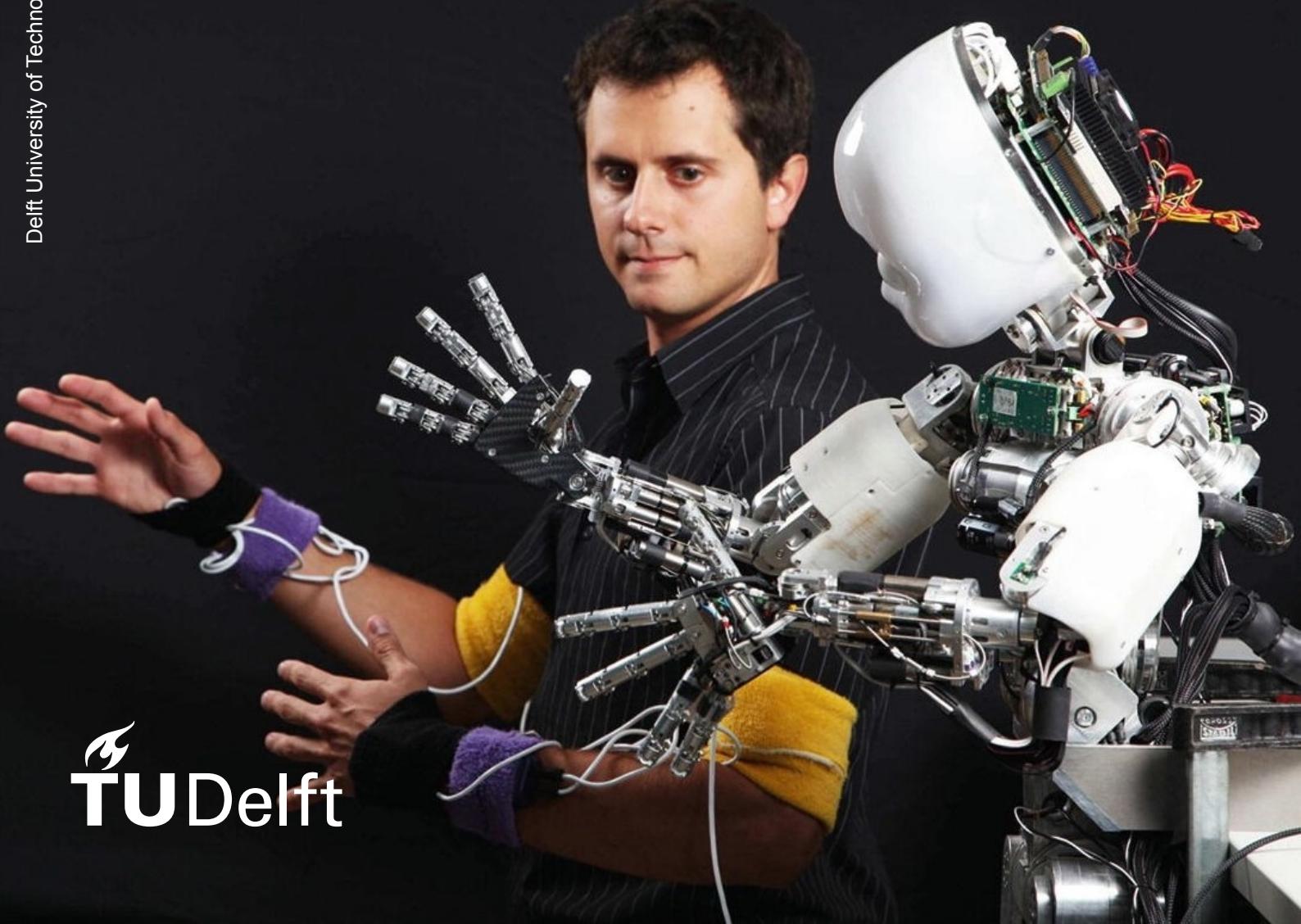


Robotic skill learning

How a novice robot becomes an expert

Literature Study

R.E.S. Maesen



Robotic skill learning

How a novice robot becomes an expert

by

R.E.S. Maesen

Author: R.E.S. Maessen
Student number: 4564200
Supervisor: Dr. L. Peternel
Institution: Delft University of Technology
Study: Msc. Robotics
Date: Wednesday 16th February, 2022

Cover Image: Demonstration of the skill with simultaneous reproduction on the robot [1].



Summary

In the past, the robot's behavior was directly programmed. This has the disability of not being able to adapt to changes in the environment, as each action had to be predefined. Therefore, the topic of robotic learning has increased in popularity. In this report, the idea of robotic learning will be examined. This is done by trying to answer the following research question:

How do methods, found in literature, enable continuous skill learning of robots, and how does the learned skill change during the learning process?

The scope of this research is limited using two aspects. Firstly, only online learning method will be examined, meaning the model representing the behavior of the robot will somehow incrementally change. And secondly, no deep learning approaches will be taken into an account.

To answer the research question, two topics are discussed, each relating to a part of this question. Firstly, in chapter 2 method enabling continuous skill learning will be examined. And secondly, in chapter 3, the changes of the learned skill during the learning process will be discussed.

Two types of skill learning methods were identified: *Imitation Learning (IL)* and *Reinforcement Learning (RL)*. IL learns by imitation, an expert. This is done based on required demonstrations. *Behavioral Cloning (BC)* and *Inverse Reinforcement Learning (IRL)* were found to be the two broad IL categories. The former tries to learn a policy by directly mapping from the state/context to the trajectories/action or the control input. Whereas the latter uses the demonstrated data to compute a reward/cost function which then can be used (using RL) to compute a policy. RL learns by using a trial-and-error approach. By interacting with the environment and using a reward function, the policy can be improved. The two main categories are *value function* and *policy search (PS)*. The value function approach tries to find an optimal policy by iteratively optimizing the value function. In contrast, PS directly learns the policy. Recently, a new type of RL method has been introduced: actor-critic RL, which combines the advantages of both methods by separating the policy (actor) and the value function (critic). Both IL and RL have their disadvantage. The skill obtained using IL is limited by the skill level of the human expert. And RL is often computationally expensive as a lot of iterations are required to obtain a solution. Therefore, a possible solution would be to combine these two methods. While research on this topic is still limited, some examples have been found. One example would be to obtain an initial policy using IL and use this as a starting point of the RL method. A comparison was made based on the different IL and RL approaches (table 2.4). They were compared on computational efficiency, stability, smoothness, generalizability, implementation abilities, ability to deal with high-dimensional input, and whether they were originally designed for online learning.

For continuous learning, multiple learning outcomes were identified. The one further examined is policy. To obtain policy different types of state inputs can be used, e.g. position, joint angle, force, or torque. Position and joint angles work well for static environments, but some problems occur when needing to adapt to the environment (or when they need to generalize). If required, some additional feedback is required. Force and torque are better at dealing with small changes in the environment, as they can adapt to it. However, for large changes also additional feedback is necessary. In addition, research was conducted where a cost/reward function was used. This showed the importance of choosing a learning setting for a specific task. Where one method can work well for certain cases, this does not necessarily have to be the same for other situations. Additionally, the usefulness of combining IL and RL was examined. Showing that combining them would increasingly reduce the computational time when compared to a traditional RL approach. In addition, the same research has shown the relevance of choosing an initial policy for RL, again resulting in faster convergence.

Based on the literature, it was concluded IL and RL allow for continuous skill learning. Here, the IL methods provide a reliable result that is evenly good as that of the expert. Whereas RL is less reliable but has the potential of finding better solutions than IL.

Contents

Summary	i
List of Abbreviations	iii
List of Figures	v
List of Tables	vi
1 Introduction	1
2 Skill Learning Methods	3
2.1 Markov Decision Process	4
2.2 Imitation Learning	4
2.2.1 Problem definition Imitation Learning	5
2.2.2 Behavioural Cloning	6
2.2.3 Inverse Reinforcement Learning.	15
2.3 Reinforcement Learning	18
2.3.1 Problem Definition Reinforcement Learning.	19
2.3.2 Value-function approaches.	20
2.3.3 Policy Search	22
2.3.4 Actor-critic method	24
2.4 Combining Imitation Learning and Reinforcement Learning	25
2.5 Comparison	25
2.6 Summary	28
3 Continuous Skill Learning	29
3.1 Learning Outcome	29
3.2 State input	30
3.2.1 Position	30
3.2.2 Joint angles.	31
3.2.3 Force	33
3.2.4 Torques	36
3.3 Reward or Cost.	37
3.3.1 Regular learning	37
3.3.2 Human Advice	39
3.4 Summary	40
4 Discussion	41
4.1 Limitation of robotic learning methods.	41
4.2 Generalizability	42
5 Conclusion & Future research	43
5.1 Conclusion	43
5.2 Future research.	44
References	49

List of Abbreviations

AC actor-critic.

BC Behavioral Cloning.

COACH Corrective Advice Communicated by Humans.

DMP Dynamic Movement Primitive.

DoF Degrees of Freedom.

DP Dynamic Programming.

DS Dynamic System.

EM Expectation Maximization.

EMG electromyography.

GHMM Growing Hidden Markov Model.

GMM Gaussian Mixture Model.

GMR Gaussian Mixture Regression.

GP Gaussian Process.

GPR Gaussian Process Regression.

HBM Hierarchical Bayesian Model.

HHMM Hierarchy Hidden Markov Model.

HMM Hidden Markov Model.

HSMM Hidden Semi-Markov Model.

i.i.d. independent and identically distributed.

IL Imitation Learning.

ILWR incremental LWR.

Inf.Th. information-theoretic insights.

IRL Inverse Reinforcement Learning.

LfD Learning from (human) Demonstration.

LGP Local GP.

LWPR Locally Weighted Projection Regression.

LWR Locally Weighted Regression.

MC Monte Carlo.

MDP Markov Decision Process.

MMP maximum margin planning.

MP Motion Primitive.

MPC model predictive control.

PbD Programming by Demonstration.

PG policy gradient.

PiH peg-in-hole.

ProMP Probabilistic Movement Primitive.

PS policy search.

RBF Radial Basis Function.

RL Reinforcement Learning.

SOGP Sparse Online Gaussian Processes.

TD Temporal-Difference.

TP task-parametrized.

TP-GMM TP GMM.

TP-HSMM TP HSMM.

List of Figures

2.1	Illustration of the different levels of representation for describing the skill [19].	6
2.2	Position, velocity and acceleration of an executed Dynamic Movement Primitive (DMP) without any external force ($f = 0$) [25].	7
2.3	Example of the learning of DMP using Locally Weighted Regression (LWR)	8
2.4	The trajectory distribution showing the joint position (upper plot) and velocity (lower plot).	10
2.5	The trajectory distribution is encoded using Probabilistic Movement Primitive (ProMP) for multiple goal positions.	11
2.6	Example of trajectory encoding using Gaussian Mixture Model (GMM).	11
2.7	Trajectory encoding and generalization using GMM and Gaussian Mixture Regression (GMR) [36].	12
2.8	Graphical representation of a generative and discriminative model [44].	13
2.9	Graphical Representation of Hidden Markov Model (HMM) [48].	14
2.10	Geometric description of policy iteration [58].	16
2.11	The agent-environment interaction in reinforcement learning as defined by Sutton & Barto [68]	18
2.12	Taxonomy of Reinforcement Learning (RL) methods [70].	20
2.13	CComparison of the backup of three value-function approaches: Monte Carlo (MC), Temporal-Difference (TD) learning, and Dynamic Programming (DP) [72].	21
2.14	Graphical representation of the DP methods [73].	21
2.15	Schematic overview of an actor-critic algorithm [76].	24
3.1	Planar example of learning and control with the task-parametrized Hidden Semi-Markov Model (TP-HSMM) [89].	31
3.2	Planar example [89] with (left) the graphical representation of the transition matrix and duration probabilities for each state, where the three paths are formed according to the three demonstrations. And (right) the executed motion on the remote side starting from a new position, using TP-HSMM and the model predictive control (MPC) controller.	31
3.3	Experimental setup of peg-in-hole (PiH) task [89].	31
3.4	Updating of the GMM using the direct update incremental learning approach [38].	32
3.5	Reproduction of different gesture movements using GMR in a 3D latent space of motion [38].	32
3.6	Learning scheme for a grasping task Pastor et al.	33
3.7	Result of grasping attempts after misplacing a flash-light [83]	34
3.8	Learning scheme of PiH tasks [30].	34
3.9	Results of PiH task with a rotated and translated baseplate [30].	35
3.10	The skill learning of a novice robot [26]	35
3.11	Results of robot-robot collaboration experiment on two-person sawing task [26].	36
3.12	Control scheme of an exoskeleton as presented in [90].	36
3.13	Results of the experiment on elbow exoskeleton using a trained subject [90].	37
3.14	The sum of reward (return) for the squared trajectory following [95].	38
3.15	The sum of reward (return) for the circle trajectory following [95].	38
3.16	Snapshot of ball-in-cup experiment [80].	39
3.17	Convergence curve of <i>ball-in-cup</i> when learning from an initial demonstration [80]	40
3.18	Convergence curve of <i>ball-in-cup</i> when learning from scratch [80]	40

List of Tables

2.1 Comparison of three methods based on which robot can require skills [9]	3
2.2 Advantages and disadvantage of the two main Imitation Learning (IL) approaches based on multiple papers [15–18].	5
2.3 Advantages and disadvantage of model-free and model-based Reinforcement Learning (RL) methods [67]	20
2.4 Comparison of robot teaching methods.	26
3.1 Overview of references using different types of state input.	30

1

Introduction

In the last few years, robotic learning has increased in popularity. Traditionally, robotic skills were directly programmed. This code, controlling the behavior of the robot, would need to be created such that the robot is prepared for all possible scenarios. This usually meant that the code would exist of hundreds or even thousands of parts, each consisting of a certain type of action [2]. While for some shielded environments this might still be useful, robots nowadays often have to execute alongside humans. This demands a high safety level and the ability of robots to adapt to uncertainties in the environment. The latter is especially relevant for environments with humans, as humans are unpredictable.

There thus is a need for robots who can adapt to their environment. In other words, there is a need for robots who can learn. In this report, robot learning will be examined. This will be done by trying to answer the following research question:

How do methods, found in literature, enable continuous skill learning of robots, and how does the learned skill change during the learning process?

To narrow the scope of this research, two decisions were made. Starting, robotic learning can either be on- or offline. In this research, the focus will only be on online learning. Before explaining why, a description of this term should be given, as different literature often uses different descriptions. Online learning is seen in this report as any method which allows for some change in the model (dictating the behavior of the robot), throughout the learning process. This can both mean that the robot executes a behavior during the learning stage where it is seen that the behavior in iteration 1 is different from that in iteration 2. And also, the updating of a model which cannot directly be seen in the real world.

To elaborate why only online learning will be examined, first offline learning should be explained. In offline learning, data can be acquired based on which a model, representing the robotic skill, is trained. This model is then used in an autonomous stage to let the robot execute certain tasks. The advantages of this method are that the operator can select the data or adjust it, to achieve the best learning performance [3]. A downside to pure offline learning is that no feedback about the behavior is provided. This can partly be addressed by providing a variety of data (e.g. different tasks or different starting points) or by correcting the skill during the autonomous stage. But this again requires a high level of robot intelligence to recognize the human device or advise a strategy to correct its performance on the sensorimotor level. Especially for complex tasks, this can be time-consuming and sometimes even unsuccessful. Therefore, online learning might be more appropriate. This method assumes that the robotic skill is gradually formed during the learning stage [3]. An advantage of this approach is that the transition between the learning and execution stage can be direct and autonomous. This approach also allows for useful feedback about the performance while the robotic skill is constructed. Another motivation for online model learning is that it is not possible to cover the complete state-space with data beforehand, but only the interesting state-space regions are known during execution [4]. Therefore, online learning will require the incremental acquisition of knowledge and, possibly, even partial forgetting of the recorded information in order to cope with errors as well as change.

In addition to online learning, deep learning will be excluded from this research. The main reason for this is that it was decided to be of importance to make the learning process of the robot understandable for humans. Deep learning uses a model, which consists of some structured internal scheme (a neural network), to represent the behavior of a robot [5]. While being updated, it is not completely understood what is going on. Therefore, if a robot executes some unexpected or undesired behavior, it cannot be easily explained. In addition, it can be preferable for some real-life tasks to be represented by a simple policy, instead of a complex model [6]. In contrast, the methods examined in this paper, are often based on some human-specified optimal behavior. For example, some reward function [7], or some behavior demonstrated by an expert to imitate [8].

To answer the research question, two topics will be discussed in this paper. Each of these topics relates to a part of the research question. The first topic is the enabling of continuous skill learning of robots, and the second is the skill change during the learning process. These will, respectively, be discussed in chapters 2 and 3. The main findings will be discussed in chapter 4, and in chapter 5 a conclusion, answering the research question, will be given.

2

Skill Learning Methods

An important goal of robotics is to allow robots to perform certain tasks in a smooth, natural way. In the past, this behavior was often programmed using a direct approach, also known as a "*hard-coded*" approach. While this approach is useful for well-known environments, problems arise when this is not the case. Directly programmed robots are, for example, unable to adapt to changing environments. This is especially important when these robots work alongside humans, as safety becomes a top priority. But even when this is not the case, adapting to the environment prevents the robot from harming itself and/or the environment. Two learning approaches that have proved to be successful are Imitation Learning (IL) and Reinforcement Learning (RL). The former uses demonstration, exerted by an expert, to train the robot. The latter takes a trial-and-error approach, where the environment is explored to perform the desired task as good as possible. An overview of the advantages and disadvantages of these three approaches is shown in table 2.1.

Table 2.1: Comparison of three methods based on which robot can require skills [9]. The first method, *direct programming*, is "*hard-coded*", whereas the other two, *IL* and *RL*, allow for continuous skill learning.

Method	Advantages	Disadvantages
Direct Programming	<ul style="list-style-type: none"> • Complete control of movement robot to the lowest level 	<ul style="list-style-type: none"> • Time-consuming • Error-prone • Not scalable • Not reusable
Imitation Learning	<ul style="list-style-type: none"> • No manual program of behavior • Solving problems which can be demonstrated, but not succinctly specified in mathematical form • Low number of demonstrations needed • Save design • High reliability 	<ul style="list-style-type: none"> • Correspondence problem: demonstrators actuation differs from human • Cannot deal properly with suboptimal data • Covariant shift/compounding error when dealing with other input distributions • Performance robot is limited by skill of expert
Reinforcement Learning	<ul style="list-style-type: none"> • No manual program of behavior • Robot can learn tasks that humans cannot demonstrate • Novel ways to reach a goal can be discovered 	<ul style="list-style-type: none"> • Requires lots of trials • No control over the actions of the robot (unsafe) • Robot has only indirect information about the goal • Need to specify reward function, policy parameterization, exploration magnitude/strategy and initial policy

In this chapter, the focus will only be on learning a policy directly. Another approach for robotic learning is by skill transfer Kroemer, Niekum, and Konidaris. This enables the robot to transfer a skill, learned in one task, to another task. An advantage of this method is that it can increase the efficiency of learning. As skill transfer does not allow detecting how the skill of a robot has been learned, and how it has improved through the learning process, it was decided that it would be excluded for further research.

This chapter is dedicated to presenting the theoretical background behind the two different learning approaches which can adapt to the environment: IL and RL. To start, in section 2.1 the Markov Decision Process (MDP) is explained, which is often used in robot teaching. This is followed by an overview of the IL approach in section 2.2 and the RL approach in section 2.3. For both approaches, this includes both a problem definition and an explanation of some methods. This is followed by a short description about the potential of combining IL and RL methods, in section 2.4. In section 2.5, an overview of the different robot teaching approaches is given. Lastly, a summary of the chapter is presented in section 2.6.

2.1. Markov Decision Process

Skill learning often requires the learning of a policy that best represents the desired behavior of a robot. To learn this policy, different methods could be employed. Most of these methods use the same type of problem. Looking at skill learning, often MDP is used as underlying decision-making [5, 10]. MDP is described as a tuple $\langle S, A, R, T, \gamma \rangle$, where:

- $S \subseteq \mathbb{R}^n$ is a finite set of N states
- $A = \{a_1, \dots, a_k\} \subseteq \mathbb{R}^m$ is a set of k actions
- $R : S \mapsto \mathbb{R}$ is the reward function expressing the immediate reward of executing an action a in state s and transitioning into state s' , this is bounded by absolute value R_{\max} [11]
- $T = P(s \mid a, s')$ is the transition function, giving the probability distribution that states s' are reached after executing action a in state s .
- $\gamma \in [1, 0)$ a discount factor expressing the agent's preference for immediate over future rewards.

To solve these MDPs, a policy $\pi : S \times A \rightarrow [0, 1]$ should be found, that maximizes the expected sum of rewards for a specific problem [10]. For the MDP to hold, it should satisfy the Markov property. This property states, that the next state s' and the reward only depend on the previous state s and action a , and not on any additional information about the past states or actions.

2.2. Imitation Learning

One frequently used method for learning the skill policies, is IL. In literature, IL is also referred to as Programming by Demonstration (PbD), Learning from (human) Demonstration (LfD), and apprenticeship learning [2]. For the remaining of this paper, the term IL will be used. The choice of IL over other robot learning methods is particularly compelling when ideal behavior can neither be scripted (as done in traditional robot programming) nor be easily defined as an optimization of a known reward function (as done in RL), but it can be demonstrated. A disadvantage of IL techniques, is that - due to only learning from demonstrations - the performance is limited by the abilities of the teacher.

The process of learning a skill using IL, consists of three steps: observing an action, representing the action, and reproducing the action [12]. For the first step, a demonstration of the desired task, performed by an expert, is required. The expert is an agent who is assumed to know a certain skill. There are different types of demonstration methods that have been used in literature [12–14]:

- *Kinesthetic Teaching*: where the human holds the robot limb and physically guides it to perform the desired task.
- *Teleoperation*: where the human performs and teaches the task by controlling the robot through human-robot interfaces.
- *Passive Observation*: where the human performs the task on his/her own, while the robot then learns from the observation.

- *Shadowing*: where the executed behavior of the human is recorded using sensors and the robot tries to match or mimic the teacher's motions while executing the task.

The remaining of this section consist of a general problem definition of the IL approaches in section 2.2.1. Followed by two broadly used IL approach: Behavioral Cloning (BC) and Inverse Reinforcement Learning (IRL), discussed in sections 2.2.2 and 2.2.3, respectively. This also includes some methods which can be used for online learning.

2.2.1. Problem definition Imitation Learning

There are numerous methods to tackle IL, which all try to solve the same problem. Starting, the goal of IL is to learn a policy π_θ which can reproduce the behavior demonstrated by the experts for performing a certain task [8]. The behavior of the expert (and eventually the learner) can be described as a trajectory $\tau = [(s_1, a_1), \dots, (s_n, a_n)]$, which consist of a sequence of state-action pairs. This can also be represented using the feature ϕ . To learn a policy, it is assumed that an expert is exerting a certain policy π^E . As this cannot directly be transferred to the robot, it is instead done by demonstrations. The demonstrations are represented by a dataset $\mathcal{D} = (\tau_i, r_i)_{i=1}^N$, where the reward signal r is optional. Using this dataset, the IL problem can be reframed as an optimization problem:

$$\pi^* = \arg \min (D(q(\phi), p(\phi))), \quad (2.1)$$

where the optimal policy π^* is learned, such that the difference between the learned sequence and the demonstrated trajectory, is minimized. Here, $q(\phi)$ is the distribution of features induced by the experts' policy, $p(\phi)$ the feature distribution induced by the learner, and $D(q, p)$ is the similarity measure between q and p .

