

A Survey on Policy Search for Robotics

Marc Peter Deisenroth^{*1}, Gerhard
Neumann^{*2} and Jan Peters³

¹ *Technische Universität Darmstadt, Germany*
Imperial College London, UK

² *Technische Universität Darmstadt, Germany*

³ *Technische Universität Darmstadt, Germany*
Max Planck Institute for Intelligent Systems, Germany

Abstract

Policy search is a subfield in reinforcement learning which focuses on finding good parameters for a given policy parametrization. It is well suited for robotics as it can cope with high-dimensional state and action spaces, one of the main challenges in robot learning. We review recent successes of both model-free and model-based policy search in robot learning.

Model-free policy search is a general approach to learn policies based on sampled trajectories. We classify model-free methods based on their policy evaluation strategy, policy update strategy and exploration strategy and present an unified view on existing algorithms. Learning a policy is often easier than learning an accurate forward model, and, hence, model-free methods are more frequently used in practice. However, for each sampled trajectory, it is necessary to interact with the robot, which can be time consuming and challenging in practice.

* Both authors contributed equally.

Model-based policy search addresses this problem by first learning a simulator of the robot's dynamics from data. Subsequently, the simulator generates trajectories that are used for policy learning. For both model-free and model-based policy search methods, we review their respective properties and their applicability to robotic systems.

Contents

1	Introduction	1
1.1	Robot Control as a Reinforcement Learning Problem	2
1.2	Policy Search Taxonomy	5
1.2.1	Model-free and Model-based Policy Search	7
1.3	Typical Policy Representations	8
1.4	Outline	11
2	Model-free Policy Search	12
2.1	Exploration Strategies	14
2.1.1	Exploration in Action Space versus Exploration in Parameter Space	14
2.1.2	Episode-based versus Step-based Exploration	16
2.1.3	Uncorrelated versus Correlated Exploration	17
2.1.4	Updating the Exploration Distribution	18
2.2	Policy Evaluation Strategies	18
2.2.1	Step-Based Policy Evaluation	19
2.2.2	Episode-Based Policy Evaluation	20

ii *Contents*

2.2.3	Comparison of Step- and Episode-based Evaluation	22
2.3	Important Extensions	22
2.3.1	Generalization to Multiple Tasks	22
2.3.2	Learning Multiple Solutions for a Single Motor Task	23
2.4	Policy Update Strategies	24
2.4.1	Policy Gradient Methods	25
2.4.1.1	Finite Difference Methods	26
2.4.1.2	Likelihood-Ratio Policy Gradients	26
2.4.1.3	Natural Gradients	32
2.4.2	Expectation Maximization Policy Search Approaches	41
2.4.2.1	The Standard Expectation Maximization Algorithm	41
2.4.2.2	Policy Search as an Inference Problem	43
2.4.2.3	Monte-Carlo EM-based Policy Search.	45
2.4.2.4	Variational Inference-based Methods	52
2.4.3	Information-Theoretic Approaches	54
2.4.3.1	Episode-based Relative Entropy Policy Search	55
2.4.3.2	Episode-Based Extension to Multiple Contexts.	57
2.4.3.3	Learning Multiple Solutions with REPS	59
2.4.3.4	Step-based REPS for Infinite Horizon Problems	62
2.4.4	Miscellaneous Important Methods	67
2.4.4.1	Stochastic Optimization	67
2.4.4.2	Policy Improvement by Path Integrals	69
2.5	Real Robot Applications with Model-Free Policy Search	78
2.5.1	Learning Baseball with eNAC	78
2.5.2	Learning Ball-in-the-Cup with PoWER	79
2.5.3	Learning Pan-Cake Flipping with PoWER/RWR	80
2.5.4	Learning Dart Throwing with CRKR	81
2.5.5	Learning Table Tennis with CRKR	82
2.5.6	Learning Tetherball with HiREPS	83

3	Model-based Policy Search	85
3.1	Probabilistic Forward Models	91
3.1.1	Locally Weighted Bayesian Regression	91
3.1.2	Gaussian Process Regression	93
3.2	Long-Term Predictions with a Given Model	95
3.2.1	Sampling-based Trajectory Prediction: PEGASUS	95
3.2.1.1	Trajectory Sampling and Policy Evaluation	95
3.2.1.2	Practical Considerations	96
3.2.2	Deterministic Long-Term Predictions	97
3.2.2.1	Practical Considerations	100
3.3	Policy Updates	101
3.3.1	Model-based Policy Updates without Gradient Information	101
3.3.2	Model-based Policy Updates with Gradient Information	101
3.3.2.1	Sampling-based Policy Gradients	102
3.3.2.2	Analytic Policy Gradients	102
3.3.3	Discussion	104
3.4	Model-based Policy Search Algorithms with Robot Applications	105
3.4.1	Sampling-based Trajectory Prediction	105
3.4.1.1	Locally Weighted Regression Forward Models and Sampling-based Trajectory Prediction	106
3.4.1.2	Gaussian Process Forward Models and Sampling-based Trajectory Prediction	109
3.4.2	Deterministic Trajectory Predictions	109
3.4.2.1	Gaussian Process Forward Models and Deterministic Trajectory Prediction	109
3.4.3	Overview of Model-based Policy Search Algorithms	113
3.5	Important Properties of Model-based Methods	113
3.5.1	Deterministic and Stochastic Long-Term Predictions	114
3.5.2	Treatment of Model Uncertainty	115

iv *Contents*

3.5.3	Extrapolation Properties of Models	116
3.5.4	Huge Data Sets	116
4	Conclusion and Discussion	118
4.1	Conclusion	118
4.2	Current State of the Art	121
4.3	Future Challenges and Research Topics.	122
	Appendices	135
A	Gradients of Frequently Used Policies	135
B	Weighted ML Estimates of Frequently Used Policies	136
C	Derivations of the Dual Functions for REPS	138

1

Introduction

From simple house-cleaning robots to robotic wheelchairs and general transport robots the number and variety of robots used in our everyday life are rapidly increasing. To date, the controllers for these robots are largely designed and tuned by a human engineer. Programming robots is a tedious task that requires years of experience and a high degree of expertise. The resulting programmed controllers are based on assuming exact models of both the robot's behavior and its environment. Consequently, hard-coding controllers for robots has its limitations when a robot has to adapt to new situations or when the robot/environment cannot be modeled sufficiently accurately. Hence, there is a gap between the robots currently used and the vision of incorporating fully autonomous robots. In *robot learning*, machine learning methods are used to automatically extract relevant information from data to solve a robotic task. Using the power and flexibility of modern machine learning techniques, the field of robot control can be further automated, and the gap toward autonomous robots, e.g., for general assistance in households, elderly care, and public services can be narrowed substantially.

1.1 Robot Control as a Reinforcement Learning Problem

In most tasks, robots operate in a high-dimensional state space \mathbf{x} composed of both internal states (e.g., joint angles, joint velocities, end-effector pose, and body position/orientation) and external states (e.g., object locations, wind conditions, or other robots). The robot selects its motor commands \mathbf{u} according to a control policy π . The control policy can either be stochastic, denoted by $\pi(\mathbf{u}|\mathbf{x})$, or deterministic, which we will denote as $\mathbf{u} = \pi(\mathbf{x})$. The motor commands \mathbf{u} alter the state of the robot and its environment according to the probabilistic transition function $p(\mathbf{x}_{t+1}|\mathbf{x}_t, \mathbf{u}_t)$. Jointly, the states and actions of the robot form a *trajectory* $\boldsymbol{\tau} = (\mathbf{x}_0, \mathbf{u}_0, \mathbf{x}_1, \mathbf{u}_1, \dots)$, which is often also called a *roll-out* or a *path*.

We assume that a numeric scoring system evaluates the performance of the robot system during a task and returns an accumulated reward signal $R(\boldsymbol{\tau})$ for the quality of the robot’s trajectory. For example, the reward $R(\boldsymbol{\tau})$ may include a positive reward for a task achievement and negative rewards, i.e., costs, that punish energy consumption. Many of the considered motor tasks are stroke-based movements, such as returning a tennis ball or throwing darts. We will refer to such tasks as *episodic learning tasks* as the execution of the task, the *episode*, ends after a given number T of time steps. Typically, the accumulated reward $R(\boldsymbol{\tau})$ for a trajectory is given as

$$R(\boldsymbol{\tau}) = r_T(\mathbf{x}_T) + \sum_{t=0}^{T-1} r_t(\mathbf{x}_t, \mathbf{u}_t), \quad (1.1)$$

where r_t is an instantaneous reward function, which might be a punishment term for the consumed energy, and r_T is a final reward, such as quadratic punishment term for the deviation to a desired goal posture. For many episodic motor tasks the policy is modeled as time-dependent policy, i.e., either a stochastic policy $\pi(\mathbf{u}_t|\mathbf{x}_t, t)$ or a deterministic policy $\mathbf{u}_t = \pi(\mathbf{x}_t, t)$ is used.

In some cases, the infinite-horizon case is considered

$$R(\boldsymbol{\tau}) = \sum_{t=0}^{\infty} \gamma^t r(\mathbf{x}_t, \mathbf{u}_t), \quad (1.2)$$

where $\gamma \in [0, 1)$ is a discount factor that discounts rewards further in the future.

Many tasks in robotics can be phrased as choosing a (locally) optimal control policy π^* that maximizes the expected accumulated reward

$$J_\pi = \mathbb{E}[R(\boldsymbol{\tau})|\pi] = \int R(\boldsymbol{\tau})p_\pi(\boldsymbol{\tau})d\boldsymbol{\tau}, \quad (1.3)$$

where $R(\boldsymbol{\tau})$ defines the objectives of the task, and $p_\pi(\boldsymbol{\tau})$ is the distribution over trajectories $\boldsymbol{\tau}$. For a stochastic policy $\pi(\mathbf{u}_t|\mathbf{x}_t, t)$, the trajectory distribution is given as

$$p_\pi(\boldsymbol{\tau}) = p(\mathbf{x}_0) \prod_{t=0}^{T-1} p(\mathbf{x}_{t+1}|\mathbf{x}_t, \mathbf{u}_t)\pi(\mathbf{u}_t|\mathbf{x}_t, t), \quad (1.4)$$

where $p(\mathbf{x}_{t+1}|\mathbf{x}_t, \mathbf{u}_t)$ is given by the system dynamics of the robot and its environment. For a deterministic policy, $p_\pi(\boldsymbol{\tau})$ is given as

$$p_\pi(\boldsymbol{\tau}) = p(\mathbf{x}_0) \prod_{t=0}^{T-1} p(\mathbf{x}_{t+1}|\mathbf{x}_t, \pi(\mathbf{x}_t, t)). \quad (1.5)$$

With this general reinforcement learning (RL) problem set-up, many tasks in robotics can be naturally formulated as *Reinforcement Learning* (RL) problems. However, robot RL poses three main challenges, which have to be solved: The RL algorithm has to manage (i) high-dimensional continuous state and action spaces, (ii) strong real-time requirements, and (iii) the high costs of robot interactions with its environment.

