rvs: What is essential for offline RL via supervised learning? In *International Conference on Learning Representations*, 2022.
- [143] Vanya Aziz, Ivo Nowak, and Eligius M. T. Hendrix. Distributional RL Using the Cramér distance. *J. Mach. Learn. Res.*, submitted, 2024.
- [144] Vanya Aziz, Ivo Nowak, Ouyang Wu, Jan Kronqvist, and Eligius M. T. Hendrix. Boosting via Latent Representation. In *Proceedings OLA’2023 International Conference on Optimization and Learning*, pages 150–153, 2023.