Learning approaches

There are different IL approaches that can be used to find the policy. Often these approaches are categorized into two groups: Behavioral Cloning (BC) and Inverse Reinforcement Learning (IRL). BC obtains the policies using supervised learning [15]. It does this by directly mapping the state. An advantage of this technique is that it is quite simple. However, it only tends to work successfully with large amounts of data, due to compounding error [16]. IRL describes the policy as the solution of an optimization or planning problem [8]. As this method learns a cost function that prioritizes the entire trajectory over others, compounding errors is not an issue [15]. Given a reward signal, IRL can obtain the policy such that it maximizes the expected return $J(\pi)$. As the reward function is unknown, it needs to be recovered from the expert demonstrations (here we assume that the behavior exerted in the demonstrations is (approximately) optimal). By maximizing the expectation of the accumulated rewards $J(\hat{\pi})$, the policy can then be found. An overview of the advantages and disadvantages of both approaches is given in table 2.2.

Table 2.2: Advantages and disadvantage of the two main IL approaches based on multiple papers [15–18].

	Advantages	Disadvantages
BC	<ul style="list-style-type: none"> • Computational efficient • Performance in high-dimensional space 	<ul style="list-style-type: none"> • Compounding error caused by covariant shift • Often large amount of data required for stable policy • Performance in dynamical and suboptimal space
IRL	<ul style="list-style-type: none"> • No compounding error • Good performance in suboptimal space 	<ul style="list-style-type: none"> • Expensive to run • Can be impractical in real-world • Performance in high-dimensional space • Matches at high level outcomes instead of actions

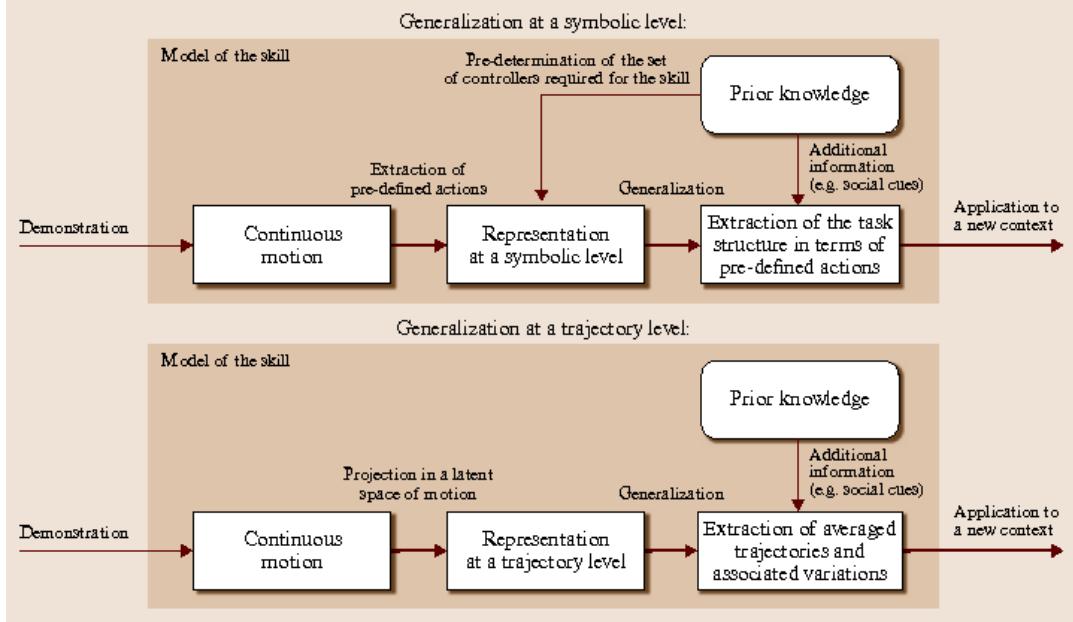


Figure 2.1: Illustration of the different levels of representation for describing the skill [19]. The upper figure shows the model of the skill for a symbolic level, whereas the lower illustrated the trajectory level.

Policy Representation

As it is the eventual goal of IL to define a policy that best represents the behavior of the expert, an important consideration is how the policy should be represented. This should be chosen such that it can capture the desired behavior properly [8]. Besides, the complexity of the policy also should be taken into account when choosing a representation; If the complexity of the representation is increased, the model will probably capture the desired behavior better, however, it will in most cases also lead to an increase of the learning time, and it would also require more training data.

The representation can either be on symbolic-level (high-level) or trajectory level (low-level) [19]. In the symbolic representation, the predefined motion elements are organized sequentially. Therefore, it allows for efficient learning of hierarchy, rules, and loop, through an interactive process. On the downside, this method depends on a large amount of prior knowledge to predefine important cues and to segment them efficiently. The trajectory-level representation uses a non-linear mapping between the sensory and motor information. It allows for the results to be a generic representation of the movement. The advantage of this is that it allows for the encoding of different types of signals or gestures. This method does, however, not allow for the reproduction of complicated high-level skills. An illustration of both methods is given in fig. 2.1. Besides using one of these policy representations, research has also been successfully conducted where both types were combined. In [20], high-level learning is done to learn a sequence of actions for an assembly task. And low-level learning is added to learn every path as needed for object manipulation.

2.2.2. Behavioural Cloning

Behavioral Cloning (BC) is the simplest form of IL. This method tries to learn the policy, by directly mapping from state/context to trajectories/actions or the control input [8]. An overview of the BC approach is given in algorithm 1. The first step is to obtain the experts' demonstrations \mathcal{D} . Different methods for doing so were discussed earlier. Next, an appropriate policy representation has to be chosen. In addition, an objective function \mathcal{L} , which represents the similarity between the demonstrated behaviors and the learner's policy, should be chosen. Some well-known and frequently used loss functions are: quadratic loss function (also known as ℓ_2 -loss function), ℓ_1 -loss function (also known as absolute loss function), log-loss function, hinge loss, and the Kullback-Leibler Divergence. Now, using the dataset of demonstration, the policy parameters θ can be learned.

Algorithm 1 Abstract of BC

-
- 1: Collect a set of trajectories demonstrated by the expert \mathcal{D}
 - 2: Select a policy representation π_θ
 - 3: Select an objective function \mathcal{L}
 - 4: Optimize \mathcal{L} w.r.t. the policy parameter θ using \mathcal{D}
 - 5: **return** optimized policy parameters θ
-

BC problems can be formulated as a supervised learning problem [8]. For such problems, the policy is obtained using a regression problem. However, it should in all cases be ensured that the trajectory is physically feasible, and therefore naive regression methods are not always the best choice. To ensure feasibility, constraints should be added. In some cases, it is possible that the constraints implicitly satisfy using a regression method. However, in most cases, it is more convenient to use a policy that explicitly satisfies some constraints. To learn these policies, the regression methods are used in a way that the constraints satisfy.

There are different methods that apply the principle of BC for learning the policy. In the next part of this section some methods, frequently used in literature, are explained. It should be noted that most of these methods require some adaptation of the original method, in order to be applicable for online learning.

Dynamic Movement Primitives

Motion Primitives (MPs) are often used to represent and learn basic movements in robotics. A frequently used method for this representation is Dynamic Movement Primitive (DMP), which was originally introduced by Ijspeert, Nakanishi, and Schaal [21, 22]. DMP provides a framework for the motor representation, based on a nonlinear dynamic system [23].

There are two general types of DMPs: discrete and rhythmic [24]. Discrete DMPs are used to encode a point-to-point motion into a stable dynamical system. Whereas, rhythmic DMPs are used to encode motion followed from a rhythmic motion. DMPs are defined using two sets of equations: the transformation system and the canonical system. The transformation system can be written as a first-order notation of damper spring model [22]

$$\begin{aligned} \tau \dot{z} &= \alpha_z (\beta_z(g - y) - z) + f(x), \\ \tau \dot{y} &= z, \end{aligned} \quad (2.2)$$

where y is the state/position of the system, \dot{y} the velocity of the joint trajectory, \ddot{y} the acceleration, z the generalized velocity, τ a time constant, g the goal state, α_z and β_z the gain terms representing the stiffness and damping, respectively, and f forced item used to restrain the repair of the trajectory. Looking at the forcing term $f(x)$, if this is equal to 0, the transformation equations represent a globally stable second-order linear system with $(z, y) = (0, g)$ as a unique point attractor [22]. An example of this has been given in fig. 2.2. To allow more complex trajectories to the goal, f can be defined as a time-dependent, non-linear function. This results in a nonlinear Dynamic System (DS) which can be difficult to solve. Therefore, an additional differential equation, called the *canonical system*, is introduced. For discrete DMPs, this system can be written as:

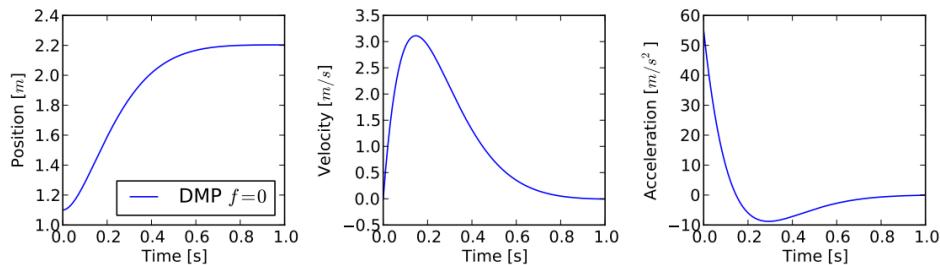


Figure 2.2: Position, velocity and acceleration of an executed DMP without any external force ($f = 0$) [25].

$$\tau \dot{x} = \alpha_x x, \quad (2.3)$$

where α_x is a predefined constant and x the phase variable. Defining $f(x)$ as a linear combination of N nonlinear Radial Basis Functions (RBFs), will ensure that the robot will follow any smooth trajectory from an initial position y_0 to the final configuration g [24]

$$f(x) = \frac{\sum_{i=1}^N w_i \Psi_i(x)}{\sum_{i=1}^N \Psi_i(x)} x, \quad (2.4)$$

where the basis function is defined as $\Psi_i(x) = \exp(-h_i(x - c_i)^2)$, c_i are the centers of the function distribution along the phase of the movement and h_i their widths. The parameter w , can here be adjusted using learning algorithms, in order to produce complex trajectories.

Rhythmic DMPs, use a transformation system which is quite similar to its discrete variant (eq. (2.2)). Instead of the gain α_z , a frequency Ω is used. And the forcing item is instead dependent on a phase angle of the oscillator in polar coordinates ϕ : $f(\phi)$. The canonical system is then given as:

$$\tau \dot{\phi} = 1, \quad (2.5)$$

Similar to the discrete variant in eq. (2.4), $f(\phi)$ is defined with N kernels

$$f(\phi) = \frac{\sum_{i=1}^N \Psi_i(\phi) w_i r}{\sum_{i=1}^N \Psi_i(\phi)}, \quad (2.6)$$

where $\Psi_i(\phi) = \exp(h(\cos(\phi - c_i) - 1))$, the weights are uniformly distributed along the phase's pace, and r is used to modulate the amplitude of the periodic signal.

To learn the encoded trajectories, a regression method should be used. One method which is frequently used for online learning problems is *Locally Weighted Regression (LWR)* [24]. This method is suitable for these types of problems, as it updates the existing models fast [26]. In essence, the goal is to find the appropriate weight w_i of f . To formulate the function approximation problem, the target function f_{target} should be defined. This can be done by rewriting eq. (2.2), resulting in

$$f_{target} = \tau^2 \ddot{y}_{demo} - \alpha_z (\beta_z (g - y_{demo}) - \tau \dot{y}_{demo}). \quad (2.7)$$

The weight of f can then be determined by minimizing the error criterion $J_i = \sum_{t=1}^P \phi_i(t)(f_{target}(t) - w_i \xi(t))^2$, with $\xi(t) = x(t)(g - y_0)$ for a discrete system and $\xi(t) = r$ for a rhythmic system [22]. A graphical example of this has been given in fig. 2.3.

Some extension on the original DMP approach exist. The classical DMPs approach is meant for single Degrees of Freedom (DoF) motions. In order to obtain multidimensional motions, the transformation

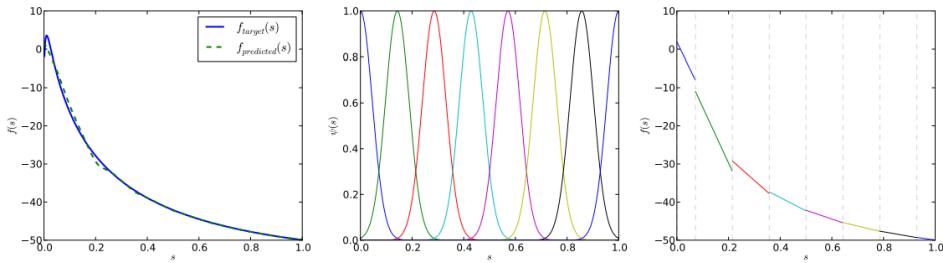


Figure 2.3: Example of the learning of DMP using LWR, where: (left) the target function f_{target} (blue) and the predicted function $f_{predicted}$ (dashed green), (middle) 8 Gaussian functions meaning 8 weights ($|w| = 8$), and (right) the local linear model [25].

system eq. (2.2) have to be repeated while the canonical system eq. (2.3) is shared. This is not a problem when the evolution of the different DoFs are independent. However, a problem arises for orientation representation, where the elements are constraint. Extensions which could be applied in these types of situation are Quaternion MDP [27], rotation matrix MDP [28], and symmetric positive definite matrices [29].

Typically, DMP is not used for online learning. A reason for this is the usage of only one demonstration. Some papers have extended this original method to allow for online learning. An example is given in [30], where a learning scheme was presented which iteratively updates the DMP.

Probabilistic Movement Primitives

A disadvantage of DMP is that it can only represent the mean solution, which is known to be suboptimal [31]. Instead of using such a deterministic approach, a probabilistic approach called *Probabilistic Movement Primitive (ProMP)* has been introduced by Paraschos et al. [32]. This allows for the direct encoding of an optimal behavior for a system with linear dynamics, quadratic costs, and Gaussian noise. This method provides a single framework for all the desirable properties of a MP: co-activation, modulation, optimality, coupling, learning, temporal scaling, and rhythmic movements [32].

ProMPs uses multiple trajectories to define a single movement, which automatically results in distribution over the trajectories. This is different from DMP, which only uses one trajectory. A single trajectory τ can be compactly represented by a weight vector w . Here, the trajectory is represented by $\{q_t\}_{t=0 \dots T}$, which is defined by the joint angle q_t over time. Using this will reduce the number of model parameters. The state vector y_t for a single time step can now be defined

$$y_t = \begin{bmatrix} q_t \\ \dot{q}_t \end{bmatrix} = \Psi_t^T w + \epsilon_y, \quad (2.8)$$

where $\Psi_t = [\psi_t, \dot{\psi}_t]$ describes a $n \times 2$ time-dependent basis matrix for the joint position q_t and velocities \dot{q}_t , n defines the number of basis functions, and $\epsilon_y \sim \mathcal{N}(0, \Sigma_y)$ is a zero-mean independent and identically distributed (i.i.d.) Gaussian noise. The probability of observing the whole trajectory $p(\tau | w)$ is then defined as

$$p(\tau | w) = \prod_{t=1}^T \mathcal{N}(y_t | \Psi_t w, \Sigma_y). \quad (2.9)$$

As the expectation is that different demonstrations of the same movements, will be slightly different, there are different weight vectors w_n needed to represent the n different instances of a movement [33]. To capture this variance in the trajectories, the distribution $p(w; \theta)$ over the weight vector w , with parameter θ , is introduced. In most cases, this distribution will be Gaussian, where the parameter vector $\theta\{\mu_w, \Sigma_w\}$ specifies the mean μ_w and variance Σ_w of w [31]. Besides using a Gaussian distribution, more complex distributions like a Gaussian Mixture Model (GMM) could also be used for this. By marginalizing out the weights w , the trajectory distribution $p(\tau; w)$ can now be computed

$$p(\tau; \theta) = \int p(\tau | w)p(w; \theta)dw. \quad (2.10)$$

This distribution $p(\tau; w)$ defines a Hierarchical Bayesian Model (HBM), whose parameters are given by the observation noise variance Σ_y and the parameters θ of $p(w; \theta)$.

The choice of basis function depends on the type of movement. The movement can either be stroke-based or rhythmic [32]. For stroke-based movements, a Gaussian basis function is used, while for rhythmic movements, a Von-Mises basis function is used to model periodicity in the phase variable. These basis functions are in most cases normalized, resulting in the normalized basis function ψ_t .

Up to this point, it was considered that each DoF was modeled independently. However, often the movement of the multiple joints has to be coordinated. A method for doing this is to introduce the phase variable z_t , which couples the mean of the trajectory distribution [32]. As it is often also desired

to encode higher-order moments of the coupling, the model is extended to multiple dimensions. Here, each dimension maintains a parameter vector w_i . They can be combined into a weight vector $w = [w_1^T, \dots, w_n^T]$. The basis matrix Ψ_t extends now to a block-diagonal matrix, which contains the basis functions and their derivatives for each dimension. The observation vector y_t consists of the angles and velocities of all joints. The probability of an observation y_t at time t is given by

$$p(y_t | w) = \mathcal{N} \left(\begin{bmatrix} y_{1,t} \\ \vdots \\ y_{d,t} \end{bmatrix} \mid \begin{bmatrix} \Psi_t^T & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & \Psi_t^T \end{bmatrix} w, \Sigma_y \right) = \mathcal{N}(y_t | \Psi_t w, \Sigma_y) \quad (2.11)$$

For MP representations, the parameters of a single primitive must be easy to acquire from demonstration [32]. To facilitate the estimations of these parameters, a Gaussian distribution for $p(w; \theta) = \mathcal{N}(w | \mu_w, \Sigma_w)$ over the parameters w is assumed. This could, however, also be done using other types of probability distributions, e.g. a mixture model. The distribution of the state $p(y_t | \theta)$ for time step t is then given by

$$p(y_t; \theta) = \int \mathcal{N}(y_t | \Psi_t^T w, \Sigma_y) \mathcal{N}(w | \mu_w, \Sigma_w) dw = \mathcal{N}(y_t | \Psi_t^T \mu_w, \Psi_t^T \Sigma_w \Psi_t + \Sigma_y) \quad (2.12)$$

Now, the mean and variance for any point t can easily be evaluated. As ProMP represents multiple ways to execute a specific movement, multiple demonstrations are needed to learn $p(w; \theta)$. The parameters $\theta\{\mu_w, \Sigma_w\}$ can be learned from these demonstrations by the maximum likelihood estimation for HBM with Gaussian distribution.

In fig. 2.4, a trajectory distribution of a single movement, generated using different approaches, is seen. Figure 2.4a shows the demonstrated trajectory distribution, which was generated by a stochastic optimal control algorithm for a via-point task. The variability in this distribution is due to the noise of the system. In contrast to the trajectory distribution generated using DMP (Figure 2.4d), does the one generated using ProMP (Figure 2.4b) almost perfectly match the demonstrated distribution. The example shown in fig. 2.4 presents a single movement distribution. In fig. 2.5, the implementation of ProMP for multiple goal positions is given. The left figure presents the distribution of the movements for five different targets.

ProMPs are a relatively new topic of research in the literature. Therefore, not much research can be found implementing ProMP for online learning. However, one extension on ProMP seems to be promising: interaction ProMP. Interaction ProMP is a framework which can be used for collaboration with a human [34]. By taking inspiration from Interaction Primitives, ProMPs which realizes the primitives that capture the correlation between trajectories of multiple agents, is proposed. Interaction ProMPs uses the distribution over the trajectories of at least two interacting agents. Koert et al. use this principle in an online learning setting [35].

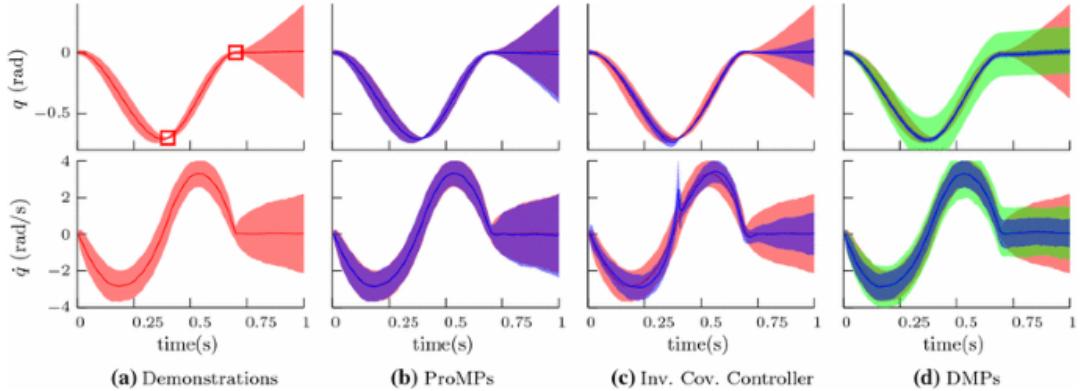


Figure 2.4: The trajectory distribution, showing the joint position (upper plot) and velocity (lower plot). The shaded area is two times the standard variant. The red shaded area shows the demonstrated trajectory [31].

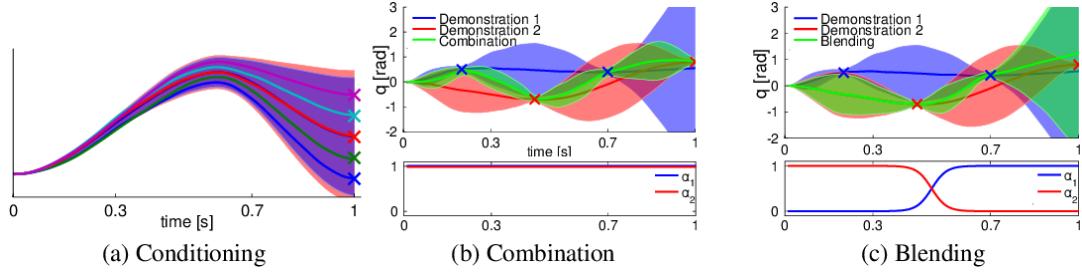


Figure 2.5: The trajectory distribution is encoded using ProMP for multiple goal positions. Here, two methods of combining the different demonstrations are indicated: combination and blending.