Traditional methods in RL, such as TD-learning [80], typically try to estimate the expected long-term reward of a policy for each state \mathbf{x} and time step t , also called the *value function* $V_t^\pi(\mathbf{x})$. The value function is used to calculate the quality of an executing action \mathbf{u} in state \mathbf{x} . This quality assessment is subsequently utilized to directly compute the policy by action selection or to update the policy π . However, value function methods struggle with the challenges encountered in robot RL, as these approaches require filling the complete state-action space with data. In addition, the value function is computed iteratively by the use of bootstrapping, which often results in a bias in

4 Introduction

the quality assessment of the state action pairs if we need to resort to value function approximation techniques as it is the case for continuous state spaces. Consequently, value function approximation turns out to be a very difficult problem in high-dimensional state and action spaces. Another major issue is that value functions are often discontinuous, especially when the non-myopic policy differs from a myopic policy. For instance, the value function of the under-powered pendulum swing-up is discontinuous along the manifold where the applicable torque is just not sufficient to swing the pendulum up [22]. Any error in the value function will eventually propagate through to the policy.

In a classical RL set-up, we seek a policy without too specific prior information. Key to successful learning is the exploration strategy of the learner to discover rewarding states and trajectories. In a robotics context, arbitrary exploration is not desired if not discouraged since the robot can easily be damaged. Therefore, the classical RL paradigm in a robotics context is not directly applicable since exploration needs to take hardware constraints into account. Two ways of implementing cautious exploration are to either avoid significant changes in the policy [57] or to explicitly discourage entering undesired regions in the state space [21].

In contrast to value-based methods, *Policy Search* (PS) methods use parametrized policies π_{θ} . They directly operate in the parameter space Θ , $\theta \in \Theta$, of parametrized policies, and typically avoid learning a value function. Many methods do so by directly using the experienced reward to come from the rollouts as quality assessment for state action pairs instead of using the rather dangerous bootstrapping used in value-function approximation. The usage of parametrized policies allows for scaling RL into high dimensional continuous action spaces by reducing the search space of possible policies.

Policy search allows task-appropriate pre-structured policies, such as movement primitives [71], to be integrated straightforwardly. Additionally, imitation learning from an expert’s demonstrations can be used to obtain an initial estimate for the policy parameters [58]. Finally, by selecting a suitable policy parametrization, stability and robustness guarantees can be given [11]. All these properties simplify the robot learning problem and permit the successful application of

reinforcement learning to robotics. Therefore, PS is often the RL approach of choice in robotics since it is better at coping with the inherent challenges of robot reinforcement learning. Over the last decade, a series of fast policy search algorithms have been proposed and shown to work well on real systems [38, 53, 58, 86, 7, 21, 17]. In this review, we provide a general overview, summarize the main concepts behind current policy search approaches, and discuss relevant robot applications of these policy search methods. We focus mainly on those aspects of RL that are predominant for robot learning, i.e., learning in high-dimensional continuous state and action spaces and a high data-efficiency and local exploration. Other important aspects of RL, such as the exploration-exploitation trade-off, feature selection, using structured models or value function approximation are not covered in this monograph.

1.2 Policy Search Taxonomy

Numerous policy search methods have been proposed in the last decade, and several of them have been used successfully in the domain of robotics. In this monograph, we review several important recent developments in policy search for robotics. We distinguish between model-free policy search methods (Section 2), which learn policies directly based on sampled trajectories, and model-based approaches (Section 3), which use the sampled trajectories to first build a model of the state dynamics, and, subsequently, use this model for policy improvement.

Figure 1.1 categorizes policy search into model-free policy search and model-based policy search and distinguishes between different policy update strategies. The policy updates in both model-free and model-based policy search (green blocks) are based on either policy gradients (PG), expectation maximization (EM)-based updates, or information-theoretic insights (Inf.Th.). While all three update strategies are fairly well explored in model-free policy search, model-based policy search almost exclusively focuses on PG to update the policy.

Model-free policy search uses stochastic trajectory generation, i.e., the trajectories are generated by “sampling” from the robot $p(\mathbf{x}_{t+1}|\mathbf{x}_t, \mathbf{u}_t)$ and the policy π_θ . This means, a system model is not

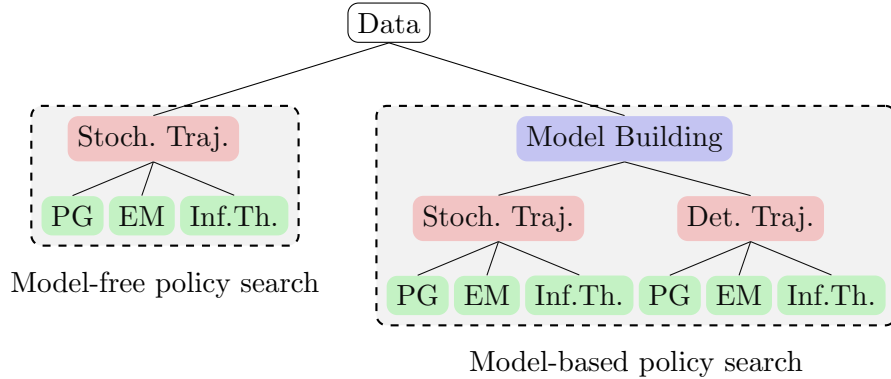


Fig. 1.1 Categorization of policy search into model-free policy search and model-based policy search. In the model-based case (right sub-tree), data from the robot is used to learn a model of the robot (blue box). This model is then used to generate trajectories. Here, we distinguish between stochastic trajectory generation and deterministic trajectory prediction. Model-free policy search (left sub-tree) uses data from the robot directly as a trajectory for updating the policy. The policy updates in both model-free and model-based policy search (green blocks) are based on either policy gradients (PG), expectation maximization (EM)-based updates, or information-theoretic insights (Inf.Th.).

explicitly required; We just have to be able to sample trajectories from the real robot. In the model-based case (right sub-tree), we can either use stochastic trajectory generation or deterministic trajectory prediction. In the case of stochastic trajectory generation, the learned models are used as simulator for sampling trajectories. Hence, learned models can easily be combined with model-free policy search approaches by exchanging the “robot” with the learned model of the robot’s dynamics. Deterministic trajectory prediction does not sample trajectories, but analytically predicts the trajectory distribution $p_{\theta}(\tau)$. Typically, deterministic trajectory prediction is computationally more involved than sampling trajectories from the system. However, for the subsequent policy update, deterministic trajectory prediction can allow for analytic computation of gradients, which can be advantageous over stochastic trajectory generation, where these gradients can only be approximated.

1.2.1 Model-free and Model-based Policy Search

Model-free policy search methods use real robot interactions to create sample trajectories $\tau^{[i]}$. While sampling trajectories is relatively straightforward in computer simulation, when working with robots, the generation of each “sample” typically needs some level of human supervision. Consequently, trajectory generation with the real system is considerably more time consuming than working with simulated systems. Furthermore, real robot interactions causes wear and tear in non-industrial robots. However, in spite of the relatively high number of required robot interactions for model-free policy search, learning a policy is often easier than learning accurate forward models, and, hence, model-free policy search is more widely used than model-based methods.

Model-based policy search methods attempt to address the problem of sample inefficiency by using the observed trajectories $\tau^{[i]}$ to learn a forward model of the robot’s dynamics and its environment. Subsequently, this forward model is used for *internal* simulations of the robot’s dynamics and environment, based on which the policy is learned. Model-based PS methods have the potential to require fewer interactions with the robot and to efficiently generalize to unforeseen situations [6]. While the idea of using models in the context of robot learning is well-known since the 1980s [2], it has been limited by its strong dependency on the quality of the learned models. In practice, the learned model is *not* exact, but only a more or less accurate approximation to the real dynamics. Since the learned policy is inherently based on internal simulations with the learned model, inaccurate models can, therefore, lead to control strategies that are not robust to model errors. In some cases, learned models may be physically implausible and contain negative masses or negative friction coefficients. These implausible effects are often exploited by the policy search algorithm, resulting in a poor quality of the learned policy. This effect can be alleviated by using models that explicitly account for model errors [72, 20]. We will discuss such methods in Section 3.

1.3 Typical Policy Representations

Typical policy representations, which are used for policy search can be categorized into time-independent representations $\pi(\mathbf{x})$ and time-dependent representations $\pi(\mathbf{x}, t)$. Time-independent representations use the same policy for all time steps, and, hence, often require a complex parametrization. Time-dependent representations can use different policies for different time steps, allowing for a potentially simpler structure of the individual policies can be used.

We will describe all policy representations in their deterministic formulation $\pi_{\boldsymbol{\theta}}(\mathbf{x}, t)$. In stochastic formulations, typically a zero-mean Gaussian noise vector $\boldsymbol{\epsilon}_t$ is added to $\pi_{\boldsymbol{\theta}}(\mathbf{x}, t)$. In this case, the parameter vector $\boldsymbol{\theta}$ typically also includes the (co)variance matrix used for generating the noise $\boldsymbol{\epsilon}_t$. In robot learning, the three main policy representations are linear policies, radial basis function networks, and dynamic movement primitives [71].

Linear Policies. Linear controllers are the most simple time-independent representation. The policy π is a linear policy

$$\pi_{\boldsymbol{\theta}}(\mathbf{x}) = \boldsymbol{\theta}^T \boldsymbol{\phi}(\mathbf{x}) \quad (1.6)$$

where $\boldsymbol{\phi}$ is a basis function vector. This policy only depends linearly on the policy parameters. However, specifying the basis functions by hand is typically a difficult task, and, hence, the application of linear controllers is limited to problems where appropriate basis functions are known, e.g., for balancing tasks, the basis functions are typically given by the state variables of the robot.

Radial Basis Functions Networks. A typical nonlinear time-independent policy representation is a radial basis function (RBF) network. An RBF policy $\pi_{\boldsymbol{\theta}}(\mathbf{x})$ is given as

$$\pi_{\boldsymbol{\theta}}(\mathbf{x}) = \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}), \quad \phi_i(\mathbf{x}) = \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_i)^T \mathbf{D}_i(\mathbf{x} - \boldsymbol{\mu}_i)\right), \quad (1.7)$$

where $\mathbf{D}_i = \text{diag}(\mathbf{d}_i)$ is a diagonal matrix. Unlike in the linear policy case, the parameters $\boldsymbol{\beta} = \{\boldsymbol{\mu}_i, \mathbf{d}_i\}_{i=1, \dots, n}$ of the basis functions themselves are now considered as free parameters that need to be learned.

Hence, the parameter vector θ of the policy is given by $\theta = \{\mathbf{w}, \beta\}$. While RBF-networks are powerful policy representations, they are also difficult to learn due to the high number of nonlinear parameters. Furthermore, as RBF networks are local representations, they are hard to scale to high-dimensional state spaces.

Dynamic Movement Primitives. Dynamic Movement Primitives (DMPs) are the most widely used time-dependent policy representation in robotics [71, 31]. DMPs use non-linear dynamical systems for generating the movement of the robot. The key principle of DMPs is to use a linear spring-damper system which is modulated by a non-linear forcing function f_t , i.e.,

$$\ddot{y}_t = \tau^2 \alpha_y (\beta_y (g - y_t) - \dot{y}_t) + \tau^2 f_t, \quad (1.8)$$

where the variable y_t directly specifies the desired joint position of the robot. The parameter τ is the time-scaling coefficient of the DMP, the coefficients α_y and β_y define the spring and damping constants of the spring-damper system and the goal-parameter g is the unique point-attractor of the spring-damper system. Note that the spring-damper system is equivalent to a standard linear PD-controller that operates on a linear system with zero desired velocity, i.e.,

$$\ddot{y}_t = k_p (g - y_t) - k_d \dot{y}_t,$$

where the P-gain is given by $k_p = \tau^2 \alpha_y \beta_y$ and the D-gain by $k_d = \tau^2 \alpha_y$. The forcing function f_t changes the goal attractor g of the linear PD-controller.

One key innovation of the DMP approach is the use of a phase variable z_t to scale the execution speed of the movement. The phase variable evolves according to $\dot{z} = -\tau \alpha_z z$. It is initially set to $z = 1$ and exponentially converges to 0 as $t \rightarrow \infty$. The parameter α_z specifies the speed of the exponential decline of the phase variable. The variable τ can be used to temporally scale the evolution of the phase z_t , and, thus, the evolution of the spring-damper system as shown in Equation (1.8). For each degree of freedom, an individual spring damper system, and, hence, an individual forcing function f_t is used. The function f_t depends on the phase variable, i.e., $f_t = f(z_t)$ and is constructed by the

weighted sum of K basis functions ϕ_i

$$f(z) = \frac{\sum_{i=1}^K \phi_i(z) w_i}{\sum_{i=1}^K \phi_i(z)} z, \quad \phi_i(z) = \exp\left(-\frac{1}{2\sigma_i^2}(z - c_i)^2\right). \quad (1.9)$$

The parameters w_i are denoted as ‘shape-parameters’ of the DMP as they modulate the acceleration profile, and, hence, indirectly specify the shape of the movement. From Equation (1.9), we can see that the basis-functions are multiplied with the phase variable z , and, hence, f_t vanishes as $t \rightarrow \infty$. Consequently, the non-linear dynamical system is globally stable as it behaves like a linear spring damper system for $t \rightarrow \infty$. From this argument, we can also conclude that the goal parameter g specifies the final position of the movement while the shape parameters w_i specify how to reach this final position.

Integrating the dynamical systems for each DoF results in a desired trajectory $\boldsymbol{\tau}^* = \{\mathbf{y}_t\}_{t=0..T}$ that is, subsequently, followed by feedback control laws [56]. The policy $\pi_{\boldsymbol{\theta}}(\mathbf{x}_t, t)$ that is specified by a DMP, directly controls the acceleration of the joint, and, hence, is given by

$$\pi_{\boldsymbol{\theta}}(\mathbf{x}_t, t) = \tau^2 \alpha_y (\beta_y (g - y_t) - \dot{y}_t) + \tau^2 f(z_t).$$

Note that the DMP policy is linear in the shape parameters \mathbf{w} and the goal attractor g , but non-linear in the time scaling constant τ .

The parameters $\boldsymbol{\theta}$ used for learning a DMP are typically given by the weight parameters w_i , but might also contain the goal parameters g as well as the temporal scaling parameter τ . In addition, the DMP approach has been extended in [36] such that the desired final velocity \dot{g} of the joints can also be modulated. Such modulation is, for example, useful for learning hitting movements in robot table-tennis. Typically, $K = 5$ to 20 basis functions are used, i.e., 5 to 20 shape weights per degree of freedom of the robot are used.

Miscellaneous Representations. Other representations that have been used in the literature include central pattern generators for robot walking [24] and feed-forward neural networks, which have been used mainly in simulation [30, 89].

1.4 Outline

The structure of this monograph is as follows: In Section 2, we give a detailed overview of model-free policy search methods, where we classify policy search algorithms according to their policy evaluation, policy update, and exploration strategy. For the policy update strategies, we will follow the taxonomy in Figure 1.1 and discuss policy gradient methods, EM-based approaches, information-theoretic approaches. Additionally, we will discuss miscellaneous important methods such as stochastic optimization and policy search approaches based on the path integral theory. Policy search algorithms can either use a step-based or episode-based policy evaluation strategy. Most policy update strategies presented in Figure 1.1 can be used for both, step-based and episode-based policy evaluation. We will present both types of algorithms if they have been introduced in the literature. Subsequently, we will discuss different exploration strategies for model-free policy search and conclude this section with robot applications of model-free policy search. Section 3 surveys model-based policy search methods in robotics. Here, we introduce two models that are commonly used in policy search: locally weighted regression and Gaussian processes. Furthermore, we detail stochastic and deterministic inference algorithms to compute a probability distribution $p_\pi(\boldsymbol{\tau})$ over trajectories (see the red boxes in Figure 1.1). We conclude this section with examples of model-based policy search methods and their application to robotic systems. In Section 4, we give recommendations for the practitioner and conclude this monograph.

2

Model-free Policy Search

Model-free policy search (PS) methods update the policy directly based on sampled trajectories $\tau^{[i]}$, where i denotes the index of the trajectory, and the obtained immediate rewards $r_0^{[i]}, r_1^{[i]}, \dots, r_T^{[i]}$ for the trajectories. Model-free PS methods try to update the parameters θ such that trajectories with higher rewards become more likely when following the new policy, and, hence, the average return

$$J_{\theta} = \mathbb{E}[R(\tau)|\theta] = \int R(\tau)p_{\theta}(\tau)d\tau \quad (2.1)$$

increases. Learning a policy is often easier than learning a model of the robot and its environment, and, hence, model-free policy search methods are used more frequently than model-based policy search methods. We categorize model-free policy search approaches based on their policy evaluation strategies, their policy update strategies [58, 57] and their exploration strategies [67, 38].

The exploration strategy determines how new trajectories are created for the subsequent policy evaluation step. The exploration strategy is essential for efficient model-free policy search, as, we need variability in the generated trajectories to determine the policy update, but an excessive exploration is also likely to damage the robot. Most model-free

Algorithm 1 Model-Free Policy Search

repeat

 Explore: Generate trajectories $\tau^{[i]}$ using policy π_k

 Evaluate: Assess quality of trajectories or actions

 Update: Compute π_{k+1} given trajectories $\tau^{[i]}$ and evaluations

until Policy converges $\pi_{k+1} \approx \pi_k$

methods therefore use a stochastic policy for exploration which explores only locally. Exploration strategies can be categorized into step-based and episode-based exploration strategies. While step-based exploration uses an exploratory action in each time step, episode-based exploration directly changes the parameter vector θ of the policy only at the start of the episode.

The policy evaluation strategy decides how to evaluate the quality of the executed trajectories. Here we can again distinguish between step-based and episode-based evaluations. Step-based evaluation strategies decompose the trajectory τ in its single steps $(\mathbf{x}_0, \mathbf{u}_0, \mathbf{x}_1, \mathbf{u}_1, \dots)$ and aim at evaluating the quality of single actions. In comparison, episode-based evaluation directly uses the returns of the whole trajectories to evaluate the quality of the used policy parameters θ .

Finally, the policy update strategy uses the quality assessment of the evaluation strategy to determine the policy update. Update strategies can be classified according to the optimization method employed by the PS algorithm. While the most common update strategies are based on gradient ascent, resulting in policy gradient methods [89, 58, 61], inference-based approaches use expectation maximization [38, 48] and information theoretic approaches [57, 17] use insights from information theory to update the policy. We will also cover additional important methods such as path-integral approaches and stochastic optimization. Model-free policy search can be applied to policies with a moderate number of parameters, i.e., up to a few hundred parameters. Most applications use linear policy representations such as linear controllers or dynamical movement primitives that have been discussed in Section 1.3.

In the following section, we will discuss the used exploration strate-

gies in current algorithms. Subsequently, we will cover policy evaluation strategies in more detail. Finally, we will review policy update methods such as policy gradients, inference/EM-based, and information theoretic policy updates as well as update strategies based on path integrals. Many policy update strategies have been implemented for both policy evaluation approaches, and, hence, we will discuss the combinations that have been explored so far. We will conclude with presenting the most interesting experimental results for policy learning for robots.

2.1 Exploration Strategies

The exploration strategy is used to generate new trajectory samples $\tau^{[i]}$ which are subsequently evaluated by the policy evaluation strategy and used for the policy update. An efficient exploration is, therefore, crucial for the performance of the resulting policy search algorithm. All exploration strategies considered for model-free policy search are local and use stochastic policies to implement exploration. Typically, Gaussian policies are used to model such stochastic policies.

We distinguish between exploration in action space versus exploration in parameter space, step-based versus episode-based exploration strategies and correlated versus uncorrelated exploration noise.

2.1.1 Exploration in Action Space versus Exploration in Parameter Space

Exploration in the action space is implemented by adding an exploration noise $\epsilon_{\mathbf{u}}$ directly to the executed actions, i.e. $\mathbf{u}_t = \boldsymbol{\mu}(\mathbf{x}, t) + \epsilon_{\mathbf{u}}$. The exploration noise is in most cases sampled independently for each time step from a zero-mean Gaussian distribution with covariance $\boldsymbol{\Sigma}_{\mathbf{u}}$. The policy $\pi_{\boldsymbol{\theta}}(\mathbf{u}|\mathbf{x})$ is, therefore, given as

$$\pi_{\boldsymbol{\theta}}(\mathbf{u}|\mathbf{x}) = \mathcal{N}(\mathbf{u}|\boldsymbol{\mu}_{\mathbf{u}}(\mathbf{x}, t), \boldsymbol{\Sigma}_{\mathbf{u}}).$$

Exploration in the action space is one of the first exploration strategies used in the literature [89, 79, 9, 62] and used for most policy gradient approaches such as the REINFORCE algorithm [89] or the eNAC algorithm [61].

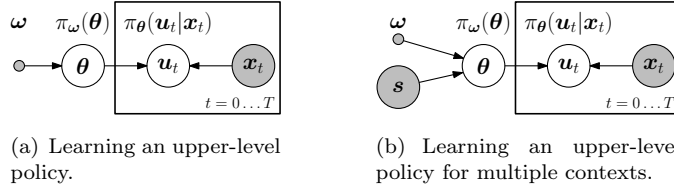


Fig. 2.1 (a) Graphical model for learning an upper-level policy $\pi_\omega(\theta)$. The upper level policy chooses the parameters θ of the lower-level policy $\pi_\theta(\mathbf{u}|\mathbf{x})$ that is executed on the robot. The parameters of the upper-level policy are given by ω . (b) Learning an upper-level policy $\pi_\omega(\theta|\mathbf{s})$ for multiple contexts \mathbf{s} . The context is used to select the parameters θ , but typically not be the lower-level policy itself. The lower-level policy $\pi_\theta(\mathbf{u}|\mathbf{x})$ is typically modeled as a deterministic policy in both settings.

Exploration strategies in parameter space perturb the parameter vector θ . This exploration can either only be added in the beginning of an episode, or, a different perturbation of the parameter vector can be used at each time step [67, 38].

Learning Upper-Level Policies Both approaches can be formalized with the concept of an upper level policy $\pi_\omega(\theta)$ which selects the parameters of the actual control policy $\pi_\theta(\mathbf{u}|\mathbf{x})$ of the robot. Hence, we will denote the latter in this hierarchical setting as lower-level policy. The upper level policy $\pi_\omega(\theta)$ is typically modeled as a Gaussian distribution $\pi_\omega(\theta) = \mathcal{N}(\theta|\mu_\theta, \Sigma_\theta)$. The lower level control policy $\mathbf{u} = \pi_\theta(\mathbf{x}, t)$ is typically modeled as deterministic policy as exploration only takes place in the parameter space of the policy.

Instead of directly finding the parameters θ of the lower-level policy, we want to find the parameter vector ω which now defines a distribution over θ . Using a distribution over θ has the benefit that we can use this distribution to directly explore in parameter space. The optimization problem for learning upper-level policies can be formalized as maximizing

$$J_\omega = \int_{\theta} \pi_\omega(\theta) \int_{\tau} p(\tau|\theta) R(\tau) d\tau d\theta = \int_{\theta} \pi_\omega(\theta) R(\theta) d\theta. \quad (2.2)$$

The graphical model for learning an upper-level policy is given in Figure 2.1(a).