Gaussian Mixture Model and Gaussian Mixture Regression

Gaussian has received a lot of attention in the literature. Therefore, it does not come as a surprise that this has also been frequently used for IL problems. Using just a Gaussian Regression method to solve IL problems would in most cases not provide enough complexity. Therefore, a combination of Gaussian Mixture Model (GMM) with Gaussian Mixture Regression (GMR) is frequently used. This method is robust to noisy data, has a low computation cost, and can capture the correlation of human motion. Therefore, it is a powerful tool for robot analysis [36]. In this case, GMM is used to encode the behavior of an expert and GMR to reproduce it.

Just as any other IL method, GMM requires a dataset \mathcal{D} , describing the motion of the demonstrations. Defining $\xi_t = (s_t, a_t)$, the joint distribution of s_{t+1} and ξ_t can be modelled as a Gaussian distribution

$$p(s_{t+1}, \xi_t) = \sum_k p(k) \mathcal{N}(\mu_k, \Sigma_k) \quad (2.13)$$

where $p(k)$ is the prior. The k^{th} Gaussian component can be described by

$$\begin{aligned} p(s_{t+1}, \xi_t | k) &= \mathcal{N}\left(\begin{bmatrix} \xi_t \\ x_{t+1} \end{bmatrix} | \begin{bmatrix} \mu_{\xi,k} \\ \mu_{x,k} \end{bmatrix}, \begin{bmatrix} \Sigma_{\xi,k} & \Sigma_{\xi,x,k} \\ \Sigma_{x\xi,k} & \Sigma_{x,k} \end{bmatrix}\right) \\ &= \frac{1}{\sqrt{(2\pi)^D |\Sigma_k|}} e^{-\frac{1}{2}((\xi_t - \mu_k)^T \Sigma_k^{-1} (\xi_t - \mu_k))}, \end{aligned} \quad (2.14)$$

where $\{\pi_k, \mu_k, \Sigma_k\}$ (the prior probability, mean vector, and covariance matrix) are the parameters of the Gaussian component k . Traditionally, such a model is iteratively trained using Expectation Maximization (EM). Other examples are spectral clustering, online learning, or self-refinement. The learned model can then be used to reproduce movements [37]. An example of the encoding using GMM is visualized in fig. 2.6.

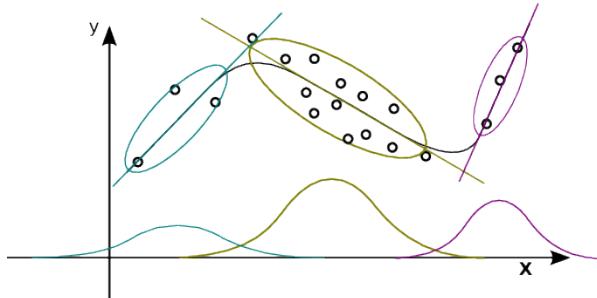


Figure 2.6: Example of trajectory encoding using GMM. The data points are indicated with a dot. The data points are clustered using three Gaussian components. Based on this, the Gaussian distributions are determined (bottom part of the plot), which results in a trajectory that is indicated with a grey, continuous line.

GMR can be used to reconstruct a general form for the signal [37]. GMR first models the joint probability density function of the data and then derives the regression function from this model, which is different from other regression functions which directly derive the regression function. In GMR, the estimation of the model parameters is done in the offline phase, making the regression independent of the number of data points. Therefore, the regression can be calculated very rapidly.

For each iteration step t , the data points can be composed into two parts, the temporal values, denoted with a t , and the spatial values, denoted with a s . This results in the following notation for the data points ξ , mean vectors μ , and covariance matrices Σ

$$\xi = \begin{bmatrix} \xi_t \\ \xi_s \end{bmatrix}, \quad \mu_k = \begin{bmatrix} \mu_{t,k} \\ \mu_{s,k} \end{bmatrix}, \quad \Sigma_k = \begin{bmatrix} \sum_{t,k} & \sum_{ts,k} \\ \sum_{st,k} & \sum_{s,k} \end{bmatrix}. \quad (2.15)$$

For each component k , the expected distribution of $\xi_{s,k}$, given ξ_t , and the estimated covariance, can be determined:

$$\begin{aligned} \hat{\xi}_{s,k} &= \mu_{s,k} + \Sigma_{st,k} (\Sigma_{t,k})^{-1} (\xi_t - \mu_{t,k}), \\ \hat{\Sigma}_{s,k} &= \Sigma_{s,k} - \Sigma_{st,k} (\Sigma_{t,k})^{-1} \Sigma_{ts,k}, \end{aligned} \quad (2.16)$$

where $\xi_{s,k}$ and $\hat{\Sigma}_{s,k}$ are mixed according to the probability that the Gaussian component $k \in \{1, \dots, K\}$, has been responsible for ξ_t . For a mixture of K components, the condition expectation of ξ_s , given ξ_t , and the conditional covariance of ξ_s , given ξ_t , can be described as

$$\hat{\xi}_s = \sum_{k=1}^K h_k \hat{\xi}_k, \quad \hat{\Sigma}_s = \sum_{k=1}^K h_k^2 \hat{\Sigma}_k, \quad (2.17)$$

where h_k is the probability of the component k to be responsible for t . An example of trajectory encoding using GMM and GMR is shown in fig. 2.7. This figure shows that at the beginning of the task execution, the constraints are loose (distribution has a high variance), whereas at the end it is quite strict.

In [38] two methods for incrementally training the model were given. The first method, the direct update method, reformulates the problem for a generic observation of multiple data points. This is done using

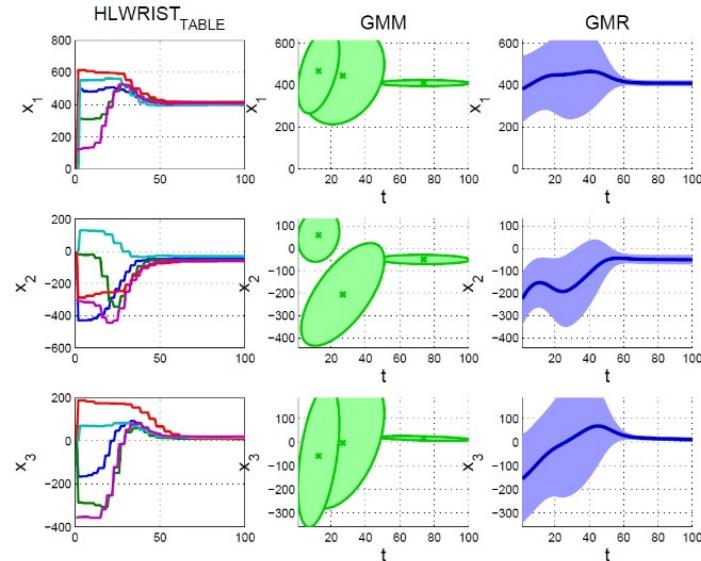


Figure 2.7: Trajectory encoding and generalization using GMM and GMR [36]. It is visible that at the start of execution the constraints are loose as the distribution has a high variance. In contrast, the final part shows that the constraints are quite strict, as the variance of the distribution is low.

an adapted EM method, where the parts belonging to the already used data and those belonging to the newly available data are separated. The second method, the generative method, used EM performed on data generated by GMR. Another paper [39] proposed three methods for incrementally updating a set of already existing trajectories, which were obtained using task-parametrized Gaussian Mixture Model (TP-GMM) [37]. TP-GMM in essence, models local (or relative) trajectories and corresponding local patterns, therefore endowing GMM with better extrapolation performance [40]. While these techniques do assume that previously trajectories were already obtained (either online or offline), this does show the ability to adapt in an online manner. The first technique estimates a new model by accumulating the new trajectory and a set of trajectories using the old model. The second techniques allow for adding parameters for the new trajectories, to the already existing parameters corresponding to the old trajectories. The last techniques update the model by the usage of a modified EM algorithm, using the information of the new parameters. Cederborg et al. [41] proposed incremental local online GMR. This method allows for the robot to learn incrementally online a new motor task, by modeling them locally as dynamic systems. Sensorimotor context is used to cope with the absence of categorical information both during demonstrations and when reproduction is asked of the system.

Gaussian Process

In the recent years, Gaussian Processes (GPs) as a solution for IL problems, has received more attention. GPs can be defined as a distribution over function, where inference directly takes place in the space of function [42]. GP aims to learn a deterministic input-output relationship, up to observation noise, based on a Gaussian prior over potential objective functions [43]. This is in contrast to GMR which provides a generative model. The difference between discriminative and generative models is graphically visualized in see fig. 2.8. While generative models, in most cases, require fewer data, their generalization performance is often poorer than that of a discriminative model.

GP models a mapping from the input $z = [x_t^T, u_t^T]$ to the output $x_{t+1} = f(z_t)$ as

$$f(z_t) \sim \mathcal{GP}(m(z_t), k(z_t, z'_t)) \quad (2.18)$$

where $k(z_t, z'_t)$ is the covariance matrix [8]. Popular approaches for the covariance matrix are the squared exponential covariance function [43] and a Gaussian kernel [45]. The joint distribution of the given target value and the function value x_{t+1}^* at the test input z_t^* can be written as

$$\begin{bmatrix} x_{t+1} \\ x_{t+1}^* \end{bmatrix} \sim \mathcal{N} \left(0, \begin{bmatrix} K(Z, Z) + \sigma_n^2 I & K(Z, z_t^*) \\ K(z_t^*, Z) & K(z_t^*, z_t^*) \end{bmatrix} \right), \quad (2.19)$$

where Z is a matrix, where all input factor z_t of the training samples are combined. The mean μ and variance σ^2 are in this case dedicated by

$$\begin{aligned} \mu(z_t^*) &= \mathbf{K}(z_t^*, Z)^\top (\mathbf{K}(Z, Z) + \sigma_n^2 I)^{-1} \mathbf{x}_{t+1} = \mathbf{K}(z_t^*, Z)^\top \alpha \\ \sigma^2(z_t^*) &= k(z_t^*, z_t^*) - \mathbf{K}(z_t^*, Z)^\top (\mathbf{K}(Z, Z) + \sigma_n^2 I)^{-1} \mathbf{K}(z_t^*, Z), \end{aligned} \quad (2.20)$$

where α denotes a prediction vector.

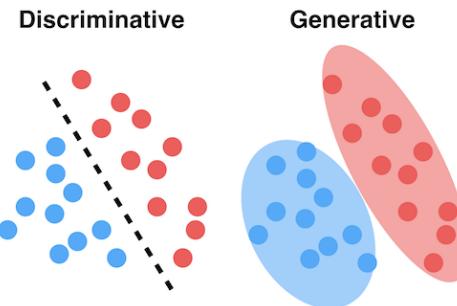


Figure 2.8: Graphical representation of a generative and discriminative model [44].

The explanation above presents the broad outline of GP. There have been multiple attempts to implement this approach for online learning. Nguyen-Tuong and Peters [45] proposed Local Gaussian Process (LGP), which combined the fast computation of local regression with a more accurate regression method. In [42] a GP-based learning approach was presented, which allows for generalization over multiple demonstrations, and encodes variability along the different phases of the to be executed task. Jaquier et al. [43] proposed a GMR-based GP. This combination allows combining the advantages of GP to encode prior beliefs through the mean and kernel functions and of GMR to retrieve the variability information, to be encapsulated in the uncertainty estimated by the GP. Another advantage is that this approach has the properties of a generative model, which allows for new trajectories to be easily generated through sampling and conditioning. In [46] the Sparse Online Gaussian Processes (SOGP) is compared to Locally Weighted Projection Regression (LWPR). SOGP allows for data to be processed while it arrives. An advantage of this method, compared to the original Gaussian Process Regression (GPR) approach, is that it only a subset of the data and their associated kernel distances should be stored. Making the storage space much smaller than the one required for GPR ($\mathcal{O}(N^2)$).

Hidden Markov Model

The usage of Hidden Markov Model (HMM) for IL can be seen as a modified version of GMM, in which the choice of mixture component for each observation also depends on the choice of the component for the previous one. It allows for the modeling of a probabilistic transition between discrete states [8]. A reason why HMM is a useful representation of human skill, is because of its ability to discover the nature of the skill [47]. It does this by representing the training data in a statistical way, by its parameters, which allows for the retrieving of the skill model. This is of importance as it is expected that each demonstration, given by a human expert, will differ from the next.

HMM models are characterized by a Markov chain of sequence, consisting of (unobserved) hidden state variables s_t and a corresponding sequence of observation variable o_t [48]. A graphical representation of this is given in fig. 2.9. The different states are connected using a state transition matrix $A = \{a_{ij}\}$, where $a_{1,2}$ is the transition probability of state 1 to state 2. The output probability matrix is defined as $B = \{b_{ij}\}$, which denotes the probability of observing a certain output while being in a certain state. The initial state distribution is also defined and denoted by d . Now, the HMM can be denoted as a set of matrices $\lambda = \{A, B, d\}$ [47]. Given this set of matrices λ and an observation sequence $O = \{o_1, \dots, o_T\}$, the likelihood of observing a given sequence $P(O | \lambda)$, can be computed. The observed motion can then be found using

$$\lambda^* = \arg \max_{\lambda} P(O | \lambda). \quad (2.21)$$

At the lowest level, each single MP is represented by a HMM. Given a set of MP $[\lambda_1, \dots, \lambda_T]$ and an observation sequence $[O_1, \dots, O_T]$, the state transition model and the probability distribution, can be learned, using the *Baum-Welch* algorithm [49].

To allow for incremental learning during observation, the system must be able to deal with new hidden states each time a new MP is learned. In addition, it should also have the ability to learn transition rules incrementally, from partial observations sequences. To do so, an incremental transition rule was proposed by Kulič et al. [49], which is applied each time a partial observation is available.

Using HMM allows for the encoding of the behavior, demonstrated by an expert. There is, however, still a need for a regression method that allows for the reproduction of this behavior. One method for doing so is GMR [50–52]. An advantage of using this, in contrast to other regression methods such

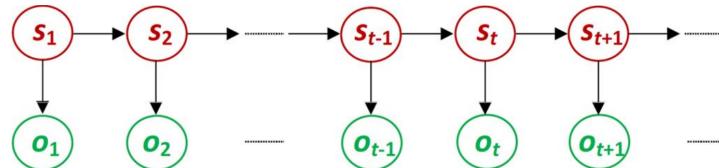


Figure 2.9: Graphical Representation of HMM [48]. With the s_n indicating the states and o_n the observations.

as LWR, LWPR or GPR, is that this does not model the regression function directly, but it models a joint probability density function of the data. After which it derives the regression function from the joint probability function.

A drawback of the HMM representation is discreteness [8]. When the number of states is relatively low, HMM works well. Too few states can however result in the model not being able to reproduce a motion sequence. For robotic applications, this is a problem, as HMMs are often used to describe discrete, high-level states. To overcome this, other techniques combined with HMM can be used. An example of this is the state-specific Gaussian models to represent continuous values [51].

In literature, some extensions of HMM have been found. Hidden Semi-Markov Model (HSMM) is defined as allowing the underlying process to be a semi-Markov chain [53]. A survey paper [53] identified multiple online implementations off HSMM, including an adaptive EM algorithm, an online algorithm based on recursive prediction error techniques, and a maximum likelihood estimation algorithm. Another extension is Hierarchy Hidden Markov Model (HHMM). In [54, 55], this framework is used to autonomously segment, cluster, and learn to a sequence of full-body motion primitives from online observations of the full-body human motion. This was expanded on in [49], where a higher abstraction HMM was used to learn higher-level ordering between movement primitives. Lastly, Growing Hidden Markov Model (GHMM) [56] was proposed to learn motion patterns incrementally, and in parallel with prediction. While this method seemed to be a promising extension of the original HMM method, most research seemed to be limited to implementations for trajectory estimation for vehicles.

2.2.3. Inverse Reinforcement Learning

Another form of IL is Inverse Reinforcement Learning (IRL). In IRL, the learner tries to recover a reward function of the environment, based on the expert's demonstration. Based on this reward function, it then tries to find the optimal policy using RL [11]. The used algorithm tries to find a reward that leads the learner to act similarly to the expert, while generalizing well to situations where the expert data is not available [10]. The assumption is made that the expert tries to maximize a certain reward function during the execution of the demonstrations, making it the aim to estimate this function. To summarize, IRL consist of two steps: 1) reward function estimation and 2) policy optimization based on this reward function. There are some variations to this general approach. For one, in some cases instead of finding a reward function, the goal is to find a cost function. This again makes the goal not to find a policy that minimizes this cost function. Another variation, found in literature [10], is to learn from failures instead of success. This approach is also referred to as IRL from Failure.

The IRL problem was originally defined by Russel [57]. He described it as the following:

Given:

- 1) measurements of an agent's behavior over time, in a variety of circumstances;
- 2) measurements of the sensory inputs to that agent;
- 3) a model of the physical environment (including the agent's body).

Determine: the reward function that the agent is optimizing.

Ng, Russell, et al. [11] defined three main types of IRL algorithms, which are still the most commonly used approaches [58]; 1) Finite-state MDP with known optimal policy; 2) Infinite-state MDP with known optimal policy; 3) Infinite-state MDP with unknown optimal policy, but demonstrations are given. The last of these approaches is the closest to practical problems, as it seems more realistic that only the expert's demonstrations are available rather than the desired policy. Therefore, only this formulation will be considered. By obtaining data from an expert interacting with a MDP, this approach tries to find a reward function R^* , of which it was assumed that the expert was trying to maximize [10, 59]. This is the opposite of RL, which tries to find a policy that tries to maximize the expectation for a given reward function, by interacting with a MDP. Therefore, IRL is formalized as an *incomplete* MDP, also known as an MDP/R or in mathematically terms \mathcal{M}/R . This is a tuple of the form $\langle S, A, T, \gamma \rangle$ [59].

Most IRL approaches assume that there is a set of M features associated with every state, which fully determine the value of the reward function. As finding a general form solution for R can be challenging, most approaches approximate it as a linear combination of the feature [11]. The reward function can then be defined as $R^* = w^{*T} \phi$. In the literature, examples have also been found where the reward

function is not approximated using this linear approach, however, for now, we will assume that it is. Assuming there is a feature mapping ϕ and the expert's feature expectation μ_E , the goal becomes finding a policy whose performance is close to that of the expert. The policy to be found $\tilde{\pi}$ should be found, such that $\|\mu(\tilde{\pi}) - \mu_E\|_2 \leq \epsilon$. For such $\tilde{\pi}$, for any $w \in \mathbb{R}^k$ ($\|w\|_1 \leq 1$, the following should hold

$$\begin{aligned} & \left| E \left[\sum_{t=0}^{\infty} \gamma^t R(s_t) \mid \pi_E \right] - E \left[\sum_{t=0}^{\infty} \gamma^t R(s_t) \mid \tilde{\pi} \right] \right| \\ &= |w^T \mu(\tilde{\pi}) - w^T \mu_E| \\ &\leq \|w\|_2 \|\mu(\tilde{\pi}) - \mu_E\|_2 \\ &\leq 1 \cdot \epsilon = \epsilon. \end{aligned} \quad (2.22)$$

Now the problem is reduced to finding a policy $\tilde{\pi}$, which induces feature expectation $\mu(\tilde{\pi})$ closest to μ_E . Here the policy can be found using policy iteration. A geometric description of policy iteration can be seen in fig. 2.10. The process of finding the policy $\tilde{\pi}$ can be described as:

1. Randomly pick some policy $\pi^{(0)}$ as initial policy, compute (or approximate via Monte Carlo (MC)) $\mu^{(0)} = \mu(\pi^{(0)})$, and set $i = 1$.
2. Compute $t^{(i)} = \max_{w: \|w\|_2 \leq 1} \min_{j \in \{0..(i-1)\}} w^T (\mu_E - \mu^{(j)})$, and let $w^{(i)}$ be the value of w that attains this maximum.
3. If $t^{(i)} \leq \epsilon$, then terminate.
4. Using the RL algorithm, compute the optimal policy $\pi^{(i)}$ for the MDP using rewards $R = (w^{(i)})^T \phi$.
5. Compute (or estimate) $\mu^{(i)} = \mu(\pi^{(i)})$.
6. Set $i = i + 1$, and go back to step 2.

A big problem in IRL is that multiple reward functions could explain the observations [60]. The reason for this is that the input is often a finite and small set of trajectories, resulting in many reward functions being able to realize the demonstrated data. Therefore, IRL suffers from an ambiguous solution. The literature often also refers to this problem as the result of IRL problems being "ill-posed" [11]. A method that tries to tackle this problem is *maximum margin*. It does this by converging on a solution that maximizes some margin. A downside of this method is that it introduces a bias into the learned reward function, thus resulting in exclusions of meaningful solutions [59]. To avoid this, the *maximum entropy* principle could be used. This approach is attractive as it is probabilistic and thus robust to noise and randomness in the demonstrations of the expert [10]. In the remainder of this section, both methods will be discussed.

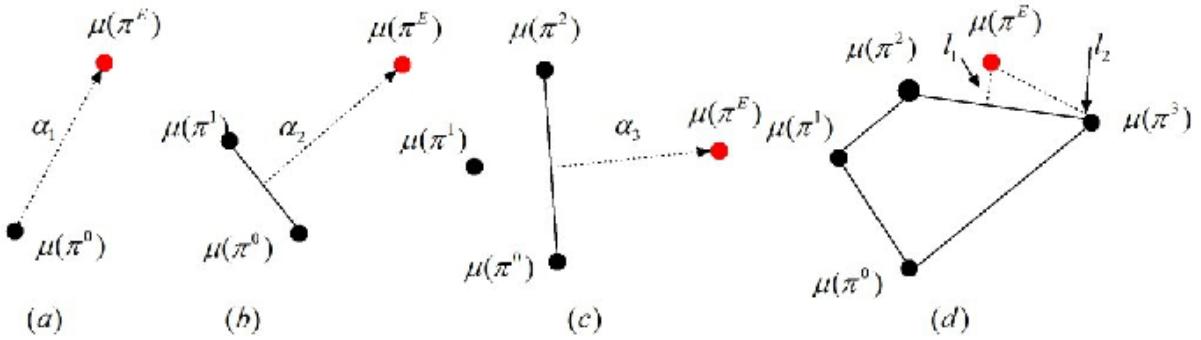


Figure 2.10: Geometric description of policy iteration [58]. The goal is to find a feature expectation $\mu(\tilde{\pi})$, with a different smaller than ϵ compared to the expert's feature expectation μ_E . This difference is indicated as l_1 . Each figure (a), (b), etc., is one step further in the iteration process.

Max-Margin Inverse Reinforcement Learning

In *maximum margin IRL*, the reward function is identified by maximizing the difference between the best policy and all other policies [13, 59]. The maximization step could be written as the following

$$\min_{t,w} \quad t \quad (2.23a)$$

$$\text{subject to} \quad w^T \mu_E \geq w^T \mu^{(j)} + t, j = 0, \dots, i-1 \quad (2.23b)$$

$$\|w\|_2 \leq 1. \quad (2.23c)$$

This shows that the goal is to find a reward function $R = w^{(i)}\phi$, where the expert does better, by a margin of t , than any of the policies previously found. It assumes that there is already an initial estimate of the reward function $j = 0$. The constraint shown in eq. (2.23c) shows a 2-norm constraint, which means that the problem cannot be seen as a linear program, but only as a quadratic program [59].