In the case of using a different parameter vector for each time step, typically, a linear control policy $u_t = \phi_t(\mathbf{x})^T \boldsymbol{\theta}$ is used. We can also rewrite the deterministic lower level control policy $\pi_{\boldsymbol{\theta}}(\mathbf{x}, t)$ in combination with the upper-level policy $\pi_{\boldsymbol{\omega}}(\boldsymbol{\theta})$ as a single, stochastic policy

$$\pi_{\boldsymbol{\theta}}(u_t | \mathbf{x}_t, t) = \mathcal{N}(u_t | \phi_t(\mathbf{x})^T \boldsymbol{\mu}_{\boldsymbol{\theta}}, \phi_t(\mathbf{x})^T \boldsymbol{\Sigma}_{\boldsymbol{\theta}} \phi_t(\mathbf{x})), \quad (2.3)$$

which follows from the properties of the expectation and the variance operators. Such exploration strategy is, for example, applied by the PoWER [38] and PI² [81] algorithms and was also suggested to be used for policy gradient algorithms in [67].

In contrast to exploration in action space, exploration in parameter space is able to use more structured noise and can adapt the variance of the exploration noise in dependence of the state features $\phi_t(\mathbf{x})$.

2.1.2 Episode-based versus Step-based Exploration

Step-based exploration strategies use different exploration noise at each time-step and can explore either in action space or in parameter space as we know from the discussion in the previous section. Episode-based exploration strategies use exploration noise only at the beginning of the episode, which naturally leads to an exploration in parameter space. Typically, episode-based exploration strategies are used in combination with episode-based policy evaluation strategies which are covered in the next section. However, episode-based exploration strategies are also realizable with step-based evaluation strategies such as with the PoWER [38] or with the PI² [81] algorithm.

Step-based exploration strategies can be problematic as they might produce action sequences which are not reproducible by the noise-free control law, and, hence, might again affect the quality of the policy updates. Furthermore, the effects of the perturbations are difficult to estimate as they are typically washed out by the system dynamics which acts as a low pass filter. Moreover, a step-based exploration strategy causes a large parameter variance which grows with the number of time steps. Such exploration strategies may even damage the robot as random exploration in every time step leads to jumps in the controls of the robot. Hence, fixing exploration for the whole episode or only

slowly vary the exploration by low-pass filtering the noise, is often beneficial in real robot applications. On the other hand, the stochasticity of the control policy $\pi_{\theta}(\mathbf{u}|\mathbf{x})$ also smoothens out the expected return, and, hence, in our experience, policy search is sometimes less prone to getting stuck in local minima using step-based exploration.

Episode-based exploration always produces action sequences which can be reproduced by the noise free control law. Fixing the exploration noise in the beginning of the episode might also decrease the variance of the quality assessment estimated by the policy evaluation strategy, and, hence, might produce more reliable policy updates [77].

2.1.3 Uncorrelated versus Correlated Exploration

As most policies are represented as Gaussian distributions, uncorrelated exploration noise is obtained by using a diagonal covariance matrix while we can achieve correlated exploration by maintaining a full representation of the covariance matrix. Exploration strategies in action space typically use a diagonal covariance matrix. For exploration strategies in parameter space, many approaches can also be used to update the full covariance matrix of the Gaussian policy. Such an approach was first introduced by the stochastic optimization algorithm CMA-ES [28] and was later also incorporated into more recent policy search approaches such as REPS [57, 17], PoWER [38], and PI² [77].

Using the full covariance matrix often results in a considerably increased learning speed for the resulting policy search algorithm [77]. Intuitively, the covariance matrix can be interpreted as a second order term. Similar to the Hessian in second order optimization approaches, it regulates the step-width of the policy update for each direction of the parameter space. However, estimating the full covariance matrix can also be difficult [68] if the parameter space becomes high dimensional ($|\Theta| > 50$) as the covariance matrix has $O(|\Theta|^2)$ elements. In this case, too many samples are needed for an accurate estimate of the covariance matrix.

2.1.4 Updating the Exploration Distribution

Many model-free policy search approaches also update the exploration distribution, and, hence, the covariance of the Gaussian policy. Updating the exploration distribution often improves the performance as different exploration rates can be used throughout the learning process. Typically, a large exploration rate can be used in the beginning of learning which is then gradually decreased to fine tune the policy parameters. In general, the exploration rate tends to be decreased too quickly for many algorithms, and, hence, the exploration policy might collapse to almost a single point. In this case, the policy update might stop improving prematurely. This problem can be alleviated by the use of an information theoretic update metric, which limits the relative entropy between the new and the old trajectory distribution. Such an information theoretic measure is, for example, used by the natural policy gradient methods as well as by the REPS algorithm [57, 17]. Peters and Schaal [61] showed that, due to the bounded relative entropy to the old trajectory distribution, the exploration does not collapse prematurely, and hence, a more stable learning progress can be achieved. Still, the exploration policy might collapse prematurely if initialized only locally. Some approaches artificially add an additional variance term to the covariance matrix of the exploration policy after each policy update to sustain exploration [77], however, a principled way of adapting the exploration policy in such situations is missing so far in the literature.

2.2 Policy Evaluation Strategies

Policy evaluation strategies are used to assess the quality of the executed policy. Policy search algorithms either try to assess the quality of single state-action pairs $\mathbf{x}_t, \mathbf{u}_t$, which we will refer to as step-based evaluations, or the quality of a parameter vector θ that has been used during the whole episode, which we will refer to as episode-based policy evaluation. The policy evaluation strategy is used to transform sampled trajectories $\tau^{[i]}$ into a data-set \mathcal{D} that contains samples of either the state-action pairs $\mathbf{x}_t^{[i]}, \mathbf{u}_t^{[i]}$ or the parameter vectors $\theta^{[i]}$ and an evalua-

tion of these samples, as will be described in this section. The data-set \mathcal{D} is subsequently processed by the policy update strategy to determine the new policy.

2.2.1 Step-Based Policy Evaluation

In step-based policy evaluation, we decompose the sampled trajectories $\tau^{[i]}$ into its single time steps $\mathbf{x}_t^{[i]}, \mathbf{u}_t^{[i]}$, and estimate the quality of the single actions. The quality of an action is given by the expected accumulated future reward when executing $\mathbf{u}_t^{[i]}$ in state $\mathbf{x}_t^{[i]}$ at time step t and subsequently following policy $\pi_{\theta}(\mathbf{u}|\mathbf{x})$,

$$Q_t^{[i]} = Q_t^{\pi}(\mathbf{x}_t^{[i]}, \mathbf{u}_t^{[i]}) = \mathbb{E}_{p_{\theta}(\tau)} \left[\sum_{h=t}^T r_h(\mathbf{x}_h, \mathbf{u}_h) \middle| \mathbf{x}_t = \mathbf{x}_t^{[i]}, \mathbf{u}_t = \mathbf{u}_t^{[i]} \right],$$

which corresponds to the state-action value function Q^{π} . However, estimating the state-action value function is a difficult problem in high-dimensional continuous spaces and often suffers from approximation errors or a bias induced by the bootstrapping approach used by most value function approximation methods. Therefore, most policy search methods rely on Monte-Carlo estimates of $Q_t^{[i]}$ instead of using value function approximations. Monte-Carlo estimates are unbiased, however, they typically exhibit a high variance. Fortunately, most methods can cope with noisy estimates of $Q_t^{[i]}$, and, hence, solely the reward to come for the current trajectory $\tau^{[i]}$ is used $Q_t^{[i]} \approx \sum_{h=t}^T r_h^{[i]}$, which avoids the additional averaging that would be needed for an accurate Monte-Carlo estimate. Algorithms based on step-based policy evaluation use a data set $\mathcal{D}_{\text{step}} = \{\mathbf{x}^{[i]}, \mathbf{u}^{[i]}, Q^{[i]}\}$ to determine the policy update step. Some step-based policy search algorithms [89, 58] also use the returns $R^{[i]} = \mathbb{E}_{p_{\theta}(\tau)} \left[\sum_{h=0}^T r_h^{[i]} \right]$ of the whole episode as evaluation for the single actions of the episode. However, as the estimate of the returns suffers from a higher variance than the reward to come, such a strategy is not preferable. Pseudo-code for a general step-based policy evaluation algorithm is given in Algorithm 2.

Algorithm 2 Policy Search with Step-Based Policy Evaluation

repeat**Exploration:**Create samples $\tau^{[i]} \sim \pi_{\theta_k}(\tau)$ following $\pi_{\theta_k}(\mathbf{u}|\mathbf{x})$, $i = 1 \dots N$ **Policy Evaluation:**Evaluate actions: $Q_t^{[i]} = \sum_{h=t}^T r_h^{[i]}$ for all t and all i Compose data set: $\mathcal{D}_{\text{step}} = \left\{ \mathbf{x}_t^{[i]}, \mathbf{u}_t^{[i]}, Q_t^{[i]} \right\}_{i=1 \dots N, t=1 \dots T-1}$ **Policy Update:**Compute new policy parameters θ_{k+1} using \mathcal{D} .Algorithms: REINFORCE, G(PO)MDP, NAC, eNAC,
PoWER, PI²**until** Policy converges $\theta_{k+1} \approx \theta_k$

2.2.2 Episode-Based Policy Evaluation

Episode-based update strategies [77, 78, 17] directly use the expected return $R^{[i]} = R(\theta^{[i]})$ to evaluate the quality of a parameter vector $\theta^{[i]}$. Typically, the expected return is given by the sum of the future immediate rewards, i.e.,

$$R(\theta^{[i]}) = \mathbb{E}_{p_{\theta}(\tau)} \left[\sum_{t=0}^T r_t | \theta = \theta^{[i]} \right]. \quad (2.4)$$

However, episode-based algorithms are not restricted to this structure of the return, but can use any reward function $R(\theta^{[i]})$ which depends on the resulting trajectory of the robot. For example, when we want to learn to throw a ball to a desired target location, the reward $R(\theta^{[i]})$ can intuitively be defined as the negative minimum distance of the ball to the target location [41]. Such reward function can not be described by a sum of immediate rewards.

The expected return $R^{[i]}$ for $\theta^{[i]}$ can be estimated by performing multiple roll-outs on the real system. However, in order to avoid such an expensive operation, a few approaches [68, 38] can cope with noisy estimates of $R^{[i]}$, and, hence, can directly use the return $\sum_{t=0}^T r_t^{[i]}$ of a single trajectory $\tau^{[i]}$ to estimate $R^{[i]}$. Other algorithms, such as stochas-

Algorithm 3 Episode-based Policy Evaluation for Learning an Upper-Level Policy

repeat

Exploration:

Sample parameter vector $\boldsymbol{\theta}^{[i]} \sim \pi_{\boldsymbol{\omega}_k}(\boldsymbol{\theta})$, $i = 1 \dots N$

Sample trajectory $\boldsymbol{\tau}^{[i]} \sim p_{\boldsymbol{\theta}^{[i]}}(\boldsymbol{\tau})$

Policy Evaluation:

Evaluate policy parameters $R^{[i]} = \sum_{t=0}^T r_t^{[i]}$

Compose data set $\mathcal{D}_{\text{ep}} = \left\{ \boldsymbol{\theta}^{[i]}, R^{[i]} \right\}_{i=1 \dots N}$

Policy Update:

Compute new policy parameters $\boldsymbol{\omega}_{k+1}$ using \mathcal{D}_{ep}

Algorithms: Episode-based REPS, Episode-based PI²

PEPG, NES, CMA-ES, RWR

until Policy converges $\boldsymbol{\omega}_{k+1} \approx \boldsymbol{\omega}_k$

tic optimizers, require a more accurate estimate of $R^{[i]}$, and, thus, either require multiple roll-outs, or suffer from a bias in the subsequent policy update step. Episode-based policy evaluation produces a dataset $\mathcal{D}_{\text{ep}} = \{\boldsymbol{\theta}^{[i]}, R^{[i]}\}_{i=1 \dots N}$, which is subsequently used for the policy updates. Episode-based policy evaluation is typically connected with parameter-based exploration strategies, and, hence, such algorithms can be formalized by the problem of learning an upper-level policy $\pi_{\boldsymbol{\omega}}(\boldsymbol{\theta})$, see Section 2.1.1. The general algorithm for policy search with episode-based policy evaluation is given in Algorithm 3.

An underlying problem of episode-based evaluation is the variance of the $R^{[i]}$ estimates. The variance depends on the stochasticity of the system, the stochasticity of the policy, and the number of time steps, consequently, for a high number of time-steps and highly stochastic systems, step-based algorithms should be preferred. In order to reduce the variance, the policy $\pi_{\boldsymbol{\theta}}(\mathbf{u}|\mathbf{x})$ is often modeled as a deterministic policy and exploration is directly performed in parameter space.

2.2.3 Comparison of Step- and Episode-based Evaluation

Step-based policy evaluation exploits the structure that the return is typically composed of the sum of the immediate rewards. Single actions can now be evaluated by the reward to come in that episode, instead of the whole reward of the episode, and, hence, the variance of the evaluation can be significantly reduced as the reward to come contains less random variables as the total reward of the episode. In addition, as we evaluate single actions instead of the whole parameter vector, the evaluated samples can be used more efficiently as several parameter-vectors θ might produce a similar action \mathbf{u}_t at time step t . Most policy search algorithms, such as the common policy gradient algorithms [89, 61], the PoWER [38] algorithm or the PI² [81] algorithm, employ such a strategy. A drawback of most step-based updates is that they often rely on a linear parametrization of the policy $\pi_{\theta}(\mathbf{u}|\mathbf{x})$. They also cannot be applied if the reward is not decomposable into isolated time steps.

Episode-based policy evaluation strategies do not decompose the returns, and, hence, might suffer from a large variance of the estimated returns. However, episode-based policy evaluation strategies typically employ more sophisticated exploration strategies which directly explore in the parameter space of the policy [17, 77, 29], and, thus, can often compete with their step-based counter-parts. So far, there is no clear answer as to which of the strategies should be preferred. The choice of the methods often depends on the problem at hand.

2.3 Important Extensions

In this section we will cover two important extensions of model-free policy search, generalization to multiple tasks and learning multiple solutions to the same task. We will introduce the relevant concepts for both extensions, however, the detailed algorithms will be covered in Section 2.4 which covers the policy update strategies.

2.3.1 Generalization to Multiple Tasks

For generalizing the learned policies to multiple tasks, so far, mainly episode-based policy evaluation strategies have been used which learn

an upper level policy. We will characterize a task by a context vector \mathbf{s} . The context describes all variables which do not change during the execution of the task but might change from task to task. For example, the context \mathbf{s} can describe the objectives of the agent or physical properties such as a mass to lift. The upper level policy is extended to generalize the lower-level policy $\pi_{\theta}(\mathbf{u}|\mathbf{x})$ to different tasks by conditioning the upper-level policy $\pi_{\omega}(\theta|\mathbf{s})$ on the context \mathbf{s} . The problem of learning $\pi_{\omega}(\theta|\mathbf{s})$ can be characterized by maximizing the expected returns over all contexts, i.e.,

$$J_{\omega} = \int_{\mathbf{s}} \mu(\mathbf{s}) \int_{\theta} \pi_{\omega}(\theta|\mathbf{s}) \int_{\tau} p(\tau|\theta, \mathbf{s}) R(\tau, \mathbf{s}) d\tau d\theta d\mathbf{s} \quad (2.5)$$

$$= \int_{\mathbf{s}} \mu(\mathbf{s}) \int_{\theta} \pi_{\omega}(\theta|\mathbf{s}) R(\theta, \mathbf{s}) d\theta d\mathbf{s}, \quad (2.6)$$

where $R(\theta, \mathbf{s}) = \int_{\tau} p(\tau|\theta, \mathbf{s}) R(\tau, \mathbf{s}) d\tau$ is the expected return for executing the lower-level policy with parameter vector θ in context \mathbf{s} and $\mu(\mathbf{s})$ is the distribution over the contexts. The trajectory distribution $p(\tau|\theta, \mathbf{s})$ can now also depend on the context, as the context can contain physical properties of the environment. We also extend the data set $\mathcal{D}_{\text{ep}} = \left\{ \mathbf{s}^{[i]}, \theta^{[i]}, R^{[i]} \right\}_{i=1 \dots N}$ used for updating the policy, which now also includes the corresponding context $\mathbf{s}^{[i]}$ that have been active for executing the lower-level policy with parameters $\theta^{[i]}$. The graphical model for learning an upper-level policies with multiple contexts is given in Figure 2.1(b) and the general algorithm is given in Algorithm 4. Algorithms that can generalize the lower level policy to multiple contexts include the Reward Weighted Regression (RWR) algorithm [59] the Cost-Regularized Regression (CrKR) algorithm [37] and the episode-based relative entropy policy search (REPS) algorithm [17]. RWR and CrKR are covered in Section 2.4.2.3 and REPS in Section 2.4.3.1.

2.3.2 Learning Multiple Solutions for a Single Motor Task

Many motor tasks can be solved in multiple ways. For example, for returning a tennis ball, in many situations we can either use a back-hand or fore-hand stroke. Hence, it is desirable to find algorithms

Algorithm 4 Learning an Upper-Level Policy for multiple Tasks

repeat**Exploration:**Sample context $\mathbf{s}^{[i]} \sim \mu(\mathbf{s})$ Sample parameter vector $\boldsymbol{\theta}^{[i]} \sim \pi_{\boldsymbol{\omega}_k}(\boldsymbol{\theta}|\mathbf{s}^{[i]})$, $i = 1 \dots N$ Sample trajectory $\boldsymbol{\tau}^{[i]} \sim p_{\boldsymbol{\theta}^{[i]}}(\boldsymbol{\tau}|\mathbf{s}^{[i]})$ **Policy Evaluation:**Evaluate policy parameters $R^{[i]} = \sum_{t=0}^T r_t^{[i]}$ Compose data set $\mathcal{D}_{\text{ep}} = \left\{ \boldsymbol{\theta}^{[i]}, \mathbf{s}^{[i]}, R^{[i]} \right\}_{i=1 \dots N}$ **Policy Update:**Compute new policy parameters $\boldsymbol{\omega}_{k+1}$ using \mathcal{D}_{ep}

Algorithms: Episode-based REPS, CRKR, PEPG, RWR

until Policy converges $\boldsymbol{\omega}_{k+1} \approx \boldsymbol{\omega}_k$

that can learn and represent multiple solutions for one task. Such approaches increase the robustness of the learned policy in the case of slightly changing conditions, as we can resort to backup-solutions. Moreover, many policy search approaches have problems with multimodal solution spaces. For example, EM-based or information-theoretic approaches use a weighted average of the samples to determine the new policy. Therefore, these approaches average over several modes, which can considerably decrease the quality of the resulting policy. Such problems can be resolved by using policy updates which are not based on weighted averaging, see Section 2.4.2.4, or by using a mixture model to directly represent several modes in the parameter space [17, 66]. We will discuss such an approach, which is based on episode based REPS in Section 2.4.3.3.

2.4 Policy Update Strategies

In this section, we will describe different policy update strategies used in policy search, such as policy gradient methods, expectation-maximization based methods, information theoretic methods, and policy updates which can be derived from the path-integral theory. In the case where the policy update method has been introduced for both step-

based and episode-based policy evaluation strategies, we will present both resulting algorithms. Policy updates for the step-based evaluation strategy use the data set $\mathcal{D}_{\text{step}}$ to determine the policy update while algorithms based on episode-based policy updates employ the data set \mathcal{D}_{ep} . We will qualitatively compare the algorithms with respect to their sample efficiency, the number of algorithmic parameters that have to be tuned by the user, the type of reward-function that can be employed and also how safe it is to apply the method on a real robot. Methods which are safe to apply on a real robot should not allow big jumps in the policy updates, as such jumps might result in unpredictable behavior which might damage the robot. We will now discuss the different policy update strategies.

Whenever it is possible, we will describe the policy search methods for the episodic reinforcement learning formulation with time-dependent policies as most robot learning tasks are episodic and not infinite horizon tasks. However, most of the derivations also hold for the infinite horizon formulation if we introduce a discount factor γ for the return and cut the trajectories after a horizon of T time steps, where T needs to be sufficiently large such that the influence of future time-steps with $t > T$ vanishes as γ^T approaches zero.

2.4.1 Policy Gradient Methods

Policy gradient (PG) methods use gradient-ascent for maximizing the expected return J_{θ} . In gradient ascent, the parameter update direction is given by the gradient $\nabla_{\theta} J_{\theta}$ as it points in the direction of steepest ascent of the expected return. The policy gradient update is therefore given by

$$\theta_{k+1} = \theta_k + \alpha \nabla_{\theta} J_{\theta},$$

where α is a user-specified learning rate and the policy gradient is given by

$$\nabla_{\theta} J_{\theta} = \int_{\tau} \nabla_{\theta} p_{\theta}(\tau) R(\tau) d\tau. \quad (2.7)$$

We will now discuss different ways to estimate the gradient $\nabla_{\theta} J_{\theta}$.

2.4.1.1 Finite Difference Methods

The finite difference policy gradient [39, 58] is the simplest way of obtaining the policy gradient. It is typically used with the episode-based evaluation strategy. The finite difference method estimates the gradient by applying small perturbations $\delta\boldsymbol{\theta}^{[i]}$ to the parameter vector $\boldsymbol{\theta}_k$. We can either perturb each parameter value separately or use a probability distribution with small variance to create the perturbations. For each perturbation, we obtain the change of the return $\delta R^{[i]} = R(\boldsymbol{\theta}_k + \delta\boldsymbol{\theta}^{[i]}) - R(\boldsymbol{\theta}_k)$. For finite difference methods, the perturbations $\delta\boldsymbol{\theta}^{[i]}$ implement the exploration strategy in parameter space. However, the generation of the perturbations is typically not adapted during learning but predetermined by the user. The gradient $\nabla_{\boldsymbol{\theta}}^{\text{FD}} J_{\boldsymbol{\theta}}$ can be obtained by using a first order Taylor-expansion of $J_{\boldsymbol{\theta}}$ and solving for the gradient in a least-squares sense, i.e., it is determined numerically from the samples as

$$\nabla_{\boldsymbol{\theta}}^{\text{FD}} J_{\boldsymbol{\theta}} = (\delta\boldsymbol{\Theta}^T \delta\boldsymbol{\Theta})^{-1} \delta\boldsymbol{\Theta}^T \delta\mathbf{R}, \quad (2.8)$$

where $\delta\boldsymbol{\Theta} = [\delta\boldsymbol{\theta}^{[1]}, \dots, \delta\boldsymbol{\theta}^{[N]}]^T$ and $\delta\mathbf{R} = [\delta R^{[1]}, \dots, \delta R^{[N]}]$. Finite difference methods are powerful black-box optimizers as long as the evaluations $R(\boldsymbol{\theta})$ are not too noisy. From optimization, this method is also known as Least Square-Based Finite Difference (LSFD) scheme [75].