The apprenticeship learning problem as presented in [59] had as a downside that there was no mechanism for explicitly matching to the experts' behavior. In addition, the solution was a stochastic mixture of multiple policies [61]. Ratliff, Bagnell, and Zinkevich therefore, proposed maximum margin planning (MMP) [62], to address these problems. This frequently used approach produces a single deterministic solution while also ensuring an upper bound on the mismatch between the demonstrated and planned behavior. The research of Ratliff, Bagnell, and Zinkevich also provides an extension of MMP for online learning, where they also prove that the algorithm holds to a sublinear regret bound.

It was found that maximum margin IRL (or MMP), is most frequently used for path planning. Looking at [62], it was already stated that this paper also provided an extension for online learning. In addition, by the set definition, the original method also counts as online learning. This, as the model is iteratively being updated, until the margin is reached. Another example is given by Ziebart et al. [63], where MMP is used for autonomous navigation in unstructured terrain.

Besides planning problems, maximum margin IRL has not been found for much other robotic applications. A reason for this could be the existence of another IRL approach, namely maximum entropy IRL [63]. The maximum margin approach [59] has a major downside that the biases, which help to search in the ill-posed problem, can also rule out some other meaningful solution. As the maximum entropy approach makes fewer assumptions, it is often desired over maximum margin.

Maximum Entropy Inverse Reinforcement Learning

As the name states, *maximum entropy IRL* uses the maximum entropy principle to compute a reward function which best represents the demonstrated behavior [63]. This allows for obtaining a distribution over behavior, which are parameterized by the reward function [8]. By using a probabilistic approach, it tries to solve the IRL problem [58]. The probability of an observed experts' trajectory τ is weighted by the estimated reward. This results in policies with a higher reward, to be exponentially more preferred:

$$p(\tau | w) = \frac{1}{Z(w)} e^{w^T \phi(\tau)} \quad (2.24)$$

where $Z(w) = \sum_{\tau} \exp(w^T \phi(\tau))$ is the partition function. This expression of the probability of the trajectory, however, only holds for deterministic environments. In the case of stochastic environments, this probability is also affected by the transition probability:

$$p(\tau | w) = \frac{1}{Z(w)} \exp(w^T \phi(\tau)) \prod_{x_{t+1}, u_t, x_t \in \tau} p(x_{t+1} | u_t, x_t). \quad (2.25)$$

The following function is now tried to be optimized:

$$\tilde{R}(\tau) = \mathbf{w}^T \phi(\tau) + \sum \log p(x_{t+1} | \mathbf{u}_t, \mathbf{x}_t) \quad (2.26)$$

Visible is that there is now a bias term due to the stochasticity of the environment, which is the main (theoretical) drawback of this approach. This has, however, been addressed by the maximum causal entropy IRL [8]. The optimal value of the parameter vector w of the reward vector R , is given by maximizing the likelihood of the observed trajectory through maximum entropy:

$$w^* = \underset{w}{\operatorname{argmax}} \mathcal{L}_{ME}(w) = \underset{w}{\operatorname{argmax}} \sum_{\tau^{\text{demo}}} \ln p(\tau^{\text{demo}} | w) \quad (2.27)$$

where $\mathcal{L}_{ME}(w)$ describes a convex, objective function. Due to its convexity, this allows for solving the problem using a gradient-based method [8]. The gradient is then given as the empirical feature counts from demonstration and the expected feature counts from the learner's policy:

$$\nabla \mathcal{L}_{ME}(w) = \mathbb{E}_{\pi^E}[\phi(\tau)] - \sum_{\tau} p(\tau | w) \phi(\tau) = \mathbb{E}_{\pi^E}[\phi(\tau)] - \sum_{x_i} D_{x_i} \phi(x_i) \quad (2.28)$$

Maximum entropy IRL works well for MDP problems, as it assumes that the state transition distribution is known [8]. However, this is often not the case in many robotic applications. Therefore, sampling-based or model learning extensions must be applied for problems where the model is unknown. In [64] such an approach was used for autonomous driving. The algorithm allowed for directly learning the reward function in the continuous domain, while also considering the uncertainties in the expert's demonstrations.

Arora, Doshi, and Banerjee [65] proposed a framework online IRL called incremental IRL. Based on this, max-entropy [63] was adapted such that it was generalized under occlusion. The results showed similar performance to the state-of-the-art applications of maximum-entropy IRL, while significantly reducing the computational time.

2.3. Reinforcement Learning

Reinforcement Learning (RL) enables learning through interaction. The main principle of RL has been represented in a flow diagram fig. 2.11. By interacting with the environment and observing the consequence of its action, an agent can learn to change its behavior based on the rewards it has received [66]. Based on this trial-and-error approach, the agent can autonomously discover an optimal behavior for a certain task [7], also called the policy. This policy will, using a RL approach, be iteratively updated, meaning that typically RL problems take place in an online setting. As the agent will learn of its environment, there is also no need for an external teacher. The need for the agent to interact with its environment, to learn, is the main difference between RL and IL [67]. This makes RL methods able to deal with changing environments [7].

While the principle of RL used for robotic applications is the same as the general approach, applying it comes with some challenges [7]. Firstly, learning takes place in a high-dimensional space (10 to 30-dimensional space). In addition, the states and actions are also continuous. This results in the learning being very costly. Secondly, it is unreasonable to assume that the states are completely observable and noise-free. Therefore, the problem often seems partially observable. Thirdly, the algorithm needs to be robust in respect to models that do not capture all the details of the real system, also referred to

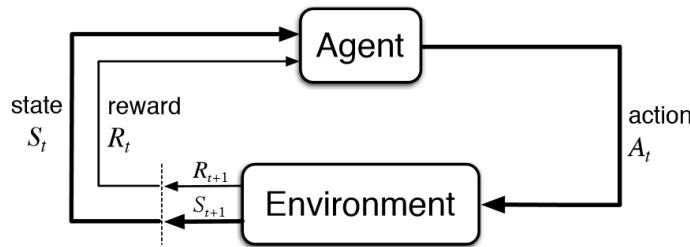


Figure 2.11: The agent-environment interaction in reinforcement learning as defined by Sutton & Barton [68]

as under-modeling, and to model uncertainty. And lastly, RL algorithms require a well-defined reward system. Therefore, it requires a fair amount of domain knowledge and may often be hard in practice. A solution for this problem is to use IRL to learn the reward function and RL to maintain the desired trajectory or to perform the desired task, once the learning has been completed. A description of this method was already given in section 2.2.3. The usage of IRL does have the downside that the advantage of RL compared to, IL of not being limited by the skill of a human expert, does no longer hold.

There are two widely used approaches which can be used to obtain the solution of a RL problem: the *value function* approach and the *policy search (PS)* approach [67]. In addition, a new approach has been developed called actor-critic. This approach combines the value function and PS method. The remaining of this section will first consist of explaining the general problem definition in section 2.3.1. Next, the three categories of RL methods will be discussed in sections 2.3.2 to 2.3.4.

2.3.1. Problem Definition Reinforcement Learning

In RL, the agent tries to maximize the accumulated reward (or minimize the cost) over a specific time. It does this by using a trial-and-error approach. This allows for the agent to learn tasks that could not be demonstrated or programmed by a human expert. The goal of RL is to find the optimal policy π^* , which maps the states (or observations) to actions, such that expected return J (defined by the reward function R) is maximized. Different expressions for the expected return have been defined. The main types of function indicate that the agent's return can be either discounted or averaged, and are calculated on a finite or infinite horizon [67]. In discount, RL the choosing of the discount factor is critical. Choosing this value too low can result in an unstable behavior, making it often not suitable for robotic application. Average rewards have the problem that it cannot distinguish between policies that initially gain a transient of larger rewards and those that do not. However, this disadvantage is far less critical than one of the discount factors, thus making the average reward setting often more desirable for robotic applications. The average reward return is then defined as:

$$J(\pi) = \sum_{s,a} \mu^\pi(s,a)R(s,a), \quad (2.29)$$

where μ^π is the stationary state distribution generated by policy π acting on the environment, i.e. the MDP. As classic RL approaches are based on the MDP, in theory, they should only work when it adheres to the Markov structure. However, it has also been discovered that in reality the approaches often still work when this is not the case [7]. The finding of a policy π can now be described as an optimization problem for $J(\pi)$. This problem has been formulated as an optimization problem

$$\max_{\pi} J(\pi) = \sum_{s,a} \mu^\pi(s,a)R(s,a), \quad (2.30a)$$

$$\text{s.t. } \mu^\pi(s') = \sum_{s,a} \mu^\pi(s,a)T(s,a,s'), \forall s' \in S, \quad (2.30b)$$

$$1 = \sum_{s,a} \mu^\pi(s,a), \quad \mu^\pi(s,a) \geq 0, \forall s \in S, a \in A. \quad (2.30c)$$

Here, eq. (2.30b) defines the stationarity of the state distribution π , which ensures that it is well-defined, and eq. (2.30c) ensured the proper state-action probability distribution. To search for a solution originally two approaches were used. Firstly, the optimal solution can be searched directly in its original, primal problem, which is known as PS. Secondly, it can be optimized using the Lagrange dual formulation, which is known as a value-function-based approach. In addition, in recent years, a third approach has been developed: the actor-critic approach.

Besides the types of algorithms used, there are two more criteria on how to categorize RL problems: model-based or model-free, and on-policy or off-policy [5, 69]. An overview of the taxonomy of RL methods is given in fig. 2.12.

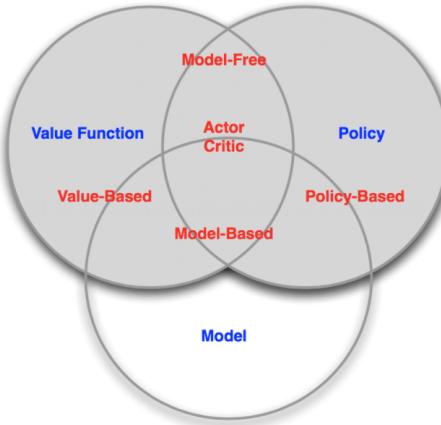


Figure 2.12: Taxonomy of RL methods [70].

Just like in most types of learning methods, both model-based and model-free methods exist for RL. Here the model-based methods rely on the model of the environment which is either known or can be explicitly learned to use the algorithm. In contrast, do the model-free methods not depend on such a model [7]. The advantage and disadvantage were described by Polydoros and Nalpantidis [67] and are shown in table 2.3.

The difference between *on-policy* and *off-policy* method is that the former tries to evaluate or improve the policy to make decisions, whereas the latter does it to generate the data [7]. On-policy require the agent to interact with the environment, meaning the policy that interacts with the environment should be the same as the one which has to be improved. In the off-policy, this is not necessary the case. For example, the experience of other agents can also be used to improve the policy.

In the remaining of this section the different approaches (value-function, PS and actor-critic) will be discussed. Here both a general notation and some classes of methods will be described.

2.3.2. Value-function approaches

This approach tries to find the optimal policy, by iteratively optimizing the value function. It assumes that each state has an associated value dependent on the reward achieved in that state. Each state also has a potential, which depends on the future reward achieved using the agent's current policy [71]. As it is often the case that a reward is only received in a specific goal state, the positions closest to this state have the highest value. By breaking down the problem into sub-problem, following the sequence of states, the value function can be found.

Value-function methods try to find a solution to the condition for optimality. This depends on the Lagrange multipliers \bar{R} and $V^\pi(s')$ and is given by

Table 2.3: Advantages and disadvantage of model-free and model-based RL methods [67]

	Advantages	Disadvantages
Model-based	<ul style="list-style-type: none"> • Small number of interactions between robot & environment • Faster convergence to optimal solution 	<ul style="list-style-type: none"> • Depend on transition models • Model accuracy has a big impact on learning task
Model-free	<ul style="list-style-type: none"> • No need for prior knowledge of transitions • Easily implementable 	<ul style="list-style-type: none"> • Slow learning convergence • High wear & tear of the robot • High risk of damage

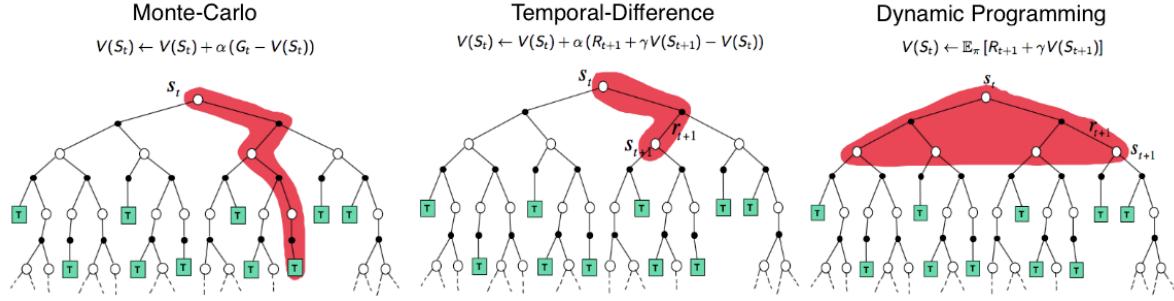


Figure 2.13: Comparison of the backup of three value-function approaches: MC, TD learning, and DP [72]. The hollow circles represent the state value, the closed circles the state-action or action values, and the lines the actions. The tree visualized in the figure represents an example of an entire action-state space. The red part represents the part examined by the method. This thus results in two directions: width and depth. The width represents the backup [66], where either a sampling of the actions is done (MC and TD) or expectation of all choices is determined (DP). The depth represents the bootstrapping, which goes from one-step TD to n-step TD learning, where the pure MC approaches do not use bootstrapping at all.

$$V^*(s) = \max_{a^*} \left[(R(s, a^*) - \bar{R} + \sum_{s'} V^*(s') T(s, a^*, s')) \right], \quad (2.31)$$

where $V^*(s)$ is the shorthand notation of $V^{\pi^*}(s)$. This formulation is equivalent to the *Bellman principle of optimality*. This approach owes its name to the value function $V^\pi(s)$. Instead of using this function, another option is to use the state-action value $Q^\pi(s, a)$. As its name suggests, in addition to the state (where the value function depends upon), does this value also depend on the action. The different methods which have attempted to estimate $V^*(s)$ or $Q^*(s, a)$, can be divided into three classes: Dynamic Programming (DP) based optimal control approaches, rollout-based MC methods and Temporal-Difference (TD) methods. A graphical representation of these methods is given in fig. 2.13. In the remainder of this section, these classes will be discussed. Before doing so, the idea of bootstrapping should be explained. You are using bootstrapping, are you make an estimation based on an earlier made estimation [66].

Dynamic Programming based methods

Dynamic Programming (DP) based methods, use a model of the transition probabilities $T(s', a, s)$ and the reward function $R(s, a)$ to calculate the value function, making it a model-based approach. Two well-known approaches are policy iteration and value iteration [7]. A graphical representation of both methods is shown in fig. 2.14.

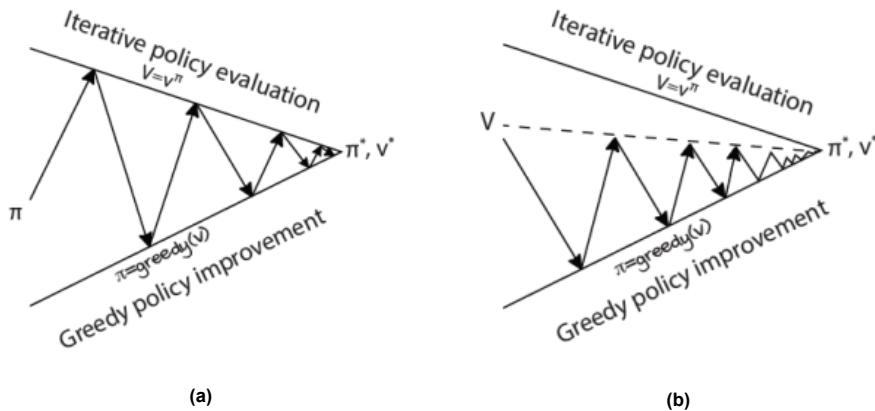


Figure 2.14: Graphical representation of the DP methods [73]. With (a) the representation for policy iteration and (b) the representation for value iteration. The greedy policy improvement uses the policy which seems to be the best in the short term (after one iteration).

Policy iteration starts with initializing an arbitrary policy [7]. After this, it alternates between two phases: policy evaluation (state or action function is calculated for a given policy) and policy improvement (the best action for each state is derived) [67]. In the evaluation phases, the value function for the current policy is evaluated. Here, each state is visited, and its value is updated based on the current value estimates of its successor states, the associated transition probabilities, as well as the policy. This procedure is repeated until the value function converges to a fixed point, which corresponds to the true value function. Next, in the improvement phases, an action is greedily selected in every state according to the value function. These two steps are repeated until the policy does not change any longer [7]. A disadvantage of policy iteration is that the iteration can be a time-consuming task [67].

In contrary to the policy iteration method, where the policy is only updated after the evaluation step has converged, the value iteration method does not wait until the convergence of the evaluation procedure for updating the policy. The value function is updated this based on eq. (2.31) every time a state is updated [7, 67].

Monte Carlo

In contrast to the DP approach, is this method model-free [74]. It can be used to replace the policy evaluation step of the value or policy iteration method [7]. The value function V^π can be estimated by randomly sampling many trajectories starting from state s , according to the given state transition matrix P [74]. This is done by keeping track and using the frequencies of the transitions and rewards. In the paper of Ding et al. [74] an example of this method is given. MC perform roll-outs by executing transition function, in other world operating on-policy [7].

Temporal-Difference

This method combines the ideas from DP and MC to estimate the value function [74]. Like DP, TD uses bootstrapping in the estimation. This is used to form a target from the observed return and an estimated state value for the next state [74]. On the other hand, just like MC, a sample approach is used to estimate the value function. To be precise, sample transition in the MDP are used [7]. To learn, TD uses the error (difference between target value and estimated value) at different time steps [74]. A basic TD method for the update is

$$V' = V(s) + \alpha (R(s, a) - \bar{R} + V(s') - V(s)), \quad (2.32)$$

with $V(s)$ the old estimate of the value function, $V'(s)$ the new one, and α the learning rate. This method is also referred to as the TD(0) or one-step TD method [7, 74], as it looks only one step ahead. N-step TD can also be developed by extending the target value with discounted rewards in the N-step future and estimate state value at the N-th step. Famous TD methods are SARSA, which is on-policy, and Q-learning, which is off-policy.

2.3.3. Policy Search

A problem with value function approximation is that it results in a difficult problem in high-dimensional state and action spaces, which can be problematic for robotic applications. Another problem is that value functions are often discontinuous, which results in the propagation of error in the value function through to the policy [75]. Therefore, policy search (PS) methods could provide a good alternative. In these methods, the optimal policy is learned directly [67]. Generally, they exist of a set of parameters, which has to be optimized such that the cumulative reward is optimized. Where value function searches in the state-action space for the best policy, does PS directly search in the parameter space. This allows for scaling RL into high-dimensional continuous action space, by reducing the search space of possible policies. Which, in other words, means that it is less computational expensive. PS has been seen as a more difficult approach, in comparison to value function, as an optimal solution cannot be directly obtained using eq. (2.30). However, PS has retrieved popularity in the field of robotics, due to better scalability and the convergence problems of the approximate value-function method [7].

Most PS methods optimize locally around existing policies π , parameterized by a set of policy parameters θ_i [7]. This is done by computing changes in the policy parameters $\Delta\theta_i$ which will increase the expected return and results in iterative updates of the form

$$\theta_{i+1} = \theta_i + \Delta\theta_i. \quad (2.33)$$

The key step is the computation of the policy update. Looking at the different policy update strategies, a distinction between the two evaluation strategies, the step-based and episode-based evaluation strategies, should be made. Step-base exploitation uses an exploratory action for each time step, whereas episode-based exploration changes the parameter vector θ only at the start of each episode [75]. The policy update strategies for both model-free and model-based methods are often based on policy gradients (PGs), EM, or information-theoretic insights (Inf.Th.) [7, 67, 75]. In most cases, these methods can both be applied for step-based and episode-based exploration. They will further be explained in the remaining of this section.

Policy gradient

The gradient-based approaches, use a gradient ascent for maximizing the expected return J_θ [7, 75]. Here, the policy is updated following the gradient of the expected return $\nabla_\theta J_\theta$ for a defined step size α . For iteration $i + 1$, the policy is given by

$$\theta_{i+1} = \theta_i + \alpha \nabla_\theta J_\theta. \quad (2.34)$$

The policy gradient can then be computed using

$$\nabla_\theta J_\theta = \int_\tau \nabla_\theta p_\theta(\tau) R(\tau) d\tau. \quad (2.35)$$

Different approaches exist to compute the gradient, and many algorithms also require tuning the step-size α . Three approaches to compute the gradient are: finite difference gradient, the likelihood ratio method, and natural gradient [75].

The *finite difference gradient* is the simplest approach and is typically used in an episode-based setting [75]. This method computes the gradient by applying a small perturbation $\delta\theta_p$, to the parameter vector θ . By using linear regression, the gradient can now be estimated. While this is a straightforward approach and even applicable for non-differentiable policies, it is often considered to be noisy and inefficient.

The second approach uses the likelihood ration trick, to compute the gradient of the episode distribution. By inserting the gradient of the episodic distribution $\nabla_\theta p_\theta(\tau) = p_\theta(\tau) \log p_\theta(\tau)$, in eq. (2.35), the PG can be computed. Here $p_\theta(\tau)$ is approximated using the sum of the sampled trajectories τ . Algorithms using this approach are REINFORCE and G(PO)MDP [75].

As convergence of the finite difference gradient and likelihood ratio method seemed to be kind of slow, the *natural gradient policy* approaches, which allow for fast convergence, might be advantageous for robotic application [7]. The idea of this method is to limit the distance between two subsequent distributions, to ensure stability [75]. The main difference with traditional gradient methods is that they quickly reduce the variance of the policy and, thus, stop exploring. In contrast, natural gradients only gradually decrease the variance, due to their limit in distance, which in the end results in finding the optimal solution faster.

Expectation Maximization

PG method often requires the setting of a learning rate, which can be problematic and result in unstable learning or slow convergence. This problem can be avoided by formulating the PS as an inference problem with latent variables and using the Expectation Maximization (EM) algorithm to derive the new policy. EM algorithms do the parameter update by computing it as a weighted maximum likelihood estimate, which has a closed-form solution for most policies [75]. This means that no learning rate is required. Some of the approaches which have been proven successful in robotics are: reward-weighted regression, policy learning by weighted exploration with returns, MC EM, and cost-regularized kernel regression [7].

Information-theoretic insights

This last approach tries to update the parameters by "staying close" to the provided data [75]. This means that the trajectory distribution after the policy update should not jump away from the trajectory distribution before the policy update. Therefore, information-theoretic insights (Inf.Th.) approaches bound themselves by the old trajectory distribution $q(\tau)$ and the newly estimated trajectory distribution $p(\tau)$, at each update step. This limits the information loss of the updates, and thus avoids the new distribution convergence to a local optimum. Some methods which have applied this idea, are the natural PG algorithm and the Relative Entropy PS.