2.4.1.2 Likelihood-Ratio Policy Gradients

Likelihood-ratio methods were among the first policy search methods introduced in the early 1990s by Williams [89], and include the REINFORCE algorithm. These methods make use of the so called ‘likelihood-ratio’ trick that is given by the identity $\nabla p_{\boldsymbol{\theta}}(\mathbf{y}) = p_{\boldsymbol{\theta}}(\mathbf{y}) \nabla \log p_{\boldsymbol{\theta}}(\mathbf{y})$ ¹. Inserting the likelihood-ratio trick into the policy gradient from Equation (2.7) yields

$$\nabla_{\boldsymbol{\theta}} J_{\boldsymbol{\theta}} = \int p_{\boldsymbol{\theta}}(\boldsymbol{\tau}) \nabla_{\boldsymbol{\theta}} \log p_{\boldsymbol{\theta}}(\boldsymbol{\tau}) R(\boldsymbol{\tau}) d\boldsymbol{\tau} = \mathbb{E}_{p_{\boldsymbol{\theta}}(\boldsymbol{\tau})} [\nabla_{\boldsymbol{\theta}} \log p_{\boldsymbol{\theta}}(\boldsymbol{\tau}) R(\boldsymbol{\tau})], \quad (2.9)$$

¹We can easily confirm this identity by using the chain rule to calculate the derivative of $\log p_{\boldsymbol{\theta}}(\mathbf{y})$, i.e., $\nabla \log p_{\boldsymbol{\theta}}(\mathbf{y}) = \nabla p_{\boldsymbol{\theta}}(\mathbf{y}) / p_{\boldsymbol{\theta}}(\mathbf{y})$.

where the expectation over $p_{\theta}(\boldsymbol{\tau})$ is approximated by using a sum over the sampled trajectories $\boldsymbol{\tau}^{[i]} = (\mathbf{x}_0^{[i]}, \mathbf{u}_0^{[i]}, \mathbf{x}_1^{[i]}, \mathbf{u}_1^{[i]}, \dots)$.

Baselines. As the evaluation $R^{[i]}$ of a parameter $\boldsymbol{\theta}^{[i]}$ or the evaluation $Q_t^{[i]}$ of an action $\mathbf{u}^{[i]}$ is typically performed by inherently noisy Monte-Carlo estimates, the resulting gradient estimates are also afflicted by a large variance. The variance can be reduced by introducing a baseline b for the trajectory reward $R(\boldsymbol{\tau})$, i.e.,

$$\nabla_{\boldsymbol{\theta}} J_{\boldsymbol{\theta}} = \mathbb{E}_{p_{\theta}(\boldsymbol{\tau})} [\nabla_{\boldsymbol{\theta}} \log p_{\theta}(\boldsymbol{\tau})(R(\boldsymbol{\tau}) - b)]. \quad (2.10)$$

Note that the policy gradient estimate remains unbiased as

$$\mathbb{E}_{p_{\theta}(\boldsymbol{\tau})} [\nabla_{\boldsymbol{\theta}} \log p_{\theta}(\boldsymbol{\tau})b] = b \int_{\boldsymbol{\tau}} \nabla_{\boldsymbol{\theta}} p_{\theta}(\boldsymbol{\tau}) d\boldsymbol{\tau} = b \nabla_{\boldsymbol{\theta}} \int_{\boldsymbol{\tau}} p_{\theta}(\boldsymbol{\tau}) d\boldsymbol{\tau} = 0, \quad (2.11)$$

where we first applied the reverse of the ‘likelihood-ratio’ trick and subsequently the identity $\int_{\boldsymbol{\tau}} p_{\theta}(\boldsymbol{\tau}) d\boldsymbol{\tau} = 1$. Since the baseline b is a free parameter, we can choose it such that it minimizes the variance of the gradient estimate. We will denote the variance-minimizing baseline as the optimal baseline. As the likelihood gradient can be estimated in different ways, the corresponding optimal baseline will change with it. We now first discuss the step-based likelihood-ratio PG algorithms, and, discuss their optimal baselines if it is given in the literature. Subsequently, we will cover the episode-based likelihood-ratio variant.

Step-Based Likelihood-Ratio Methods

Step-based algorithms exploit the structure of the trajectory distribution, i.e.,

$$p_{\theta}(\boldsymbol{\tau}) = p(\mathbf{x}_1) \prod_{t=1}^T p(\mathbf{x}_{t+1} | \mathbf{x}_t, \mathbf{u}_t) \pi_{\theta}(\mathbf{u}_t | \mathbf{x}_t, t)$$

to decompose $\nabla_{\boldsymbol{\theta}} \log p_{\theta}(\boldsymbol{\tau})$ into the single time steps. As the product is transformed into a sum by a logarithm, all terms which do not depend on the policy parameters $\boldsymbol{\theta}$ disappear during differentiation. Hence,

$\nabla_{\theta} \log p_{\theta}(\tau)$ is given by

$$\nabla_{\theta} \log p_{\theta}(\tau) = \sum_{t=0}^{T-1} \nabla_{\theta} \log \pi_{\theta}(\mathbf{u}_t | \mathbf{x}_t, t). \quad (2.12)$$

Equation (2.12) reveals a key result for policy gradients: $\nabla_{\theta} \log p_{\theta}(\tau)$ does not depend on the transition model $p(\mathbf{x}_{t+1} | \mathbf{x}_t, \mathbf{u}_t)$. Note that this result holds for any stochastic policy. However, for deterministic policies $\pi_{\theta}(\mathbf{x}_t, t)$, the gradient $\nabla_{\theta} \log p_{\theta}(\tau)$ includes

$$\nabla_{\theta} p(\mathbf{x}_{t+1} | \mathbf{x}_t, \pi_{\theta}(\mathbf{x}_t, t)) = \frac{\partial p(\mathbf{x}_{t+1} | \mathbf{x}_t, \mathbf{u}_t)}{\partial \mathbf{u}_t} \frac{\partial \mathbf{u}_t}{\partial \theta} \Big|_{\mathbf{u}_t = \pi_{\theta}(\mathbf{x}_t, t)},$$

and, hence, the transition model needs to be known. Consequently, stochastic policies play a crucial role for policy gradient methods.

The REINFORCE Algorithm. Equation (2.12) is used by one of the first PG algorithms introduced in the machine learning literature, the REINFORCE algorithm [89]. The REINFORCE policy gradient is given by

$$\nabla_{\theta}^{\text{RF}} J_{\theta} = \mathbb{E}_{p_{\theta}(\tau)} \left[\sum_{t=0}^{T-1} \nabla_{\theta} \log \pi_{\theta}(\mathbf{u}_t | \mathbf{x}_t, t) (R(\tau) - b) \right], \quad (2.13)$$

where b denotes the baseline.

To minimize the variance of $\nabla_{\theta}^{\text{RF}} J_{\theta}$, we estimate the optimal baseline b^{RF} . The optimal baseline also depends on which element h of the gradient $\nabla_{\theta} J_{\theta}$ we want to evaluate, and, hence needs to be computed for each dimension h separately. The optimal baseline b_h^{RF} for the REINFORCE algorithm minimizes the variance of $\nabla_{\theta_h}^{\text{RF}} J_{\theta}$, i.e., it satisfies the condition

$$\begin{aligned} \frac{\partial}{\partial b} \text{Var} [\nabla_{\theta_h}^{\text{RF}} J_{\theta}] &= \frac{\partial}{\partial b} \left(\mathbb{E}_{p_{\theta}(\tau)} [(\nabla_{\theta_h}^{\text{RF}} J_{\theta})^2] - \mathbb{E}_{p_{\theta}(\tau)} [\nabla_{\theta_h}^{\text{RF}} J_{\theta}]^2 \right) \\ &= \frac{\partial}{\partial b} \mathbb{E}_{p_{\theta}(\tau)} [(\nabla_{\theta_h}^{\text{RF}} J_{\theta})^2] = 0, \end{aligned} \quad (2.14)$$

where the second term disappeared as the expected gradient is not affected by the baseline, see Equations (2.10) and (2.11). Solving this

Algorithm 5 REINFORCE

Input: policy parametrization θ ,

$$\text{data-set } \mathcal{D} = \left\{ \mathbf{x}_{1:T}^{[i]}, \mathbf{u}_{1:T-1}^{[i]}, r_{1:T}^{[i]} \right\}_{i=1 \dots N}$$

Compute returns: $R^{[i]} = \sum_{t=0}^T r_t^{[i]}$ **for** each dimension h of θ **do**

Estimate optimal baseline:

$$b_h^{\text{RF}} = \frac{\sum_{i=1}^N \left(\sum_{t=0}^{T-1} \nabla_{\theta_h} \log \pi_{\theta} \left(\mathbf{u}_t^{[i]} \mid \mathbf{x}_t^{[i]}, t \right) \right)^2 R^{[i]}}{\sum_{i=1}^N \left(\sum_{t=0}^{T-1} \nabla_{\theta_h} \log \pi_{\theta} \left(\mathbf{u}_t^{[i]} \mid \mathbf{x}_t^{[i]}, t \right) \right)^2}$$

 Estimate derivative for dimension h of θ :

$$\nabla_{\theta_h}^{\text{RF}} J_{\theta} = \frac{1}{N} \sum_{i=1}^N \sum_{t=0}^{T-1} \nabla_{\theta_h} \log \pi_{\theta} \left(\mathbf{u}_t^{[i]} \mid \mathbf{x}_t^{[i]}, t \right) (R^{[i]} - b_h^{\text{RF}})$$

end forReturn $\nabla_{\theta}^{\text{RF}} J_{\theta}$

equation for b yields

$$b_h^{\text{RF}} = \frac{\mathbb{E}_{p_{\theta}(\tau)} \left[\left(\sum_{t=0}^{T-1} \nabla_{\theta_h} \log \pi_{\theta} (\mathbf{u}_t \mid \mathbf{x}_t, t) \right)^2 R(\tau) \right]}{\mathbb{E}_{p_{\theta}(\tau)} \left[\left(\sum_{t=0}^{T-1} \nabla_{\theta_h} \log \pi_{\theta} (\mathbf{u}_t \mid \mathbf{x}_t, t) \right)^2 \right]}. \quad (2.15)$$

The REINFORCE algorithm with its optimal baseline is summarized in Algorithm 5.

The G(PO)MDP Algorithm. From Equation (2.13), we realize that REINFORCE uses the returns $R(\tau)$ of the whole episode as the evaluations of single actions despite using the step-based policy evaluation strategy. As already discussed before, the variance of the returns can grow with the trajectory length, and, hence, deteriorate the performance of the algorithm even if used with the optimal baseline. However, by decomposing the return in the rewards of the single time steps, we can use the observation that rewards from the past do not depend on

actions in the future, and, hence, $\mathbb{E}_{p_{\theta}(\tau)} [\partial_{\theta} \log \pi_{\theta}(\mathbf{u}_t | \mathbf{x}_t, t) r_j] = 0$ for $j < t^2$. If we look at a reward r_j of a single time-step j , we realize that we can neglect all derivatives of future actions. This intuition has been used for the G(PO)MDP algorithm [9, 10] to decrease the variance of policy gradient estimates. The policy gradient of G(PO)MDP is given by

$$\nabla_{\theta}^{\text{GMDP}} J_{\theta} = \mathbb{E}_{p_{\theta}(\tau)} \left[\sum_{j=0}^{T-1} \sum_{t=0}^j \nabla_{\theta} \log \pi_{\theta}(\mathbf{u}_t | \mathbf{x}_t, t) (r_j - b_j) \right], \quad (2.16)$$

where b_j is a time-dependent baseline. The optimal baseline for the G(PO)MDP algorithm $b_{h,j}^{\text{GMDP}}(\mathbf{x})$ for time step j and dimension h of θ can be obtained similarly as for the REINFORCE algorithm and is given by

$$b_{h,j}^{\text{GMDP}} = \frac{\mathbb{E}_{p_{\theta}(\tau)} \left[\left(\sum_{t=0}^j \nabla_{\theta_h} \log \pi_{\theta}(\mathbf{u}_t | \mathbf{x}_t, t) \right)^2 r_j \right]}{\mathbb{E}_{p_{\theta}(\tau)} \left[\left(\sum_{t=0}^j \nabla_{\theta_h} \log \pi_{\theta}(\mathbf{u}_t | \mathbf{x}_t, t) \right)^2 \right]}. \quad (2.17)$$

The G(PO)MDP algorithm is summarized in Algorithm 6.

The Policy Gradient Theorem Algorithm. Instead of using the returns $R(\tau)$ we can also use the expected reward to come at time step t , i.e., $Q_t^{\pi}(\mathbf{x}_t, \mathbf{u}_t)$, to evaluate an action \mathbf{u}_t . Mathematically, such an evaluation can be justified by the same observation that has been used for the G(PO)MDP algorithm, i.e., that rewards are not correlated with future actions. Such evaluation is used by the Policy Gradient Theorem (PGT) algorithm [79], which states that

$$\begin{aligned} \nabla_{\theta}^{\text{PG}} J_{\theta} &= \mathbb{E}_{p_{\theta}(\tau)} \left[\sum_{t=0}^{T-1} \nabla_{\theta} \log \pi_{\theta}(\mathbf{u}_t | \mathbf{x}_t) \left(\sum_{j=t}^T r_j \right) \right] \\ &= \mathbb{E}_{p_{\theta}(\tau)} \left[\sum_{t=0}^{T-1} \nabla_{\theta} \log \pi_{\theta}(\mathbf{u}_t | \mathbf{x}_t) Q_t^{\pi}(\mathbf{x}_t, \mathbf{u}_t) \right]. \end{aligned} \quad (2.18)$$

We can again subtract an arbitrary baseline $b_t(\mathbf{x})$ from $Q_t^{\pi}(\mathbf{x}_t, \mathbf{u}_t)$, which now depends on the state \mathbf{x} as well as on the time step.

²We can follow the same argument as in Equation (2.11) for introducing the baseline to prove this identity.

Algorithm 6 G(PO)MDP Algorithm

 Input: Policy parametrization θ ,

 Data set $\mathcal{D} = \left\{ \mathbf{x}_{1:T}^{[i]}, \mathbf{u}_{1:T-1}^{[i]}, r_{1:T-1}^{[i]} \right\}_{i=1 \dots N}$
for each time step $t = 0 \dots T - 1$ **do**
for each dimension h of θ **do**

Estimate optimal time-dependent baseline:

$$b_{h,j}^{\text{GMDP}} = \frac{\sum_{i=1}^N \left(\sum_{t=0}^j \nabla_{\theta_h} \log \pi_{\theta} \left(\mathbf{u}_h^{[i]} \mid \mathbf{x}_h^{[i]}, h \right) \right)^2 r_j^{[i]}}{\sum_{i=1}^N \sum_{t=0}^{T-1} \left(\sum_{t=0}^j \nabla_{\theta_h} \log \pi_{\theta} \left(\mathbf{u}_k^{[i]} \mid \mathbf{x}_k^{[i]}, k \right) \right)^2}$$

end for

 Estimate gradient for dimension h :

$$\nabla_{\theta_h}^{\text{GMDP}} J_{\theta} = \sum_{i=1}^N \sum_{j=0}^{T-1} \left(\sum_{t=0}^j \nabla_{\theta_h} \log \pi_{\theta} \left(\mathbf{u}_t^{[i]} \mid \mathbf{x}_t^{[i]}, t \right) \right) \left(r_j^{[i]} - b_{h,j}^{\text{GMDP}} \right)$$

end for

 Return $\nabla_{\theta}^{\text{GMDP}} J_{\theta}$

While $Q_t^{\pi}(\mathbf{x}_t, \mathbf{u}_t)$ can be estimated by Monte-Carlo roll-outs, the PGT algorithm can be used in combination with function approximation as will be covered in Section 2.4.1.3.

Episode-Based Likelihood-Ratio Methods

Episode-based likelihood-ratio methods directly update the upper-level policy $\pi_{\omega}(\theta)$ for choosing the parameters θ of the lower-level policy $\pi_{\theta}(\mathbf{u}_t \mid \mathbf{x}_t, t)$. They optimize the expected return J_{ω} , as defined in Equation (2.2). The likelihood gradient of J_{ω} can be directly obtained by replacing $p_{\theta}(\tau)$ with $\pi_{\omega}(\theta)$ and $R(\tau)$ with $R(\theta) = \int_{\tau} p_{\theta}(\tau) R(\tau) d\tau$ in Equation (2.9), resulting in

$$\nabla_{\omega}^{\text{PE}} J_{\omega} = \mathbb{E}_{\pi_{\omega}(\theta)} [\nabla_{\omega} \log \pi_{\omega}(\theta) (R(\theta) - b)]. \quad (2.19)$$

Such an approach was first introduced by [73, 74] with the Parameter Exploring Policy Gradient (PEPG) algorithm. The optimal base-line

Algorithm 7 Parameter Exploring Policy Gradient Algorithm

Input: Policy parametrization ω Data set $\mathcal{D} = \left\{ \boldsymbol{\theta}^{[i]}, R^{[i]} \right\}_{i=1 \dots N}$ **for** each dimension h of ω **do**

Estimate optimal baseline:

$$b_h^{\text{PGPE}} = \frac{\sum_{i=1}^N \left(\nabla_{\omega_h} \pi_{\omega}(\boldsymbol{\theta}^{[i]}) \right)^2 R^{[i]}}{\sum_{i=1}^N \left(\nabla_{\omega_h} \pi_{\omega}(\boldsymbol{\theta}^{[i]}) \right)^2}$$

Estimate derivative for dimension h of ω :

$$\nabla_{\omega_h}^{\text{PE}} J_{\omega} = \frac{1}{N} \sum_{i=1}^N \nabla_{\omega} \pi_{\omega_h}(\boldsymbol{\theta}^{[i]}) (R^{[i]} - b_h^{\text{PGPE}})$$

end for

b_h^{PGPE} for the h -th element of the PEPG gradient is obtained similarly as for the REINFORCE algorithm and given by

$$b_h^{\text{PGPE}} = \frac{\mathbb{E}_{\pi_{\omega}(\boldsymbol{\theta})} \left[\left(\nabla_{\omega_h} \pi_{\omega}(\boldsymbol{\theta}) \right)^2 R(\boldsymbol{\theta}) \right]}{\mathbb{E}_{p_{\boldsymbol{\theta}}(\boldsymbol{\tau})} \left[\left(\nabla_{\omega_h} \pi_{\omega}(\boldsymbol{\theta}) \right)^2 \right]}. \quad (2.20)$$

The PEPG algorithm is summarized in Algorithm 7.

2.4.1.3 Natural Gradients

The natural gradient [3] is a well known concept from supervised learning for optimizing parametrized probability distributions $p_{\boldsymbol{\theta}}(\mathbf{y})$, where \mathbf{y} is a random variable, which often achieves faster convergence than the traditional gradient. Traditional gradient methods typically use an Euclidean metric $\delta \boldsymbol{\theta}^T \delta \boldsymbol{\theta}$ to determine the direction of the update $\delta \boldsymbol{\theta}$, i.e., they assume that all parameter dimensions have similarly strong effects on the resulting distribution. However, already small changes in $\boldsymbol{\theta}$ might result in large changes of the resulting distribution $p_{\boldsymbol{\theta}}(\mathbf{y})$. As the gradient estimation typically depends on $p_{\boldsymbol{\theta}}(\mathbf{y})$ due to the sampling process, the next gradient estimate can also change dramatically.

To achieve a stable behavior of the learning process, it is desirable to enforce that the distribution $p_{\boldsymbol{\theta}}(\mathbf{y})$ does not change too much in one update step. This intuition is the key concept behind the natural gradient which limits the distance between the distribution $p_{\boldsymbol{\theta}}(\mathbf{y})$ before and $p_{\boldsymbol{\theta}+\delta\boldsymbol{\theta}}(\mathbf{y})$ after the update. To measure the distance between $p_{\boldsymbol{\theta}}(\mathbf{y})$ and $p_{\boldsymbol{\theta}+\delta\boldsymbol{\theta}}(\mathbf{y})$, the natural gradient uses an approximation of the Kullback-Leibler (KL) divergence. The KL-divergence is a similarity measure of two distributions. It has been shown that the Fisher information matrix

$$\mathbf{F}_{\boldsymbol{\theta}} = \mathbb{E}_{p(\mathbf{y})} [\nabla_{\boldsymbol{\theta}} \log p(\mathbf{y}) \nabla_{\boldsymbol{\theta}} \log p(\mathbf{y})^T] \quad (2.21)$$

can be used to approximate the KL divergence between $p_{\boldsymbol{\theta}+\delta\boldsymbol{\theta}}(\mathbf{y})$ and $p_{\boldsymbol{\theta}}(\mathbf{y})$ for sufficiently small $\delta\boldsymbol{\theta}$, i.e.,

$$\text{KL}(p_{\boldsymbol{\theta}+\delta\boldsymbol{\theta}}(\mathbf{y})||p_{\boldsymbol{\theta}}(\mathbf{y})) \approx \delta\boldsymbol{\theta}^T \mathbf{F}_{\boldsymbol{\theta}} \delta\boldsymbol{\theta}. \quad (2.22)$$

The natural gradient update $\delta\boldsymbol{\theta}^{\text{NG}}$ is now defined as the update $\delta\boldsymbol{\theta}$ that is the most similar to the traditional ‘vanilla’ gradient $\delta\boldsymbol{\theta}^{\text{VG}}$ update that has a bounded distance

$$\text{KL}(p_{\boldsymbol{\theta}+\delta\boldsymbol{\theta}}(\mathbf{y})||p_{\boldsymbol{\theta}}(\mathbf{y})) \leq \epsilon$$

in the distribution space. Hence, we can formulate the following optimization program

$$\delta\boldsymbol{\theta}^{\text{NG}} = \operatorname{argmax}_{\delta\boldsymbol{\theta}} \delta\boldsymbol{\theta}^T \delta\boldsymbol{\theta}^{\text{VG}} \quad \text{s.t.} \quad \delta\boldsymbol{\theta}^T \mathbf{F}_{\boldsymbol{\theta}} \delta\boldsymbol{\theta} \leq \epsilon. \quad (2.23)$$

The solution of this program is given by $\delta\boldsymbol{\theta}^{\text{NG}} \propto \mathbf{F}_{\boldsymbol{\theta}}^{-1} \delta\boldsymbol{\theta}^{\text{VG}}$ up to a scaling factor. The proportionality factor for the update step is typically subsumed into the learning rate. The natural gradient linearly transforms the traditional gradient by the inverse Fisher matrix, which renders the parameter update also invariant to linear transformations of the parameters of the distribution, i.e., if two parametrizations have the same representative power, the natural gradient update will be identical. As the Fisher information matrix is always positive definite, the natural gradient always rotates the traditional gradient by less than 90 degrees, and, hence, all convergence guarantees from standard gradient-based optimization remain. In contrast to the traditional gradient, the natural gradient avoids premature convergence on plateaus and overaggressive steps on steep ridges due to its isotropic convergence properties [3, 78].