2.3.4. Actor-critic method

The actor-critic (AC) approach forms a special class of PG methods, as they aim to combine the advantage of both the value function and PS approaches. The idea behind this method is to separate the policy (actor) and value function (critic) as visible in fig. 2.15. Both of these entities are represented as a parametrized function of the state x , with the parameters ψ and ϑ , respectively. At time step t , the actor calculates an action u_t based on the current state x_t and applies it to the system. This results into the transition of to a new state x_{t+1} and a reward r_{t+1} . Next, the evaluation of policy can be done using any policy evaluation method. The method explained in [76] use TD for this. The TD-error δ can be calculated based on the new state and reward. This error is used to estimate the actual value function and criticize the actor, by updating the corresponding parameters. If $\delta > 0$, the action is favorable, which means that the actor should go toward the direction. In contrast, if $\delta < 0$, the actor should avoid taking that action in the next step.

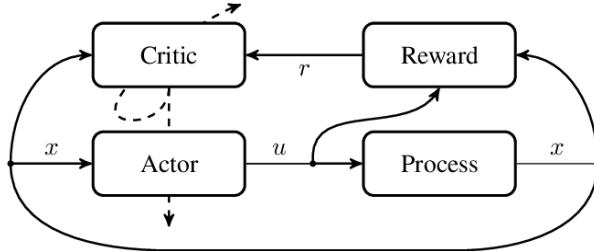


Figure 2.15: Schematic overview of an actor-critic algorithm [76].

Algorithm 2 Actor-critic Algorithm [77]

```

1: Initialize  $\lambda, \gamma, \alpha_a, \alpha_c$ 
2: Initialize  $\vartheta_0, \psi_0$ 
3: for each trial do
4:   Initialize  $x_0$ 
5:   Generate a random initial action  $u_0$ 
6:   Initialize eligibility trace  $\zeta_0 = 0$ 
7:    $t \leftarrow 0$ 
8:   repeat
9:     generate exploration  $\Delta u_t$ 
10:    calculate current action  $u_t = \hat{\pi}(x_t, \psi_t) + \Delta u_t$ 
11:    apply  $u_t$ , measure  $x_{t+1}$ 
12:    receive reward  $r_{t+1} = \rho(x_t, u_t)$ 
13:     $\delta_t = r_{t+1} + \gamma \hat{V}(x_{t+1}, \vartheta_t) - \hat{V}(x_t, \vartheta_t)$            ▷ calculated TD-error
14:     $\zeta_{t+1} = \lambda \gamma \zeta_t + \frac{\partial \hat{V}(x, \vartheta)}{\partial \vartheta} \Big|_{x=x_t, \vartheta=\vartheta_t}$            ▷ eligibility trace
15:     $\vartheta_{t+1} = \vartheta_t + \alpha_c \delta_t \zeta_{t+1}$            ▷ update critic parameter
16:     $\psi_{t+1} = \psi_t + \alpha_a \delta_t \Delta u_t \frac{\partial \hat{\pi}(x, \psi)}{\partial \psi} \Big|_{x=x_t, \psi=\psi_t}$            ▷ update actor parameter
17:     $t \leftarrow t + 1$ 
18:   until  $t = \text{maximum number of samples } T_s$ 
19: end for
  
```

Algorithm 2 shows the algorithm of the actor-critic method. While most points can be connected to the image in fig. 2.15, some additional explanation is required to understand this algorithm. First, some coefficients are initialized, which represent the trace decay rate $\lambda \in [0, 1]$, discount factor γ , the learning rate for the actor α_a , and the learning rate of the critic α_c . The next line initializes the actor and critic parameters. To increase the probability of visiting all the system states multiple times, an exploration signal Δu_t is added to the output of the actor. A standard way to update the critic is to use the TD-error δ . If the value function is approximated to be linear $V_\theta(x) = \theta^T \phi(x, u)$, then the TD-error can be written as $\delta_{t+1} = \delta_t + \alpha_{c,t} \delta_t \phi(x_t)$. This is also referred to as *TD(0)*, as no eligibility traces are used. In algorithm 2, however, the extension is shown where an eligibility trace ζ_t is used (line 14). This causes the critic to be updated based on a series of visited states instead of just the prior one, which should have the advantage that the learning is speed up. Finally, to update the parameters of the critic and actor, the TD-error is used.

Grondman et al. [76] defined two criteria for the taxonomy of actor-critic methods. Firstly, just as in other types of RL methods, a distinction between discounted return and average return can be made. In addition, two types of gradients have been identified: standard and natural.

2.4. Combining Imitation Learning and Reinforcement Learning

IL and RL both provide features that make them desirable for robotic teaching. As stated before, IL can fast learn a certain skill, however, it is limited by the performance by the skill of an expert. This problem does not occur for RL, however, its learning time is often much slower. The combination of these two methods would therefore be preferable. In this case, both best features could be combined: use an expert to quickly learn a reasonable policy by imitation and exploring how to improve the skill upon the expert by using RL.

Chang et al. tried to implement the Locally Optimal Learning to Search (LOLS) algorithm, which combines IL and RL by stochastically interleaving incremental RL and IL updates [78]. For doing so, the IRL framework L2S was used in combination with a policy-iteration based RL method. As a result, the learned policy would either perform as well as the expert policy (due to IL method) or reach a local optima (due to RL method). While it is possible that the local optima result in a better performance than the experts' policy, it is difficult to quantify this.

In [79], the idea of combining IL and RL was implemented via the idea of reward shaping. For the IL part, access is assumed to be a cost-to-go oracle, which provides an estimate of the expert's cost-to-go during training. The principle of cost-to-go (or reward-to-go) is that a cost is only received after a certain action has been taken. After shaping the cost, the planning horizon of the new MDP would be truncated, followed by a search for a policy that optimizes over the truncated planning horizon. While Sun et al. (2018) did implement their approach in a simulation, they did not use it in the real world.

An example where a combined IL and RL approach was used for policy learning, is [80]. For one of their experiments, Celemín et al. found an initial policy using the IL method ProMP. From where, they used, among other things, PS to further learn the policy. In another experiment, they did not compute this initial policy. Comparing the results shows a significant improvement while using the method which initially computes the policy using ProMP. This experiment will further be expanded on in section 3.3.2.

While research has shown promising results in combining IL and RL, it is still quite limited. Further research should be done comparing IL and RL to the combined approach.

2.5. Comparison

In this section, a comparison of the different learning approaches described earlier in this chapter, will be made. This comparison is shown in table 2.4. Here, the different robot teaching methods are compared based on multiple criteria, which will be expanded on in the remaining of this section. It should be noted that the combined IL with RL is not taken into account in this comparison, as not enough information could be found to give a comprehensive result for this method. The comparison is a summary of information found in literature and assumptions based on this. The used literature will be described in the remaining of this section. To have a more fact-based comparison, further research should be conducted.

Table 2.4: Comparison of robot teaching methods. It should be noted that this comparison is based on an interpretation of literature. To have a more conclusive analysis, research is required. The methods are compared based on some criteria, here + indicates a high score relating to the criteria, +/- an average score, and - a bad score. These scores were given with respect to each other. From left to right the different criteria are; *Computational efficiency* describing the amount necessary to learn. *Stability* describes the ability to deal with perturbation. *Smoothness* refers to the smoothness of the trajectory. *Generalizability* describes the ability of the system to deal with the situation not occurred during the learning process (e.g. new starting point or goal). The difficulty of implementing a certain system is indicated by *implementation*. *High-dimensional input* states whether a system can deal with a large amount of data. Lastly, *online* describes whether the system was originally implemented for online usage.

		Method	Comp. efficiency		Smoothness	Generalizability	Implementation	High-dim. input	Online
IL	RL		Stability						
BC	Value function	+	+/-	+	+	+	+/-	-	
		+/ - ²	+/ - ²	+	+	+/ - ¹	+/ - ¹	+	
IRL	Policy Search	+/ - ¹	-	+/ - ¹	-	-	-	+	
		-	-	+/ - ¹	-	-	+	+	
		-	+/ - ¹	+/ - ¹	+/ - ¹	-	+/ - ³	+	

¹Assumption based on information of IL and RL methods.

²Assumption based on the information of RL methods.

³Assumption based on the information of PS method.

The *computational efficiency* describes the amount of data and computational time, required to converge to a solution [13]. Especially for real-time applications, this is of the highest importance as a high computational cost would result in a slower system, thus making this system less reactive. Slow learning might not necessarily be a problem when learning in a "save" environment, it could be a problem when adaptation to the environment is required. IL method often have a higher computational efficiency than RL. This is due to the fact that fewer iterations are needed to obtain the policy. RL on the other hand, use a trial-and-error approach, resulting in a much more iteration. Looking at the different IL method, BC often only have to do one iteration per demonstration. There is, however, a difference in computational time between the different BC methods. For example DMP just have to update one parameter [8]. Whereas, GP has to compute a new covariance matrix (which can be large for high-dimensional inputs) and has to invert this [42]. As IRL typically have to update their cost function multiple times, their sample efficiency is slightly worse than those of the BC. Looking at the RL methods, the value function methods are known to have a high sample efficiency in comparison to the PS method [5]. This could be translated into the former having a higher computational efficiency than the latter. As actor-critic methods use the bootstrapped value function to reduce the variance of the gradient estimate, they gain the sample efficiency of the value function.

Stability refers to the ability of the system to deal with perturbations. In cases where perturbation would result in a big difference in the behavior, the system is said to be unstable. Instead, it is desired that the behavior only changes a small fraction or maybe not even at all. The reason for this is to ensure a safe design, which is of importance to avoid a robot damaging itself, and something or someone else. The stability of the BC methods is not always ensured. For example, DMP is statically asymptotically stable [81], meaning it has a stable attraction to a target position, whereas ProMP [8], GMR and GPR [82] do not provide this guarantee. As RL methods generally search a large part of the solution space to find a policy, one would expect good stability. However, this is not necessarily the case as most methods depend on multiple conditions. PG requires the setting of a learning rate, which, if chosen incorrectly, results in unstable behavior [75]. This instability problem was, however, tackled by formulating PS as an inference problem with latent variables and, using the EM algorithm, to determine a new policy. Another example is the Inf.Th. principle, which states that in order to provide a stable behavior, the distance between the old and new trajectory distribution, should be bounded. In [71], a solution to increase the stability was presented - replay. During replay, the robot is initiated in states which are not stored, which results in improvement of the exploration efficiency, and stability. This idea of replay

is not necessary for the actor-critic algorithm A3C to ensure stability, as it allows for multiple agents to asynchronously explore different policies in parallel [71].

With *smoothness*, the smoothness of the encoded behavior (often trajectory) is meant. To minimize the risk of damage to the robot, it is desired that the behavior is smooth, without sharp changes [13]. This criterion is hard to quantify as it has to be visually examined. Therefore, results become quite subjective. For the IL approaches, the different smoothness was based on literature. For each approach, multiple papers were used (for BC [32, 38, 42, 51, 83, 84] and for IRL [62, 85]). As IL methods try to imitate the expert, which is assumed to have a smooth behavior, the resulting behavior of the robot is in most cases also smooth. The same could not be said for RL method. Here, the smoothness of the system highly depends on the chosen reward function. Just as for stability, the usage of a bad trajectory would result in rough behavior.

Generalizability is the ability of a system to perform well in similar environments as used for the training of the system [60]. Examples of this are states and action unobserved in the demonstration, starting the task at a different initial state, and reusing skills in different problem settings [5]. BC method often generalize quite well. However, in some cases, some extensions are required [8]. An example is ProMP, which learns the distribution of demonstrated trajectories in parameter space. To generalize to a new start and goal position, the learned distribution has to be conditioned. An example of this was given in [84]. Research has also proven IRL method to generalize well to un-encountered settings [63]. Generally speaking, RL often has a lack of generalizability [71]. A solution could be to increase the sample dataset and train on multiple (random) environments [86]. This, however, comes at the cost of computational efficiency. A difference can also be seen between model-based and model-free approaches. Where the former often generalizes better compared to the latter [71]. There have been PS algorithms that generalize relatively well. Most of these are episodic-based policy evaluation strategies [75]. Even though some examples can be found in the literature on the generalizability of RL, applying this is not frequently done.

It is desired that a system is easy to implement. Difficult *implementation* can mean that a lot of time needs to be spent to implement the system, and in addition, that more can go wrong implementing it. For example, more assumptions have to be made. Overall, BC methods are quite easy to implement. The reason for this is they are often model-free [8]. In contrast, IRL methods are often model-based, meaning that if no model of the system is available, implementation can become a lot more difficult. Implementing RL methods, is often more difficult than IL. In addition to the difficulty of determining a reward function and a model of the environment, running RL methods can take multiple iterations, which is quite costly [7]. A solution for this is to first run the system in simulation. This can however not completely replace real-time learning, as even small differences in the environment can result in totally different behaviors.

The ability of BC method to deal with *high-dimensional inputs* depends on the chosen algorithm. Both DMP and ProMP are often not able to do so, whereas GMM, GPR and HMM are [40]. As RL often uses a continuous stream of data, the amount of data becomes large. This makes it often quite difficult for RL method to deal with robotic learning. Nevertheless, PS has proven to be able to scale well with high-dimensional state space [5]. The same cannot be said for value function methods.

The last criterium relates to *online* learning. For each method, an online implementation can be found. However, these implementations often require some adaptation to the original method. Therefore, this criterium states whether methods were originally meant for online implementations. BC method usually obtained the demonstrations first, after which a model was computed. Recent literature can be found on online implementations of BC methods [26, 30], but as stated they require some adaptations of the original method. IRL uses demonstration to find a reward function, and use this to compute a policy (using some RL method) [60]. As the updating of the policy is done iteratively, this method can also be stated as an online learning process. In contrast to the IL approaches, offline learning using RL not common [71]. Instead, RL is in most cases used in an online setting, where the policy is iteratively updated based on the sensory feedback [7].

2.6. Summary

In this chapter, the main goal was to give an overview of the existing methods for robot teaching. In this section, the findings will be summarized.

- **IL and RL:** Two main approaches have been identified, which both have the benefit that the behavior does not need to be manually programmed. IL uses demonstration, exerted by an expert, to learn skills. Which makes learning easy and natural. In addition, these methods are often computationally efficient as only a small amount, or in some cases just one, has to be done to find a solution. However, due to these demonstrations, the skill of the robot is limited by the skill of the expert. Often, IL also has to deal with the correspondence problem, which is due to the different embodiment of the learning and the expert. Both of these limitations do not count for RL, as it does not require an expert teacher. However, RL is often more difficult to specify, and choosing an incorrect reward function can have tremendous results. In section 2.5, a comparison between these two methods has been given.
- **BC and IRL:** IL problems can be categorized into two main approaches. The BC approach tries to optimize the policy, by directly mapping. Most implementations are quite simple and computationally efficient. However, often it only tends to work when a large amount of data is required, due to a compounding error caused by covariant shift. This is not a problem when using IRL. In addition, IRL also, performance well in a suboptimal space, which is not the case for BC. This comes at the cost of computational efficiency.
- **RL methods:** Value function tries to find an optimal policy, by iteratively optimizing the value function. A disadvantage of this approach is that it is often difficult to apply in high-dimensional state and action space, which is problematic for robotic applications. Therefore, PS can be a good alternative. This method directly learns the optimal policy. However, this method can be a bit more difficult to implement. A special type of PS method is AC. This method aims to combine the advantages of the value function and PS approach.
- **Combining IL and RL:** The idea behind combining IL and RL is to use both of their advantages in one approach: the computational efficiency of IL and the ability to exert a human skill by using RL. While research seems to be promising, the literature combining these approaches still seems to be quite limited.

3

Continuous Skill Learning

Chapter 2 focused on what methods can be used to teach a skill to a robot, either by imitation of an expert or by learning based on a (predefined) cost/reward function. The focus of this chapter will be on how the usage of these methods changes the skill of a robot throughout the learning process. As this paper focuses on online learning methods, it is assumed that the behavior which the robot will exploit at the start of the learning process, will not be the same as at the end of the learning process, i.e. it will incrementally change.

The remaining of this chapter consists of an introduction to different learning outcomes in section 3.1. Next, the effect of skill learning, using different state inputs, will be examined in section 3.2. Section 3.3 is dedicated to methods using a reward/cost function to compute a policy. Lastly, the findings will be summarized in section 3.4.

3.1. Learning Outcome

To learn a certain task, the learning outcome needs to be defined. In [13], three learning outcomes for Imitation Learning (IL) problems were defined: a policy, a reward or cost function, and a plan. The choice of the learning outcome depends on the task and the associated constraints. The first two learning outcomes have already been discussed in chapter 2, i.e. the computing of Behavioral Cloning (BC) and Reinforcement Learning (RL) method of a policy and the computing of Inverse Reinforcement Learning (IRL) methods of a reward/cost function. The latter, however, requires some additional explanation. This method includes learning at the highest level of task abstraction. Here it is assumed that the task is performed according to a structured plan, made up of several sub-tasks of primitive action. These plans typically encode the pattern and constraints in the sub-tasks or primitives and take the robot from a beginning to a goal state. For the remainder of this research, this specific topic will not be discussed separately. The reason for this is that these methods, often separate their task into multiple skills which then can be described using a skill's specific policy [87].

Looking at policy as a learning outcome, Ravichandar et al. [13] defined different taxonomies to classify the different methods: policy input, policy output, and policy class. Policy input involves identifying the appropriate input to the policy, an example being a state. The policy output is similar to the abstraction defined in section 2.2. Policy class refers to the class of mathematical function to which the policy belongs. As the goal of this chapter is to show the continuous learning of robots, using this as a classification might be a bit too abstract. Mainly, as functions do not give a clear idea of how the different implementations hold together. It was chosen when relating to policy, that comparing different inputs would seem the most logical choice. A reason for this is as the classification would be the easiest.

Table 3.1: Overview of references using different types of state input. For the method, the encoding and learning method is meant. The demonstration and task which are named in this table represent the relevant information for this chapter. In some cases, other demonstration types and/or tasks were also executed.

Author & year	Ref.	State input	Method	Demonstration	Task
Havoutis et al. (2017)	[89]	Position	TP-HSMM	Kinesthetic teaching	PiH
Calinon et al. (2007)	[38]	Joint Angle	GMM + GMR	Shadowing & Kinesthetic teaching	Gesture imitation
Pastor et al. (2011)	[83]	Force	DMP	Kinesthetic teaching	Grasping
Nemec et al. (2013)	[30]	Force	DMP + ILWR	Teleoperation	PiH
Peternel et al. (2017)	[26]	Force	DMP + LWR	Shadowing	Sawing
Peternel et al. (2016)	[90]	Torque	DMP + LWR	Shadowing	Periodic movement.

3.2. State input

Different state inputs for the learning algorithm have been identified [13]. The different types of state input which will be examined in this section are position, joint angle, force, and torque. While there are more state inputs to be found in literature, these are the ones most frequently used. An overview of the different papers which will be discussed in this section is given in table 3.1.

Before elaborating on the continuous learning of robots, based on the state input, the term state-input should be defined. In the learning scheme of the policy, there are often two segments that make use of some type of input: the policy (or controller) and the robot. For the classification of the different policy state inputs, as only online learning will be considered, the control scheme will either be a complete close-loop or it will close-loop with an additional feed-forward loop [88]. In the latter one, the feed-forward loop will, in most cases, consist of one demonstration. This is often the case for approaches using Dynamic Movement Primitive (DMP) [30]. This means that the state input for the policy and robot will not be the same. In this section, the state input is classified as the state input used for the policy.

3.2.1. Position

Position is often used for state input for policy learning using IL methods [13]. There are multiple ways to require this type of data. One method is to use the embedded sensors in the robot. This is especially effective when using kinesthetic teaching for the requirement of the data. However, often the data should be required from observing a (human) expert. In these cases, extracting the information can be done using vision and/or motion sensors [50, 91]. Most state-of-the-art approaches consider vision, as it seems to be the most natural way to capture information from the human. However, these systems must deal with problems such as occlusion, appearance change, complex human-robot kinematic mapping, and the encoding of the video material. Therefore, using a motion sensor can be a more practical approach. This allows for the tracking of the motion preciser and a simple mapping can be established.

An example, where the position was used as the state input, is in [89]. Havoutis and Calinon proposed a framework for supervised teleoperation with online learning. By using an online Bayesian non-parametric learning algorithm, the motion of the robot could be presented online by task-parametrized Hidden Semi-Markov Model (TP-HSMM). These motions can be executed autonomously using a model predictive control (MPC) approach. In the paper, the principle is illustrated using a simplistic example in a planar environment, as visible in fig. 3.1. By means of teleoperation, three demonstrations were done. This result is a model with six Gaussian components. Once it is the remote robot's turn, it is visible the system is able to create a path to the goal from a new initial position, which means that it is able to generalize. Here, it should be noted that the starting point is indicated by the demonstrator. A graphical representation of this idea has been presented in fig. 3.2. In addition to its generalizability capabilities, also smoothness is visible.

While the planar example gives a clear overview of the main principle of this approach, the eventual goal of the paper [89] was to execute a PiH task. In this experiment, the left arm of a robot is used to kinesthetically demonstrate the behavior, and the right arm to remotely reproduce it. The experimental

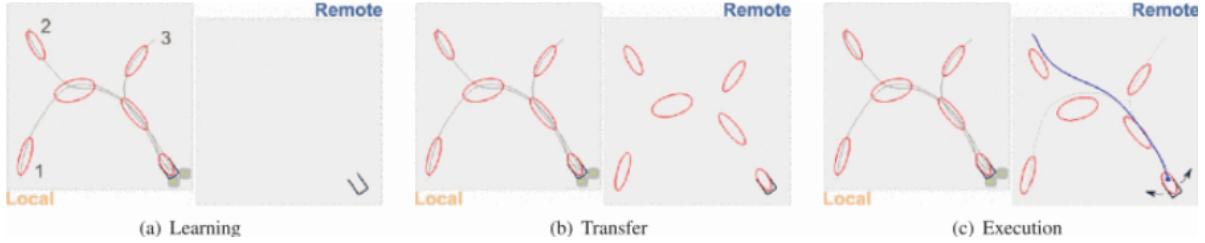


Figure 3.1: Planar example of learning and control with the TP-HSMM [89]. Here, in each plot, left represents the local system (operator's side) and right the remote system (robot side). (left) represents three motions and six corresponding Gaussian (here Gaussian are added after each demonstration). (middle) shows the transfer from the local system to the remote system. (right) shows a new motion which is generated by the MPC controller. It should be noted that this can only be created once the operator has initialized a new initial state.

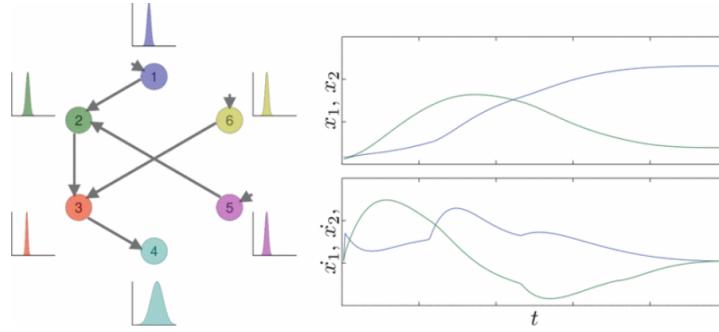


Figure 3.2: Planar example [89] with (left) the graphical representation of the transition matrix and duration probabilities for each state, where the three paths are formed according to the three demonstrations. And (right) the executed motion on the remote side starting from a new position, using TP-HSMM and the MPC controller.

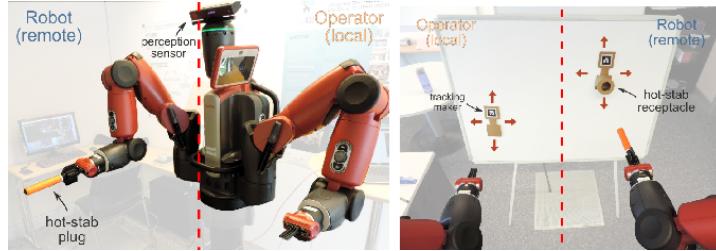


Figure 3.3: Experimental setup of peg-in-hole (PiH) task [89]. The left arm is used by the operator to provide demonstrations and directly teleoperate the right arm. During demonstrations, a TP-HSMM model is learned. This can be used to execute the PiH tasks autonomously using the right arm. Here, the robot locates the new goal positions using an RGB-D sensor.

setup is shown in fig. 3.3. The model is updated online, using five trials. In the reproduction phase, to find the goal position, an RGB-D sensor was used. This showed the ability of the system to generalize well, even when the starting and/or goal position was much different from demonstrated data. A reason why this worked well was the availability of information about the environment. For example, if no RGB-D sensor was used, the new goal positions would not be found, hence it would not generalize well.

3.2.2. Joint angles

Joint angles are used to compute the position of an end-effector. Making the idea behind using joint as a state input quite similar to using positions. A difference can be noticed in the variability in the configuration. Where position often only defines the position of the end-effector, does the joint angles define the entire configuration of the robot. Therefore, the solution space using joint angles reduces compared to using position. While this would result in a more stable system, this also has the disadvantage that it is less able to adapt to changes in the environment and is thus less safe. Another main problem arising with the usage of joint angles as state input is the correspondence problem, as the joints of an

expert might not be the same as those of the robot. However, this can be avoided by using kinesthetic teaching for demonstration. For these last cases, the embedded sensors of the robot can be used to obtain the state inputs for the learning algorithms. However, for the other situation (e.g. teleoperation or observation), external sensors are required. Just as for position, vision and motion sensors are a good option for this.

In [38], a wearable device was attached to the human expert, which included eight sensors that each provided its absolute orientation. The goal of this research is to transfer the motion from a human to a (small) robot. This consisted of two steps: an observation step and a kinesthetic teaching step. In this first step, the recorded joint angles are projected in the latent space and encoded using Gaussian Mixture Model (GMM). The training of the parameters is done using Expectation Maximization (EM). In the second step, kinesthetic teaching is used to refine the motion of the robot. The updating of the GMM is done by estimating online the new incoming data and the previous estimation of the GMM parameters. One of the approaches used for updating the GMM is the direct update approach, which uses an adapted version of EM. Results for a specific gesture are shown in fig. 3.4. The usage of both the motion sensor and kinesthetic teaching allows for dealing with the correspondence problem.

In addition to the direct updating method, another method was used to update the model: the generative method. Both methods were compared to an offline learning approach: batch learning. The results for three different gestures are visualized in fig. 3.5. For the reproduction of the motion GMR was used. The results show that each method was able to capture the essential characteristic of the demonstrated data. In addition, the results seem to be quite smooth.

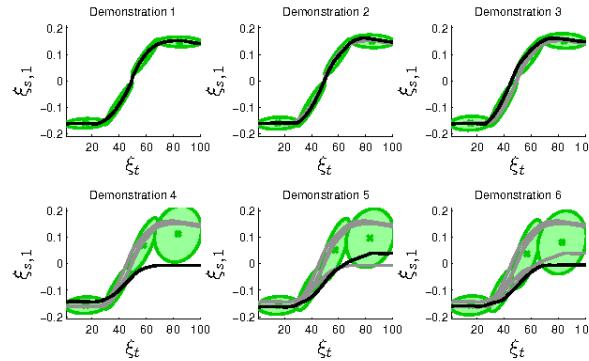


Figure 3.4: Updating of the GMM using the direct update incremental learning approach [38]. The graph shows the encoding of the data in GMM. Each algorithm uses only the latest observed trajectory (illustrated by a black line) and the previous estimation of the GMM parameters, to update the models. Demonstrations 1-3 made use of only the motion sensors, and in demonstration 4-6 the additional information was required using kinesthetic teaching

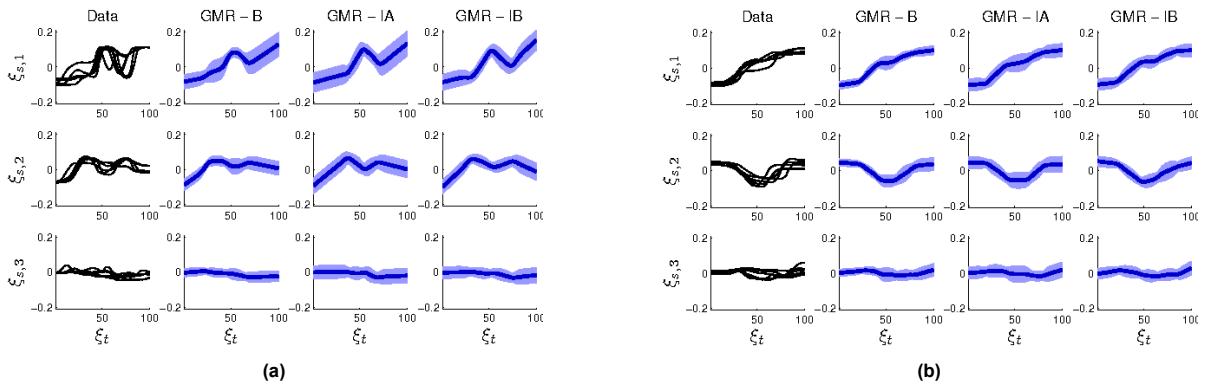


Figure 3.5: Reproduction of different gesture movements using Gaussian Mixture Regression (GMR) in a 3D latent space of motion [38]. (a & b) visualize three different gestures. For each gesture three models are trained: one with different batch B and two incremental training methods /A and /B representing the direct update and generative method.

3.2.3. Force

While position and joint angles can provide a suitable solution in some cases, often they do not account well for adaptation in the environment. The reason for this is that they are simply unaware of the changes. Take for example the grasping of a cup, if the cup is not located at the same position as it has always been and there is no additional feedback, then the robot will simply not grasp the cup. A solution for this is the usage of some force feedback. In this case, the robot could feel whether it is touching the cup and adjust its trajectory based on this. This would, however, only work in situations where the difference is quite small, as the next example will indicate.

Pastor et al. used a DMP framework to encode stereotypical movement and extended this for online movement adaptation using sensory feedback [83]. This framework was applied to a grasping task, using DMPs sensory feedback and nonlinear position and force control. The learning scheme used for this is visualized in fig. 3.6. The low-level controller was separated for the end effector (i.e. the hand) and the fingers. As visible in the graph, the used DMP is changed as a result of the demonstrated data and the previous task execution. In the experiments, the proposed control model was compared to an open-loop and a force control without DMP adaptation control mode. The task for the robot was to grasp a certain object from a table, in the case of the experiment a flash-light. For each control mode, it was tested how they would adapt if the flash-light would be displaced. The results are shown in fig. 3.7. A visible on the left, using DMP with adaptation, significantly increases the area where the robot is still able to grasp the flash-light. In the provided video, it was visible that the robot would always try to go to the position where it "thinks" the flash-light is located. However, if a force input is provided which indicates that the flash-light is not at the position where it is assumed to be, then it adapts its behavior according to this force feedback. Therefore, it could be concluded that the reason that the robot is not able to grasp the flash-light at any position in the area (as indicated with the red crosses in fig. 3.7), even when the adaptation is used, has something to do with not being in contact with the flash-light at these points (either due to going over it or not reaching it). To deal with this, another form of feedback could also be used, which does not require touching the obstacle. A good example for this is vision, based on which the position of the object could be extracted [92].

Nemec et al. [30] presented a policy learning and adaptation algorithm, using DMP, for a PiH task. Its learning scheme has been presented in fig. 3.8. The main difference, compared to other IL approaches, is that not only trajectories but also forces and torques arose during the task demonstration. To test

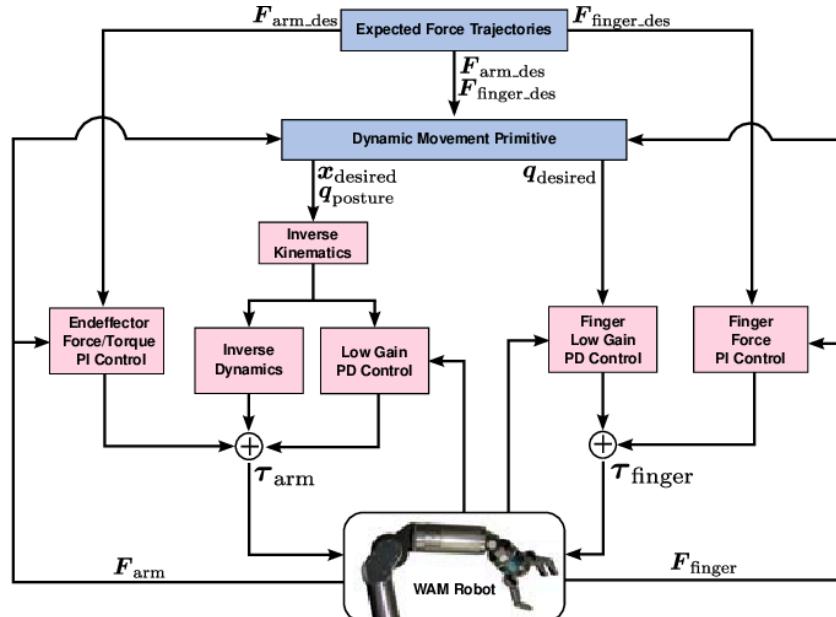


Figure 3.6: Learning scheme for a grasping task Pastor et al. The low-level controller for the end-effector (i.e. hand) and fingers are separated. The control law for the finger τ_{arm} is a combination of a position controller (consisting of a velocity-based operational space controller and an inverse dynamics law and feedback error compensation in joint space) and an end-effector force controller. The control of the fingers τ_{finger} consists of a position PD controller and a force PI controller.

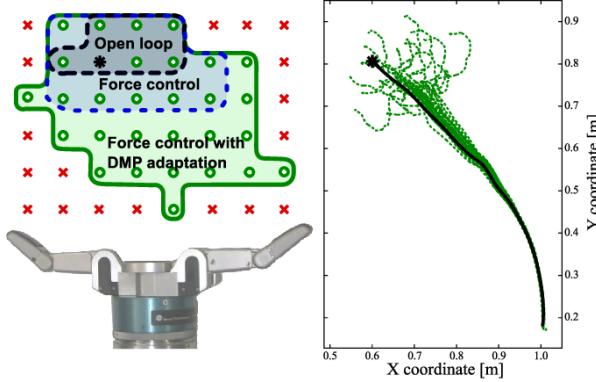


Figure 3.7: Result of grasping attempts after misplacing a flash-light [83]. (left) shows at which positions, the different control modes are still able to grasp the flash-light and (right) shows in black the open-loop trajectory and dashed-green the adapted trajectories after touching the object.

this scheme, two types of robots were used. However, as both obtained similar results, only one will be discussed. The robot, a six Degrees of Freedom (DoF) UR5 equipped with gripper and a force/torque wrist sensor, uses a high gain non-compliant controller and therefore allows for the implementation of an admittance force control law. The robot is also equipped with a magnetic tracker to track the peg's pose and orientation. The demonstration was exerted using teleoperation with a pose tracker. To test the proposed algorithm, three PiH tasks were done, of which the results shown in fig. 3.9 are of the second task where the robot has to insert a square peg when the baseplate has been rotated and translated. As expected, iteration results in a lower offset of the position. However, as visible, the offset is not completely removed. This does not necessarily have to be an issue, as multiple trajectories can be used to exert a similar task. Looking at the force, it is visible that the force used in the learned trajectories is lower than the demonstrated one. The paper names as reason an imperfect force sensor calibration. This is most likely to be accurate, as the goal is to minimize the error between the demonstrated and generated forces. Therefore, directly trying to recreate the demonstration thus leads to suboptimal performance. It was also noticed that the robot would try to replicate the learned force/torques rather than the position/orientation of the human demonstration. This, however, does not necessarily have to be a problem, as there can be multiple trajectories that result in the execution of the task.

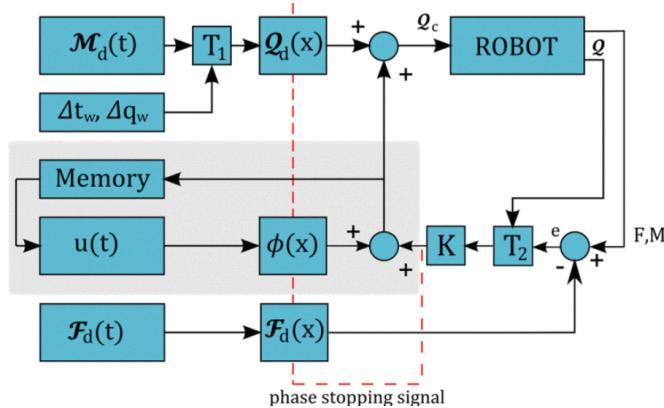


Figure 3.8: Learning scheme of PiH tasks [30], describing an iteratively updated DMP using both position/orientation and force/torque input, where $\mathcal{M}_d(t)$ denotes the original demonstrated trajectory encoded by DMPs, Q_t the output signal obtained by integrating the transformation equations of the DMPs and applying workpiece displacement (Δt_w , Δq_w) as estimated by a vision sensor (this takes place in block T_1), \mathcal{F}_d denotes the captured forces and torques captured during demonstration they are encoded using a linear combination of Radial Basis Function (RBF), the eventual goal is to minimize the difference between the measured force/torque and the desired force/torque (denoted as e), in T_2 this error is transformed into robot base coordinates. By using admittance control, the position/orientation offset is calculated and added to the existing offset of the previous iteration step $\phi(x)$. The motor commands Q_c are then a combination of the motor commands obtained from the DMPs (Q_t) and the offset. The iteration procedure is repeated until the measured forces, match the desired forces.

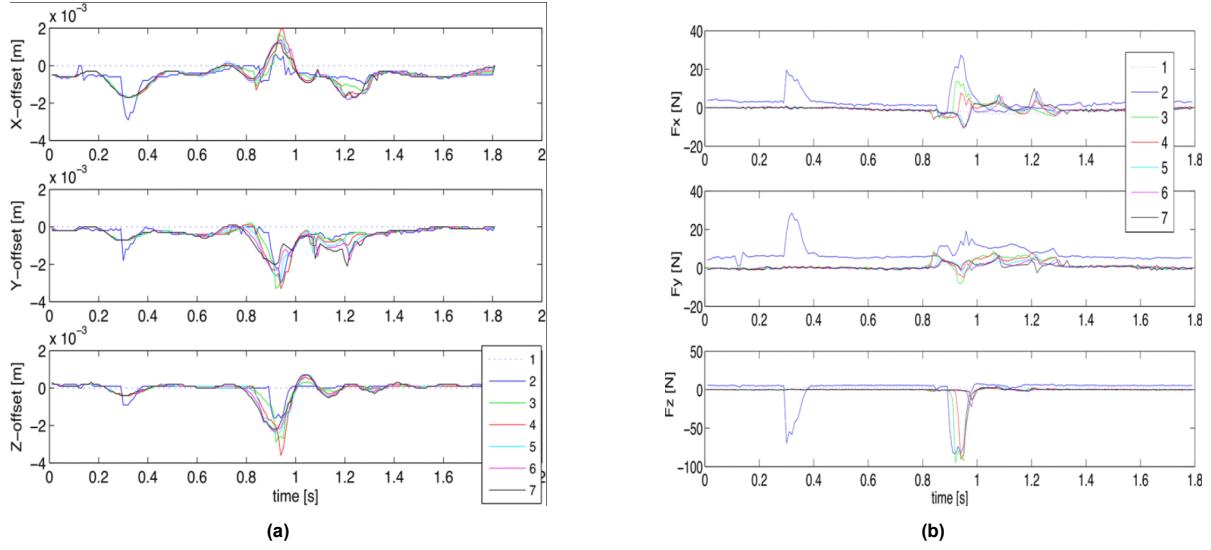


Figure 3.9: Results of PiH task with a rotated and translated baseplate [30]. (a) The learned positional offset of 6 consecutive cycles and demonstrated trajectory denoted by 1, which has no offset. (b) The force measured during the execution of the task in 6 consecutive cycles, again 1 denoting the demonstrated trajectory.

All experiments till now have shown the ability of a robot to learn from demonstrations exerted by a human expert. Peternel and Ajoudani expanded this idea by letting robots learn from each other [26]. They presented a method where one robot was assumed to be an expert. Before becoming an expert, this robot had to require its skill from a human expert. This concept is presented in fig. 3.10. The interest of the paper [26] lies in the area illustrated with red. This part consisted of three stages: a learning stage where the novice robot is completely compliant, and the expert robot has high stiffness, a stage where the novice robot learns the desired impedance behavior which encodes a follow/leader role in different phases of the task, and a third stage where both robots are assumed to be an expert. DMP was used to encode the trajectories in the first stage, and Locally Weighted Regression (LWR) to learn these encoded trajectories. In fig. 3.11 the learning process of the novice robot is shown. The interesting part occurs in the third stage. In theory, the stiffness pattern of both robots should be the same, except for a phase shift of half the period time. As visible, the stiffness pattern of the novice robot is not at all the same as the pattern of the expert robot. In the previously shown experiment [30, 38, 83, 93], the aim was always to mimic the human expert as good as possible. However, as this experiment uses collaboration between the robots to learn the skill, randomness could occur. A theory could be that the has, due to this randomness, improving the skill in comparison to the expert robot. This could therefore be an interesting topic of research. A disadvantage of randomness is the possibility that behavior becomes worse, therefore making it less reliable.

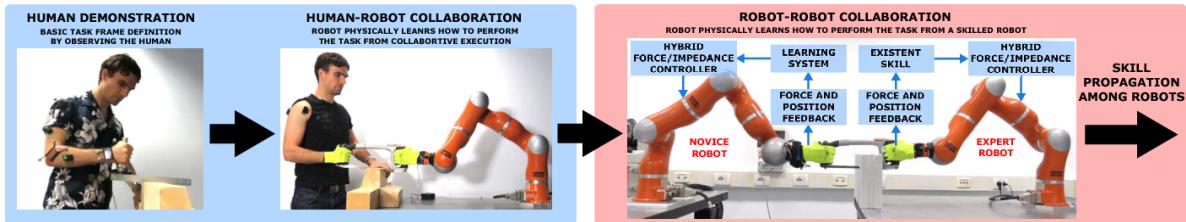


Figure 3.10: The skill learning of a novice robot [26]. The blue area illustrates the skill acquisition for an expert robot, based on the concept presented in [94]. This starts with demonstrations from the human which is observed by the robot, this allows for defining the task frame control framework. Next, the expert robot uses this knowledge to collaborate with the human and learn the task specifically, using physical interaction. The next area (indicated with red) illustrates the proposed method in the paper, where the novice robot learns from the expert robot. This part can be propagated among other robots.

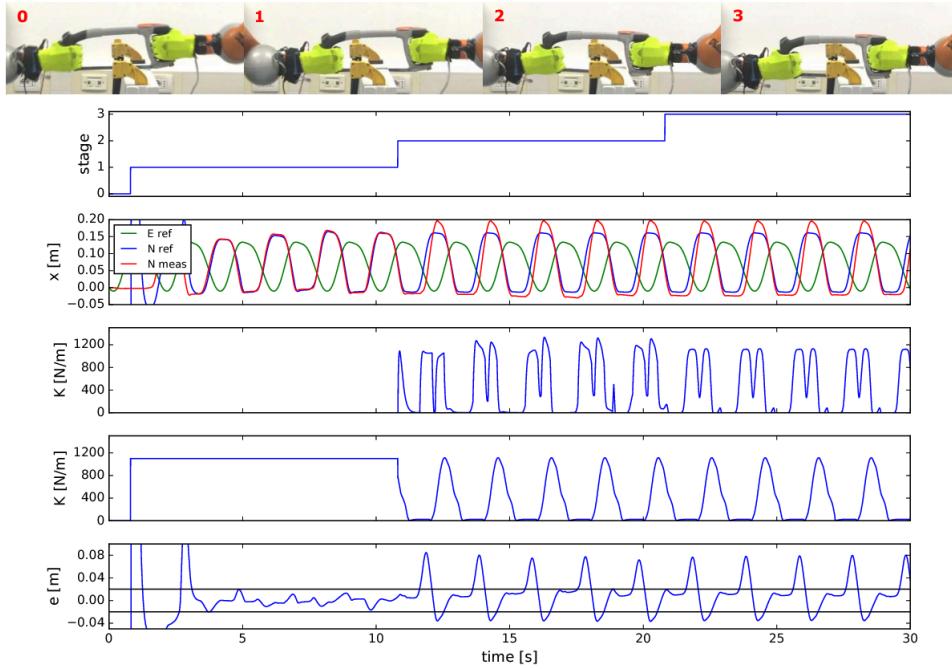


Figure 3.11: Results of robot-robot collaboration experiment on two-person sawing task [26]. The sequence of photos on top shows the sawing progress throughout the different stages. Next, from top to bottom, the graphs visualize; The different learning stages. The sawing motion in the x-axis, where the green line represents the reference motion of the expert robot, the blue line the reference motion of the novice robot, and the red line the actual motion. The stiffness of the novice robot. The stiffness of the expert robot. And lastly, the error between the reference and the actual motion used in the stiffness learning process.

3.2.4. Torques

A last frequently used state input is torque. In [90] an example is given where an exoskeleton is used to reduce the effort exerted by an expert. While this is often measured using electromyography (EMG), in essence, the goal is to reduce the torque exerted by the human. When using positional signals for this problem, the robot would as output use a learned torque trajectory. This would result in the human having to produce an opposing torque to change the joint angle. Therefore, instead, it was found to be more useful to use the torque (measured using EMG) as an input signal.

In [90] the torque was used to learn a pressure trajectory. The presented control scheme is graphically shown in fig. 3.12. The sensory feedback, resulting from the human, was measured using EMG. The system used DMP to encode the trajectory and LWR to learn it. In the paper, different experiments were conducted, each having the goal to reduce the human muscular effort as much as possible. In the tasks, the human subject had to move an object repetitively from one position to another at a certain

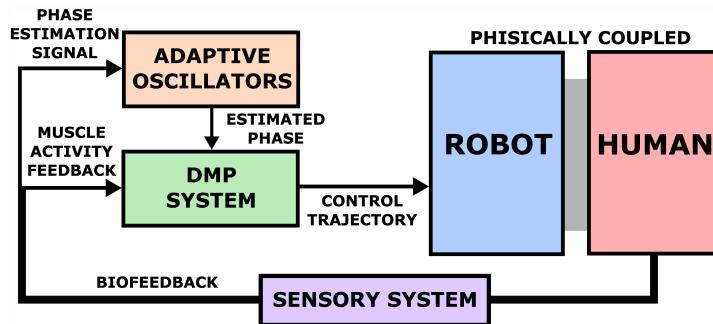


Figure 3.12: The control scheme of an exoskeleton as presented in [90]. The human is included in the control loop, by using sensory biofeedback obtained with EMG. The goal is to minimize the human muscle activity feedback by learning an appropriate assistive torque behavior in the joint. The Adaptive oscillator extracts the phase and frequency information from the human biofeedback and adjusts the phase and frequency of the learned feed-forward trajectories based on this. Based on this estimated phase and the muscle activity feedback, a DMP is created, which results in a control trajectory.

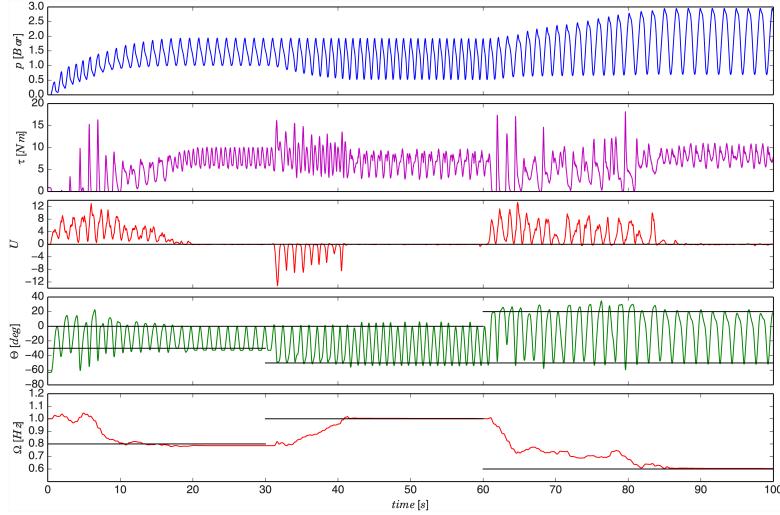


Figure 3.13: Results of the experiment on elbow exoskeleton using a trained subject [90]. From top to bottom the different plots illustrate the: command PAM pressure learned by the DMP system p , the joint torque calculated based on the force measurement of the load cell τ , the human muscle activity based on feedback obtained with the measured EMG U , the angle of the robot elbow joint Θ and the frequency of the motion as extracted by the adaptive oscillator from the muscle activity feedback Ω . The reference states are indicated by a black line.

frequency. While doing so, the exoskeleton would try to adapt to the human and, by doing so, reduce the effort. One of the experiments used an elbow exoskeleton, meaning that the subject would only be supported by torque exerted on this specific joint. The subject had to execute three different phases, where each phase consisted of a target elbow joint range and frequency. The results are visible in fig. 3.13. The third plot shows muscle activity. As visible, each phase first consists of the subject putting in much effort until it reaches the desired behavior. Once it researches this behavior, the effort is reduced to zero. To do so, assistive torque is used. The second plot shows, that it would take a bit of time to adapt to new behavior.

3.3. Reward or Cost

RL has advanced that the learned skill is not limited by the knowledge of the demonstrator. However, it has the limitation that it requires the reward or cost function to be known. The choosing of this function is critical as the resulting behavior depends on this. If this function is not known, using IRL could be an option. However, this again requires human demonstration, thus adding the drawback of being limited by its skill. In this section, the continuous learning of a robot using a reward or cost function is shown. The skill improvement will be analyzed based on the improvement in the received reward or reduced cost. In the first part, section 3.3.1, this will be done using a regular RL approach. In addition, an implementation of RL using human advice is described in section 3.3.2.

3.3.1. Regular learning

In [95], actor-critic RL (described in section 2.3.4) was implemented to execute a trajectory following task. In the paper, two control methods were introduced: a RL-based reference compensation and a RL-based input compensation. The former is formulated as a function of the joint space error and its derivative. The latter directly compensates for the reference signal fed to the system. Two of the tracking task, used to test both RL approaches, will be discussed below. For these tasks, a UR5 robot, controlled using velocity commands, is used.

The first task is to follow a discontinuous trajectory, namely a square. The goal is here to reduce the error in the z-axis. The video provided by the paper shows the learning process of the robot. The first iterations show that the robot is a bit higher above the table than required. It also shows an overshoot when trying to follow the squared shape trajectory. It is noticeable that when the square-shaped trajectory starts, there is nearly no delay in trying to follow it. During the learning process, the overshoot becomes incrementally better. The results also show that the skill improvement is the

highest at the start of the learning process, and slowly converge until the average best performance is obtained. The learning process based on the sum of rewards (return) is visualized in fig. 3.14. For the input-based method, the return of each agent (joint), slowly converges to an average best solution. In contrast, the reference-based method shows erratic behavior. The paper stated two possible reasons for this; Firstly, it can be caused by the RBFs which are initialized with an inappropriate value, causing the learning curve to deteriorate before improving. And secondly, the reference discontinuities introduce very large errors, which cause the policy value function to first behave erratically before converging to stable behavior.

The second task is to follow a continuous trajectory, namely a circle. The goal is to minimize both the x-axis and y-axis errors. The provided video shows that initially, the robot is already able to follow the circle trajectory quite well. Looking at the two different methods, fig. 3.15, it is visible that using the input-based compensation result in much faster convergence. This is a result of directly reducing the trajectory error, which means it directly compensates for the joint velocity. The response of this method is however a bit slower, as the corrected trajectory is tracked using a nominal PD controller. Lastly, as this method modifies the reference instead of the control input, the result is slightly smoother than that of the input-based controller.

In [95], multiple reasons are stated why either one of the control functions would be better for certain criteria (response time, convergence speed, smoothness, etc.). Besides these criteria, the importance of choosing a suitable control function for a specific task can be seen. Where the reference-based compensation does not even work for the following of the discontinuous trajectory, results it in better behavior for the continuous trajectory.

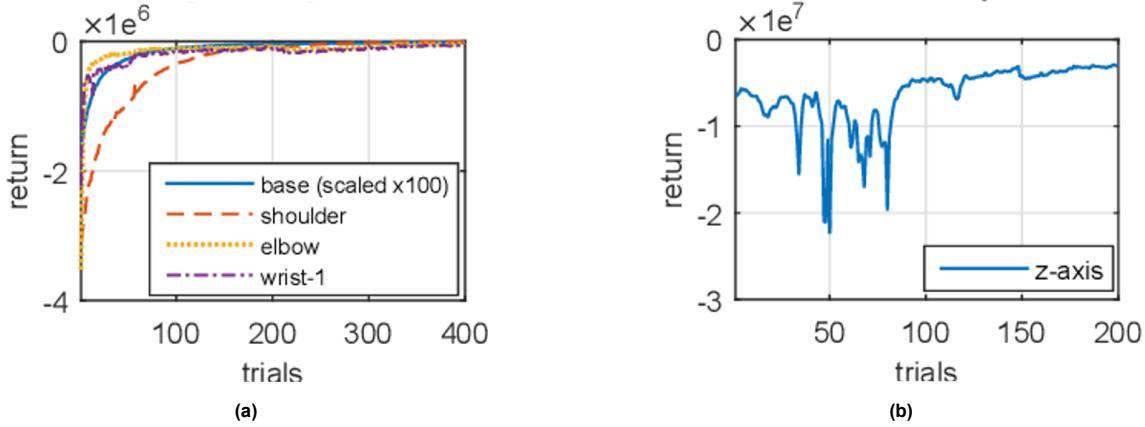


Figure 3.14: The sum of reward (return) for the squared trajectory following [95]. With (a) the RL-based input compensation, and (b) the RL-based reference compensation.

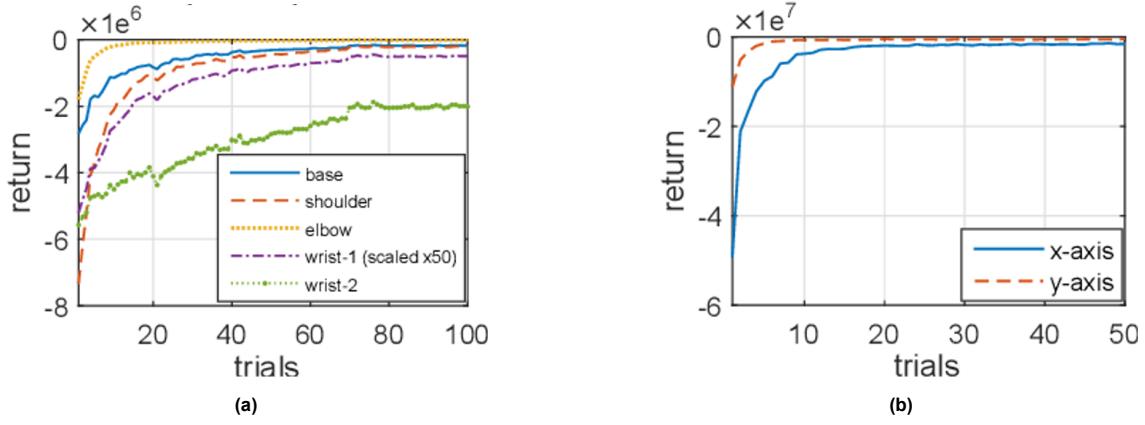


Figure 3.15: The sum of reward (return) for the circle trajectory following [95]. With (a) the RL-based input compensation, and (b) the RL-based reference compensation.

3.3.2. Human Advice

In section 2.4, the usefulness of combining IL and RL was already expanded on. This idea was adapted by Celemín et al. [80]. Here, human feedback was used to excel the policy beyond the capabilities of a human expert. In their research, two methods using feedback were used; Firstly, an adapted version of Corrective Advice Communicated by Humans (COACH) was used. This method can be used to incorporate the teacher correction in the policy search (PS) loop. During execution, the human teacher can give occasionally advice, to correct the behavior, during execution time. Secondly, they introduced *interactive PS*. This method also depends on COACH, used human guidance to influence and bias the exploration of the PS algorithm. This was done in the form of exploration noise.

To examine the performance of the approaches, a *ball-in-cup* task was used, as shown in fig. 3.16. Here the reward function was separated into two parts: one where the height of the ball is lower than the cup and one where it is higher. This was an extension of the earlier used approach [96], as that reward function would not provide any useful information if the ball was lower than the cup. Due to the extension, the task can be considered as two sub-tasks: *swinging the ball* and *catching the ball*. The policy was represented using Probabilistic Movement Primitive (ProMP) and PI^2 was the base for the PS algorithm.

The first demonstration used a policy derived using kinesthetic teaching as an initial starting point. The learning curve of this task has been presented in fig. 3.17. It is visible that both methods using human feedback converge initially much faster than the standard PS method. COACH and the interactive PS are able to achieve policies that can catch the ball at 10 and 15 episodes, respectively. Whereas the original method has this at 60. Noticeable is that the COACH implementation does not improve much after the initial catch. This can be explained as the human will stop giving advice once the ball has been caught, either because the need for improvement is not evident or it is not necessary. This shows the inability of a human to improve once they are near an optimal policy. Controversy, the interactive PS learns a bit slower at the start, which is probably due to some earliest roll-outs in the update process. However, it keeps improving until it reaches the best average performance around 70 episodes. The pure PS, might converge slower, but reaches, just as the interactive PS method, a better performance than the solely COACH method. These results show that the advantages of combining IL with RL is that it also combines its features: faster convergence (due to IL) and researching of a better performance (due to RL).

In the second set of experiments, the relevance of choosing a good initial policy is shown. In contrast, to the earlier experiment, does this one start from scratch. The results are visualized in fig. 3.18. The effect of starting from scratch is visible as the starting point is now at 1 instead of 0.1. This has especially a big effect on the PS method, which takes a lot longer to converge now. Showing its sensitivity to the initial policy. At the start, the learning process is quite slow as the random movement seems to reduce the effect of the previous ones. The results do not even show the PS receiving a behavior where it is able to catch the ball. However, if more episodes were conducted, it is reasonable to assume that it would eventually be able to. Looking at the interactive PS and the COACH method. A similar

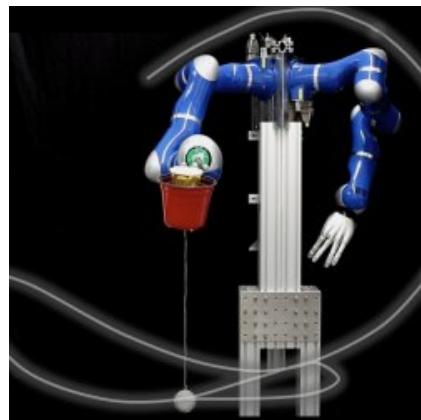


Figure 3.16: Snapshot of ball-in-cup experiment [80].

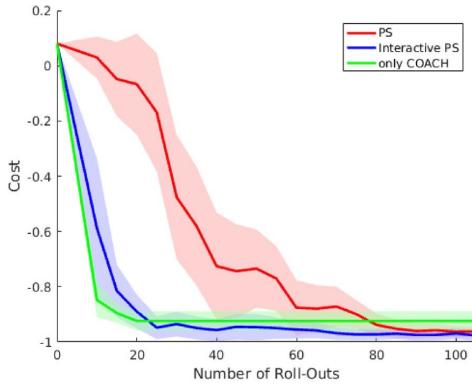


Figure 3.17: Convergence curve of *ball-in-cup* when learning from an initial demonstration [80]. This task was learned, using a standard PS method, an interactive PS method, and PS with COACH.

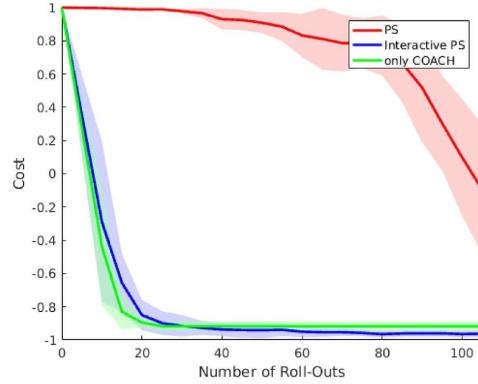


Figure 3.18: Convergence curve of *ball-in-cup* when learning from scratch [80]. This task was learned, using a standard PS method, an interactive PS method, and PS with COACH.

convergence behavior is visible as in fig. 3.17. While it takes them both a bit longer to catch the ball, they eventually achieve similar performance as before.

This research shows two important factors. Firstly, the effect of combining IL and RL on the learning time of the policy. And secondly, the usefulness of choosing a good initial policy.

3.4. Summary

The goal of this chapter was to illustrate how skills can change during the learning process. The main findings of this chapter are:

- **Learning outcome:** There are multiple types of learning outcomes, the two elaborated on are: a policy, and a reward or cost function. For online learning, the research for the latter has proven to be limited. Policy learning can be classified based on the type of state input used. The four examined in this chapter are position, joint angle, force, and torque.
- **Position and joint angles as state-input:** The benefits of position and joint angle often are quite similar. This is logical, as the position is often a translation of joint angle. While both these inputs have proven to result in good performances, they also seem to be limited when relating to adapting to environments. Changes in the environment cannot be easily adapted to by just using the configuration of the robot. Instead, some types of additional feedback need to be provided in these cases, examples are vision or some sensory information.
- **Force and torque as state input:** An advantage of using either one of these state inputs, is their ability to adapt to small changes in the environment. After slightly shifting an object from its original position, if the robot still touches it, it would be able to successfully grasp the object. However, just as for position and joint angles, additional information is required to find poses totally different from the original.
- **Collaborative task:** These types of tasks are promising for future research. Task using some kind of demonstration (teleoperation, kinaesthetic teaching, etc.) try to find a behavior that mimics the behavior of the expert as good as possible. Whereas this is also the original goal of the collaborative task, research shows mutated behavior which can - but does not necessarily have to be - better than those of the expert.
- **Reward/cost function:** Using reward/cost function to compute a policy, have the ability to come to a (good) solution after multiple iterations. However, choosing the learning system is critical. Improper definitions can result in slow convergence or even no convergence at all.
- **Human feedback:** The results in [80], show the benefit of using human feedback in addition to RL. The main advantage, is the improvement in convergence time, while as a result not losing the final performance.

4

Discussion

This literature report has aimed to give an overview of methods that could be used for continuous skill learning. In this chapter, the main finding will be discussed. In addition, this chapter aims to combine the finding of the different methods identified in chapter 2 and the effect of skill learning in chapter 3. In this discussion, two main topics will be discussed: the limitations of the robotic learning method, and robot generalizability capabilities.

4.1. Limitation of robotic learning methods

In chapter 2, multiple methods have been identified which can enable continuous skill learning. The two most broad being Imitation Learning (IL) and Reinforcement Learning (RL). Both come with their advantage and disadvantage. Starting with IL, this approach allows for a reliable and safe design, which often results in good generalization [8, 9, 63]. However, a major drawback is a fact that they are limited by the skill of the expert. In the literature, three solutions that show potential to tackle this problem have been identified.

- **Collaborative task:** Often in IL, the robot is trying to mimic the behavior of the expert as good as possible. This concept was expanded on in [26], where instead of imitating the expert, a collaborative task resulted in the robot gaining a skill. While it might not be preferable to let the robot completely learn from scratch, using a collaborative approach, could be complementary to a basic skill obtained by imitation. For example, the robot would require a basic skill by observing an expert. After which, it collaboratively performs a task. Whether the idea would work mainly depends on the type of task which should be learned. For example, the collaborative sawing task in [26] can be learned to use such an approach, whereas it might become more difficult for the peg-in-hole (PiH) task in [89]. In addition, it might not be desired to execute a task collaboratively. Therefore, an idea would be to the first approach is as a collaborative task, after which this again should be translated to a solely performed task.
- **Noise in data:** Another option would be to allow for a lot of noise in the data [30]. In essence, this results in more randomness in the solution and thus has the potential of becoming better. However, a major drawback of this idea is that it also has the potential of becoming much worse. This method can also not provide any stability guarantees, which potentially results in an unsafe system.
- **Combine with RL:** A last promising solution, is to make RL complementary to IL. Here, IL would be used to learn the demonstrated skill fast. After which the skill can be improved using RL. An example of this was given in [80]. This research showed the improvement in computational time while using a RL approach after first determining an initial policy using IL.

The idea behind each of these methods is that a larger part of the solution space will be examined to obtain good behavior. The searching of a large part of the solution space is typically done by RL. Two major drawbacks of this approach are, however, their computational time (as many iterations are needed to compute a solution), and the limitation of the skill by the choice of reward function [7]. This

first drawback was already discussed, as combining with IL would increasingly reduce this time. The choice of reward function is however an ongoing topic. While most research mentions the importance of choosing a correct reward function, the effect of choosing it incorrectly seems to be limited. Instead, most implementation has either based their function on previous research [80] or they neglected to mention why this would be the best function in comparison to any other possibility [95]. Even though this function often does not seem to be having the best reasoning behind it, it often still provides a solution for the given problems. This results in the question of what a good solution is. Some literature has used the completion of a task as a reward function. However, this binary approach often does not seem to be the best option, as not finishing a task, even if you are close, is always seemed like the same failure. Therefore, reward shaping is necessary [79]. This is a method for providing more frequent feedback on the appropriate behaviors. A method for finding rewards is Inverse Reinforcement Learning (IRL) showing to close link between RL and IL methods. This again results in the limitation of IL to only be as good as the expert.

4.2. Generalizability

One criterion discussed in section 2.5 is generalizability. The ability of a method to generalize is an ongoing topic in robotic learning [65]. As already discussed, some examples of generalizing are states and action unobserved in the demonstration, starting the task at a different initial state, and reusing skills in different problem settings [5]. Some methods were seen to potentially better generalize than others, e.g. IL methods often generalize better than RL methods. In the papers discussed in chapter 3, some generalizability capabilities also came forward.

Starting, most papers only assume the ability to start at a different position or to go to a different position, when talking about generalizability. Some examples given are [30, 83, 90, 93]. A distinction between large and small changes could be made. In [83], the system was able to grasp an object standing at a different position. While some could say that this method thus had the capabilities to generalize, I beg to differ. Instead, it could be stated that this system had the ability to adapt. Initially, it would move towards the position at which the object used to be placed. If during execution, force feedback was felt, then the robot would adapt its movement such that it would still be able to grasp the object. However, if the object was not touched, it would also not be able to grasp it. While this ability to adapt to small changes can definitely be beneficial, e.g. when the position of an object can not precisely be determined, it does not mean it is able to generalize.

Another example is given in [30], where a PiH task was performed. After using demonstrations to train the model, sensory feedback was used to detect the position where the peg had to be inserted. Using this sensory information, a position could be extracted.

So, while some methods have been stated in literature to generalize better than other methods. It was found, that in order to generalize well, it is necessary to have information about the environment. It should be mentioned, that while the usage of sensors to extract information about the world (based on which generalization and also adaptation could be done) is frequently implemented, it has not thoroughly been examined in this research. Therefore, additional research should be acquired to make better conclusions on this subject. Additionally, in this section, only the ability of a system to move from or to a new position was discussed. Eventually, it could also be desired for the system to execute its skill in a totally different environment or to adapt its skill for a different task.

5

Conclusion & Future research

This chapter aims to give a conclusion, based on the research question defined in chapter 1. This is followed by some advice for further work.

5.1. Conclusion

The research question on which will be concluded in this chapter is as follows:

How do methods, found in literature, enable continuous skill learning of robots, and how does the learned skill change during the learning process?

In the literature, two broadly applied approaches have been identified (chapter 3). Both allow for continuous skill learning of robots:

- **Imitation Learning (IL):** learning a policy by mimicking some demonstrated task.
- **Reinforcement Learning (RL):** learning a policy using a trial-and-error approach, where the behavior is either rewarded using a reward function or penalized using a cost-function.

In literature, it was found that both approaches have another skill learning process. Starting with Imitation Learning (IL), different imitation methods allow for different learning processes. The main distinction can be made between methods using one demonstration and methods using multiple. For the one demonstration, the learner can either learn a policy in one go, which cannot be seen as an online implementation, or the one demonstration can be used to compare the behavior of the robot to (as done in [30]). In this last approach, the policy is iteratively updated until the difference between the demonstrated data and the behavior of the robot is minimized. This thus assumes that the optimal behavior is not learned in one go. For multiple demonstrations, the learner can either obtain all demonstrations at once and compute a distribution based on this, which is again not an online implementation, or each demonstration is used to iteratively update the policy. The latter then takes each new iteration not all previous ones into consideration, but instead only uses the last generated policy [26, 38]

RL in most cases enables online learning, as each iteration the policy is updated. The skill learning shows that with each iteration the behavior becomes slightly better. Often the system starts with updating into one direction in the solution space (the one direction meaning towards an optimal/good behavior), until it overshoots. After the overshoot, the direction is shifted, resulting in a zigzag behavior around the optimal solution. An example is the ball-in-cup experiment [80], where it first undershoots a couple of times, next it overshoots, followed by an undershot, and so on until it eventually goes through it.

Concluding, IL method provide a reliable result that is evenly good as that of the expert. Whereas, RL is less reliable but has the potential of finding better solutions than IL.

5.2. Future research

During this literature study, multiple gaps in research have come to light. Therefore, in this section, some future research topics will be presented. These topics have been identified based on the papers' research in this report and are therefore assumed not to be existent at this moment.

Staring, research on both IL and RL has been widely applied. There is, however, a lack of research comparing both methods. To do so, suitable metrics should be used. In most cases, these metrics will be quite similar to the variables on which the reward or cost function of the RL method depends upon.

While reward functions are often mentioned as being highly critical for the result of the RL algorithm, comparing different functions for robotic learning has not been found. Therefore, it could be useful to compare different reward functions for the same task.

Literature has shown the potential of exerting an experts' skill by the usage of a collaborative task for skill learning. As the demonstrated data is not directly transferred from the expert to the novice robot, there is a potential of becoming better. This idea is of course still limited by the conditions of the task. However, applying this principle multiple times could result in a divergence of the original execution of a certain task. This again can also have positive results on the generalizability of the system.

Combining the features of IL and RL methods into one approach could result in the robot exerting the behavior of using either of these approaches singular. While some research has already shown the potential of this idea, there is a lack of comparison between the methods. Therefore, this could be an interesting research topic.

References

- [1] Sylvain Calinon et al. “A probabilistic approach based on dynamical systems to learn and reproduce gestures by imitation”. In: *IEEE Robotics and Automation Magazine* 17.2 (2010), pp. 44–54.
- [2] Aude G. Billard, Sylvain Calinon, and Rüdiger Dillmann. “Learning from humans”. In: *Springer Handbook of Robotics* (Jan. 2016), pp. 1995–2014. DOI: 10.1007/978-3-319-32552-1_74.
- [3] Luka Peterlin, Erhan Oztop, and Jan Babic. “A Shared Control Method for Online Human-in-the-Loop Robot Learning Based on Locally Weighted Regression”. In: IEEE, 2016, pp. 3900–3906. ISBN: 9781509037612.
- [4] Duy Nguyen-Tuong and Jan Peters. “Model learning for robot control: a survey”. In: *Cognitive processing* 12.4 (2011), pp. 319–340.
- [5] Oliver Kroemer, Scott Niekum, and George Konidaris. “A Review of Robot Learning for Manipulation: Challenges, Representations, and Algorithms”. In: *Journal of Machine Learning Research* 22 (2021), pp. 1–82. URL: <http://jmlr.org/papers/v22/19-804.html>.
- [6] Shixiang Gu et al. “Continuous deep q-learning with model-based acceleration”. In: *International conference on machine learning*. PMLR. 2016, pp. 2829–2838.
- [7] Jens Kober, J. Andrew Bagnell, and Jan Peters. “Reinforcement learning in robotics: A survey”. In: *International Journal of Robotics Research* 32 (11 Sept. 2013), pp. 1238–1274. ISSN: 02783649. DOI: 10.1177/0278364913495721.
- [8] Takayuki Osa et al. “An Algorithmic Perspective on Imitation Learning”. In: *Foundations and Trends ® in Robotics* 7 (2 2018), pp. 1–179. DOI: 10.1561/2300000053.
- [9] Petar Kormushev, Sylvain Calinon, and Darwin G Caldwell. “Reinforcement learning in robotics: Applications and real-world challenges”. In: *Robotics* 2.3 (2013), pp. 122–148.
- [10] Kyriacos Shiarlis, Joao Messias, and Shimon Whiteson. “Inverse Reinforcement Learning from Failure”. In: 2016, pp. 1060–1068. URL: www.ifaamas.org.
- [11] Andrew Y Ng, Stuart J Russell, et al. “Algorithms for inverse reinforcement learning.” In: *Icmi*. Vol. 1. 2000, p. 2.
- [12] Zuyuan Zhu and Huosheng Hu. “Robot Learning from Demonstration in Robotic Assembly: A Survey”. In: *Robotics* 7 (2 Apr. 2018), p. 17. ISSN: 22186581. DOI: 10.3390/ROBOTICS7020017. URL: <https://www.mdpi.com/2218-6581/7/2/17>.
- [13] Harish C. Ravichandar et al. “Recent Advances in Robot Learning from Demonstration”. In: *Annual Review of Control, Robotics, and Autonomous Systems* 3 (1 May 2020), pp. 297–330. ISSN: 2573-5144. DOI: 10.1146/ANNREV-CONTROL-100819-063206.
- [14] Brenna D Argall et al. “A survey of robot learning from demonstration”. In: *Robotics and autonomous systems* 57.5 (2009), pp. 469–483.
- [15] Jonathan Ho and Stefano Ermon. “Generative adversarial imitation learning”. In: *Advances in neural information processing systems* 29 (2016), pp. 4565–4573.
- [16] Stéphane Ross et al. “A Reduction of Imitation Learning and Structured Prediction to No-Regret Online Learning”. In: *JMLR Workshop and Conference Proceedings*, 2011, pp. 627–635.
- [17] Boyuan Zheng et al. “Imitation Learning: Progress, Taxonomies and Opportunities”. In: *arXiv preprint arXiv:2106.12177* (2021).
- [18] Yiren Lu and Jonathan Tompson. “ADAIL: Adaptive Adversarial Imitation Learning”. In: *CoRR* (2020).

- [19] Aude Billard et al. "Survey: Robot Programming by Demonstration". In: Springer, 2008, pp. 1371–1394. ISBN: 978-3-540-23957-4. DOI: 10.1007/978-3-540-30301-5_60. URL: <https://infoscience.epfl.ch/record/114050>.
- [20] Maria Kyarini et al. "Robot learning of industrial assembly task via human demonstrations". In: *Autonomous Robots* 43.1 (2019), pp. 239–257.
- [21] Auke Jan Ijspeert, Jun Nakanishi, and Stefan Schaal. "Movement imitation with nonlinear dynamical systems in humanoid robots". In: *Proceedings - IEEE International Conference on Robotics and Automation* 2 (2002), pp. 1398–1403. ISSN: 10504729. DOI: 10.1109/ROBOT.2002.1014739.
- [22] Auke Jan Ijspeert et al. "Dynamical movement primitives: learning attractor models for motor behaviors". In: *Neural computation* 25.2 (2013), pp. 328–373.
- [23] Bin Fang et al. "Skill learning for human-robot interaction using wearable device". In: *Tsinghua Science and Technology* 24.6 (2019), pp. 654–662.
- [24] Matteo Saveriano et al. "Dynamic Movement Primitives in Robotics: A Tutorial Survey". In: *arXiv preprint arXiv:2102.03861* (2021). DOI: 10.1177/ToBeAssigned. URL: www.sagepub.com/.
- [25] Karl Glatz. *Adaptive Learning from Demonstration using Dynamic Movement Primitives*. 2012.
- [26] Luka Peternel and Arash Ajoudani. "Robots learning from robots: A proof of concept study for co-manipulation tasks". In: *2017 IEEE-RAS 17th International Conference on Humanoid Robotics (Humanoids)*. IEEE. 2017, pp. 484–490.
- [27] Fares J Abu-Dakka et al. "Adaptation of manipulation skills in physical contact with the environment to reference force profiles". In: *Autonomous Robots* 39.2 (2015), pp. 199–217.
- [28] Aleš Ude et al. "Orientation in cartesian space dynamic movement primitives". In: *2014 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE. 2014, pp. 2997–3004.
- [29] Fares J Abu-Dakka and Ville Kyrki. "Geometry-aware dynamic movement primitives". In: *2020 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE. 2020, pp. 4421–4426.
- [30] Bojan Nemec et al. "Transfer of assembly operations to new workpiece poses by adaptation to the desired force profile". In: *2013 16th International Conference on Advanced Robotics (ICAR)*. IEEE. 2013, pp. 1–7.
- [31] Alexandros Paraschos et al. "Using probabilistic movement primitives in robotics". In: *Autonomous Robots* 42 (3 Mar. 2018), pp. 529–551. ISSN: 15737527. DOI: 10.1007/s10514-017-9648-7.
- [32] Alexandros Paraschos et al. "Probabilistic Movement Primitives". In: *Advances in neural information processing systems* (2013).
- [33] Daniel Schäle, Martin F Stoelen, and Erik Kyrkjebo. "Incremental Learning of Probabilistic Movement Primitives (ProMPs) for Human-Robot Cooperation". In: (2021).
- [34] Guilherme Maeda et al. "Learning interaction for collaborative tasks with probabilistic movement primitives". In: *2014 IEEE-RAS International Conference on Humanoid Robots*. IEEE. 2014, pp. 527–534.
- [35] Dorothea Koert et al. "Incremental learning of an open-ended collaborative skill library". In: *International Journal of Humanoid Robotics* 17.01 (2020), p. 2050001.
- [36] Congcong Ye, Jixiang Yang, and Han Ding. "Bagging for Gaussian mixture regression in robot learning from demonstration". In: *Journal of Intelligent Manufacturing* (2020), pp. 1–13.
- [37] Sylvain Calinon. "A tutorial on task-parameterized movement learning and retrieval". In: *Intelligent Service Robotics* 9 (2016), pp. 1–29. DOI: 10.1007/s11370-015-0187-9. URL: [http://www.idiap.ch/software/pbdlib/..](http://www.idiap.ch/software/pbdlib/)
- [38] Sylvain Calinon and Aude Billard. "Incremental Learning of Gestures by Imitation in a Humanoid Robot". In: 2007, pp. 255–262. ISBN: 9781595936172.
- [39] Jose Hoyos et al. "Incremental learning of skills in a task-parameterized gaussian mixture model". In: *Journal of Intelligent & Robotic Systems* 82.1 (2016), pp. 81–99.
- [40] Weiyong Si, Ning Wang, and Chenguang Yang. "A review on manipulation skill acquisition through teleoperation-based learning from demonstration". In: *Cognitive Computation and Systems* 3 (1 Mar. 2021), pp. 1–16. ISSN: 2517-7567. DOI: 10.1049/ccs2.12005.

- [41] Thomas Cederborg et al. "Incremental local online Gaussian Mixture Regression for imitation learning of multiple tasks". In: *IEEE/RSJ 2010 International Conference on Intelligent Robots and Systems, IROS 2010 - Conference Proceedings* (2010), pp. 267–274. DOI: 10.1109/IROS.2010.5652040.
- [42] Miguel Arduengo et al. "Gaussian-process-based robot learning from demonstration". In: *arXiv preprint arXiv:2002.09979* (2020).
- [43] Noémie Jaquier, David Ginsbourger, and Sylvain Calinon. "Learning from demonstration with model-based Gaussian process". In: *Conference on Robot Learning*. PMLR. 2020, pp. 247–257.
- [44] Chirag Goyal. *Deep understanding of discriminative and generative models*. <https://www.analyticsvidhya.com/blog/2021/07/deep-understanding-of-discriminative-and-generative-models-in-machine-learning/>. [Accessed: 14-1-2021]. July 2021.
- [45] Duy Nguyen-Tuong and Jan Peters. "Local gaussian process regression for real-time model-based robot control". In: *2008 IEEE/RSJ International Conference on Intelligent Robots and Systems*. IEEE. 2008, pp. 380–385.
- [46] Daniel H Grollman and Odest Chadwicke Jenkins. "Sparse incremental learning for interactive robot control policy estimation". In: *2008 IEEE International Conference on Robotics and Automation*. IEEE. 2008, pp. 3315–3320.
- [47] Jie Yang, Yangsheng Xu, and Chiou S. Chen. "Hidden Markov Model Approach to Skill Learning and Its Application to Telerobotics". In: *IEEE Transactions on Robotics and Automation* 10 (5 1994), pp. 621–631. ISSN: 1042296X. DOI: 10.1109/70.326567.
- [48] Aleksandar Vakanski et al. "Trajectory learning for robot programming by demonstration using hidden Markov model and dynamic time warping". In: *IEEE Transactions on Systems, Man, and Cybernetics, Part B (Cybernetics)* 42.4 (2012), pp. 1039–1052.
- [49] Dana Kulic and Yoshihiko Nakamura. "Incremental learning of human behaviors using hierarchical hidden Markov models". In: *2010 IEEE/RSJ International Conference on Intelligent Robots and Systems*. IEEE. 2010, pp. 4649–4655.
- [50] Leonel Rozo, Pablo Jiménez, and Carme Torras. "A robot learning from demonstration framework to perform force-based manipulation tasks". In: *Intelligent service robotics* 6.1 (2013), pp. 33–51.
- [51] Sylvain Calinon et al. "Learning and reproduction of gestures by imitation". In: *IEEE Robotics & Automation Magazine* 17.2 (2010), pp. 44–54.
- [52] Sylvain Calinon et al. "Learning collaborative manipulation tasks by demonstration using a haptic interface". In: *2009 International Conference on Advanced Robotics*. IEEE. 2009, pp. 1–6.
- [53] Shun-Zheng Yu. "Hidden semi-Markov models". In: *Artificial intelligence* 174.2 (2010), pp. 215–243.
- [54] Dana Kulic, Wataru Takano, and Yoshihiko Nakamura. "Incremental on-line hierarchical clustering of whole body motion patterns". In: *RO-MAN 2007-The 16th IEEE International Symposium on Robot and Human Interactive Communication*. IEEE. 2007, pp. 1016–1021.
- [55] Dana Kulic, Wataru Takano, and Yoshihiko Nakamura. "Online segmentation and clustering from continuous observation of whole body motions". In: *IEEE Transactions on Robotics* 25.5 (2009), pp. 1158–1166.
- [56] Dizan Vasquez, Thierry Fraichard, and Christian Laugier. "Growing hidden markov models: An incremental tool for learning and predicting human and vehicle motion". In: *The International Journal of Robotics Research* 28.11-12 (2009), pp. 1486–1506.
- [57] Stuart Russell. "Learning agents for uncertain environments". In: *Proceedings of the eleventh annual conference on Computational learning theory*. 1998, pp. 101–103.
- [58] Shao Zhifei and Er Meng Joo. "A survey of inverse reinforcement learning techniques". In: *International Journal of Intelligent Computing and Cybernetics* (2012).
- [59] Pieter Abbeel and Andrew Y Ng. "Apprenticeship learning via inverse reinforcement learning". In: *Proceedings of the twenty-first international conference on Machine learning*. 2004, p. 1.

- [60] Saurabh Arora and Prashant Doshi. "A survey of inverse reinforcement learning: Challenges, methods and progress". In: *Artificial Intelligence* (2021), p. 103500.
- [61] David Silver, J Andrew Bagnell, and Anthony Stentz. "Learning from demonstration for autonomous navigation in complex unstructured terrain". In: *The International Journal of Robotics Research* 29.12 (2010), pp. 1565–1592.
- [62] Nathan D Ratliff, J Andrew Bagnell, and Martin A Zinkevich. "Maximum margin planning". In: *Proceedings of the 23rd international conference on Machine learning*. 2006, pp. 729–736.
- [63] Brian D Ziebart et al. "Maximum entropy inverse reinforcement learning." In: *Aaa*. Vol. 8. Chicago, IL, USA. 2008, pp. 1433–1438.
- [64] Zheng Wu et al. "Efficient sampling-based maximum entropy inverse reinforcement learning with application to autonomous driving". In: *IEEE Robotics and Automation Letters* 5.4 (2020), pp. 5355–5362.
- [65] Saurabh Arora, Prashant Doshi, and Bikramjit Banerjee. "Online inverse reinforcement learning under occlusion". In: *Proceedings of the 18th International Conference on Autonomous Agents and Multi-Agent Systems*. 2019.
- [66] Kai Arulkumaran, Marc Peter Deisenroth, et al. "Deep reinforcement learning: A brief survey". In: *IEEE Signal Processing Magazine* 34 (6 Nov. 2017), pp. 26–38. ISSN: 10535888. DOI: 10.1109/MSP.2017.2743240.
- [67] Athanasios S. Polydoros and Lazaros Nalpantidis. "Survey of Model-Based Reinforcement Learning: Applications on Robotics". In: *J Intell Robot Syst* 86 (Jan. 2017), pp. 153–173. DOI: 10.1007/s10846-017-0468-y.
- [68] Richard S Sutton and Andrew G Barto. "Reinforcement Learning: An Introduction". In: (1998).
- [69] Hongming Zhang et al. "Taxonomy of Reinforcement Learning Algorithms". In: *Deep Reinforcement Learning: Fundamentals, Research and Applications* (Jan. 2020), pp. 125–133. DOI: 10.1007/978-981-15-4095-0_3. URL: https://link.springer.com/chapter/10.1007/978-981-15-4095-0_3.
- [70] Philip Odonkor and Kemper Lewis. "Control of shared energy storage assets within building clusters using reinforcement learning". In: *International Design Engineering Technical Conferences and Computers and Information in Engineering Conference*. Vol. 51753. American Society of Mechanical Engineers. 2018, V02AT03A028.
- [71] Andrew Lobbezoo, Yanjun Qian, and Hyock-Ju Kwon. "Reinforcement Learning for Pick and Place Operations in Robotics: A Survey". In: *Robotics* 10.3 (2021), p. 105.
- [72] David Silver. *Lecture 4: Model-free prediction*. 2015.
- [73] Miguel Morales. "Grokking deep reinforcement learning". In: Manning Publications, 2020. Chap. 6.
- [74] Zihan Ding et al. "Introduction to Reinforcement Learning". In: *Deep Reinforcement Learning: Fundamentals, Research and Applications* (Jan. 2020), pp. 47–123. DOI: 10.1007/978-981-15-4095-0_2. URL: https://link.springer.com/chapter/10.1007/978-981-15-4095-0_2.
- [75] Marc Peter Deisenroth, Gerhard Gerhard Neumann, and Jan Peters. "A Survey on Policy Search for Robotics". In: *Foundations and Trends R in Robotics* 2 (2 2011), pp. 1–142. DOI: 10.1561/2300000021.
- [76] Ivo Grondman et al. "A survey of actor-critic reinforcement learning: Standard and natural policy gradients". In: *IEEE Transactions on Systems, Man, and Cybernetics, Part C (Applications and Reviews)* 42.6 (2012), pp. 1291–1307.
- [77] Yudha P Pane, Subramanya P Nageshrao, and Robert Babuška. "Actor-critic reinforcement learning for tracking control in robotics". In: *2016 IEEE 55th conference on decision and control (CDC)*. IEEE. 2016, pp. 5819–5826.
- [78] Kai-Wei Chang et al. "Learning to search better than your teacher". In: *International Conference on Machine Learning*. PMLR. 2015, pp. 2058–2066.
- [79] Wen Sun, J Andrew Bagnell, and Byron Boots. "Truncated horizon policy search: Combining reinforcement learning & imitation learning". In: *arXiv preprint arXiv:1805.11240* (2018).

- [80] Carlos Celemín et al. “Reinforcement learning of motor skills using policy search and human corrective advice”. In: *The International Journal of Robotics Research* 38.14 (2019), pp. 1560–1580.
- [81] S. Mohammad Khansari-Zadeh and Aude Billard. “Learning stable nonlinear dynamical systems with Gaussian mixture models”. In: *IEEE Transactions on Robotics* 27 (5 Oct. 2011), pp. 943–957. ISSN: 15523098. DOI: 10.1109/TRO.2011.2159412.
- [82] Jordi Bautista-Ballester, Jaume Vergés-Llahí, and Domènec Puig. “Programming by demonstration: A taxonomy of current relevant methods to teach and describe new skills to robots”. In: *ROBOT2013: First Iberian Robotics Conference*. Springer. 2014, pp. 287–300.
- [83] Peter Pastor et al. “Online movement adaptation based on previous sensor experiences”. In: *2011 IEEE/RSJ International Conference on Intelligent Robots and Systems*. IEEE. 2011, pp. 365–371.
- [84] Guilherme J Maeda et al. “Probabilistic movement primitives for coordination of multiple human–robot collaborative tasks”. In: *Autonomous Robots* 41.3 (2017), pp. 593–612.
- [85] Abdeslam Boularias, Jens Kober, and Jan Peters. “Relative entropy inverse reinforcement learning”. In: *Proceedings of the Fourteenth International Conference on Artificial Intelligence and Statistics*. JMLR Workshop and Conference Proceedings. 2011, pp. 182–189.
- [86] Yangang Ren et al. “Improving generalization of reinforcement learning with minimax distributional soft actor-critic”. In: *2020 IEEE 23rd International Conference on Intelligent Transportation Systems (ITSC)*. IEEE. 2020, pp. 1–6.
- [87] George Konidaris et al. “Robot learning from demonstration by constructing skill trees”. In: *The International Journal of Robotics Research* 31.3 (2012), pp. 360–375.
- [88] Mariantonieta Gutierrez Soto and Hojjat Adeli. “Recent advances in control algorithms for smart structures and machines”. In: *Expert Systems* 34.2 (2017), e12205.
- [89] Ioannis Havoutis and Sylvain Calinon. “Supervisory teleoperation with online learning and optimal control”. In: *2017 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE. 2017, pp. 1534–1540.
- [90] Luka Peternel et al. “Adaptive control of exoskeleton robots for periodic assistive behaviours based on EMG feedback minimisation”. In: *PloS one* 11.2 (2016), e0148942.
- [91] Leonel Rozo, Pablo Jiménez, and Carme Torras. “Robot learning from demonstration of force-based tasks with multiple solution trajectories”. In: *2011 15th International Conference on Advanced Robotics (ICAR)*. IEEE. 2011, pp. 124–129.
- [92] David A Duque, Flavio A Prieto, and Jose G Hoyos. “Trajectory generation for robotic assembly operations using learning by demonstration”. In: *Robotics and Computer-Integrated Manufacturing* 57 (2019), pp. 292–302.
- [93] Ioannis Havoutis, Ajay Kumar Tanwani, and Sylvain Calinon. “Online incremental learning of manipulation tasks for semi-autonomous teleoperation”. In: Institute of Electrical and Electronics Engineers. 2016.
- [94] Luka Peternel, Leonel Rozo, et al. “A Method for Derivation of Robot Task-Frame Control Authority from Repeated Sensory Observations”. In: *IEEE Robotics and Automation Letters* 2 (2 2017), pp. 719–726. URL: <https://ieeexplore.ieee.org/stamp/stamp.jsp?tp=&arnumber=7812599>.
- [95] Yudha P Pane et al. “Reinforcement learning based compensation methods for robot manipulators”. In: *Engineering Applications of Artificial Intelligence* 78 (2019), pp. 236–247.
- [96] Jens Kober and Jan Peters. “Policy search for motor primitives in robotics”. In: *Learning Motor Skills*. Springer. 2014, pp. 83–117.