Natural Policy Gradients. The intuition of the natural gradients to limit the distance between two subsequent distributions is also useful for policy search. Here, we want to maintain a limited step-width in the trajectory distribution space, i.e.,

$$\text{KL}(p_{\boldsymbol{\theta}}(\boldsymbol{\tau})||p_{\boldsymbol{\theta}+\delta\boldsymbol{\theta}}(\boldsymbol{\tau})) \approx \delta\boldsymbol{\theta}^T \mathbf{F}_{\boldsymbol{\theta}} \delta\boldsymbol{\theta} \leq \epsilon.$$

The Fisher information matrix

$$\mathbf{F}_{\boldsymbol{\theta}} = \mathbb{E}_{p_{\boldsymbol{\theta}}(\boldsymbol{\tau})} \left[\nabla_{\boldsymbol{\theta}} \log p_{\boldsymbol{\theta}}(\boldsymbol{\tau}) \nabla_{\boldsymbol{\theta}} \log p_{\boldsymbol{\theta}}(\boldsymbol{\tau})^T \right]$$

is now computed for the trajectory distribution $p_{\boldsymbol{\theta}}(\boldsymbol{\tau})$. The natural policy gradient $\nabla_{\boldsymbol{\theta}}^{\text{NG}} J_{\boldsymbol{\theta}}$ is therefore given by

$$\nabla_{\boldsymbol{\theta}}^{\text{NG}} J_{\boldsymbol{\theta}} = \mathbf{F}_{\boldsymbol{\theta}}^{-1} \nabla_{\boldsymbol{\theta}} J_{\boldsymbol{\theta}}, \quad (2.24)$$

where $\nabla_{\boldsymbol{\theta}} J_{\boldsymbol{\theta}}$ can be estimated by any traditional policy gradient method.

The difference between the natural and traditional policy gradient for learning a simple linear feedback policy is shown in Figure 2.2. In this example, a scalar controller gain and the variance of the policy are optimized. While the traditional gradient quickly reduces the variance of the policy, and, hence will stop exploring, the natural gradient only gradually decreases the variance, and, in the end, finds the optimal solution faster.

Step-Based Natural Gradient Methods

Similar to the gradient, the Fisher information matrix can also be decomposed in the policy derivatives of the single time steps [8]. In [62, 61], it was shown that the Fisher information matrix of the trajectory distribution can be written as the average Fisher information matrices for each time step, i.e.,

$$\mathbf{F}_{\boldsymbol{\theta}} = \mathbb{E}_{p_{\boldsymbol{\theta}}(\boldsymbol{\tau})} \left[\sum_{t=0}^{T-1} \nabla_{\boldsymbol{\theta}} \log \pi_{\boldsymbol{\theta}}(\mathbf{u}_t | \mathbf{x}_t, t) \nabla_{\boldsymbol{\theta}} \log \pi_{\boldsymbol{\theta}}(\mathbf{u}_t | \mathbf{x}_t, t)^T \right]. \quad (2.25)$$

Consequently, as it is the case for estimating the policy gradient, the transition model is not needed for estimating $\mathbf{F}_{\boldsymbol{\theta}}$. Still, estimating $\mathbf{F}_{\boldsymbol{\theta}}$

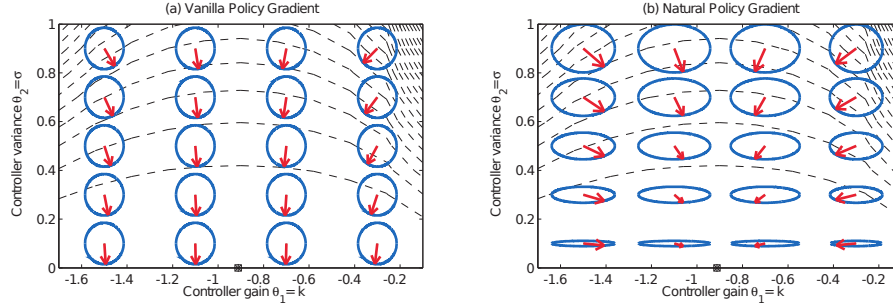


Fig. 2.2 Comparison of the natural gradient to the traditional gradient on a simple linear task with quadratic cost function. The controller has two parameters, the feedback gain k and the variance σ^2 . The main difference between the two methods is how the change in parameters is punished, i.e., the distance between current and next policy parameters. This distance is indicated by the blue ellipses in the contour plot while the dashed lines show the expected return. The red arrows indicate the resulting gradient. While the traditional gradient quickly reduces the variance of the policy, the natural gradient only gradually decreases the variance, and therefore continues to explore.

from samples can be difficult, as \mathbf{F}_θ contains $O(d^2)$ parameters, where d is the dimensionality of θ . However, the Fisher information matrix \mathbf{F}_θ does not need to be estimated explicitly if we use *compatible function approximations*, which we will introduce in the next paragraph.

Compatible Function Approximation. In the PGT algorithm, given in Section 2.4.1.2, the policy gradient was given by

$$\nabla_{\theta}^{\text{PG}} J_{\theta} = \mathbb{E}_{p_{\theta}(\tau)} \left[\sum_{t=0}^{T-1} \nabla_{\theta} \log \pi_{\theta}(\mathbf{u}_t | \mathbf{x}_t) (Q_t^{\pi}(\mathbf{x}_t, \mathbf{u}_t) - b_t(\mathbf{x}_t)) \right]. \quad (2.26)$$

Instead of using the future rewards of a single roll-out to estimate $Q_t^{\pi}(\mathbf{x}_t, \mathbf{u}_t) - b_t(\mathbf{x})$, we can also use function approximation [79] to estimate the value, i.e., estimate a function $\tilde{A}_{\mathbf{w}}(\mathbf{x}_t, \mathbf{u}_t, t) = \psi_t(\mathbf{x}_t, \mathbf{u}_t)^T \mathbf{w}$ such that $\tilde{A}_{\mathbf{w}}(\mathbf{x}_t, \mathbf{u}_t, t) \approx Q_t^{\pi}(\mathbf{x}_t, \mathbf{u}_t) - b_t(\mathbf{x})$. The quality of the approximation is determined by the choice of the basis functions $\psi_t(\mathbf{x}_t, \mathbf{u}_t)$, which might explicitly depend on the time step t . A good function approximation does not change the gradient in expectation, i.e., it does not introduce a bias. To find basis functions $\psi_t(\mathbf{x}_t, \mathbf{u}_t)$ that fulfill this condition, we will first assume that we already found a parameter vec-

tor \mathbf{w} which approximates Q_t^π . For simplicity, we will for now assume that the baseline $b_t(\mathbf{x})$ is zero. A parameter vector \mathbf{w} , which approximates Q_t^π , also minimizes the squared approximation error. Thus, \mathbf{w} has to satisfy

$$\frac{\partial}{\partial \mathbf{w}} \mathbb{E}_{p_\theta(\tau)} \left[\sum_{t=0}^{T-1} \left(Q_t(\mathbf{x}_t, \mathbf{u}_t) - \tilde{A}_\mathbf{w}(\mathbf{x}_t, \mathbf{u}_t, t) \right)^2 \right] = 0. \quad (2.27)$$

Computing the derivative yields

$$2\mathbb{E}_{p_\theta(\tau)} \left[\sum_{t=0}^{T-1} (Q_t(\mathbf{x}_t, \mathbf{u}_t) - \tilde{A}_\mathbf{w}(\mathbf{x}_t, \mathbf{u}_t, t)) \frac{\partial}{\partial \mathbf{w}} \tilde{A}_\mathbf{w}(\mathbf{x}_t, \mathbf{u}_t, t) \right] = 0 \quad (2.28)$$

with $\partial \tilde{A}_\mathbf{w}(\mathbf{x}_t, \mathbf{u}_t, t) / \partial \mathbf{w} = \psi_t(\mathbf{x}_t, \mathbf{u}_t)$. By subtracting this equation from the Policy Gradient Theorem in Equation (2.18), it is easy to see that

$$\nabla_{\text{PG}} J(\boldsymbol{\theta}) = \mathbb{E}_{p_\theta(\tau)} \left[\sum_{t=1}^T \nabla_{\boldsymbol{\theta}} \log \pi_{\boldsymbol{\theta}}(\mathbf{u}_t | \mathbf{x}_t, t) \tilde{A}_\mathbf{w}(\mathbf{x}_t, \mathbf{u}_t, t) \right] \quad (2.29)$$

if we use $\nabla_{\boldsymbol{\theta}} \log \pi_{\boldsymbol{\theta}}(\mathbf{u}_t | \mathbf{x}_t, t)$ as basis functions $\psi_t(\mathbf{x}_t, \mathbf{u}_t)$ for $\tilde{A}_\mathbf{w}$.

Using $\psi_t(\mathbf{x}_t, \mathbf{u}_t) = \nabla_{\boldsymbol{\theta}} \log \pi_{\boldsymbol{\theta}}(\mathbf{u}_t | \mathbf{x}_t, t)$ as basis functions is also called *compatible function approximation* [79], as the function approximation is compatible with the policy parametrization. The policy gradient using compatible function approximation can now be written as

$$\nabla_{\boldsymbol{\theta}}^{\text{FA}} J_{\boldsymbol{\theta}} = \mathbb{E}_{p_\theta(\tau)} \left[\sum_{t=0}^{T-1} \nabla_{\boldsymbol{\theta}} \log \pi_{\boldsymbol{\theta}}(\mathbf{u}_t | \mathbf{x}_t, t) \log \pi_{\boldsymbol{\theta}}(\mathbf{u}_t | \mathbf{x}_t, t)^T \right] \mathbf{w} = \mathbf{G}_{\boldsymbol{\theta}} \mathbf{w}. \quad (2.30)$$

Hence, in order to compute the traditional gradient, we have to estimate the weight parameters \mathbf{w} of the advantage function and the matrix $\mathbf{G}_{\boldsymbol{\theta}}$. However, as we will see in the next section, the matrix $\mathbf{G}_{\boldsymbol{\theta}}$ cancels out for the natural gradient, and, hence, computing the natural gradient reduces to computing the weights \mathbf{w} for the compatible function approximation.

Step-Based Natural Policy Gradient. The result given in Equation (2.30) implies that the policy gradient $\nabla_{\boldsymbol{\theta}}^{\text{FA}} J_{\boldsymbol{\theta}}$ using the compatible

function approximation already contains the Fisher information matrix as $\mathbf{G}_\theta = \mathbf{F}_\theta$. Hence, the calculation of the step-based natural gradient simplifies to

$$\nabla_\theta^{\text{NG}} J_\theta = \mathbf{F}_\theta^{-1} \nabla_\theta^{\text{FA}} J_\theta = \mathbf{w}. \quad (2.31)$$

The natural gradient still requires estimating the function \tilde{A}_w . Due to the baseline $b_t(\mathbf{x})$, the function $\tilde{A}_w(\mathbf{x}_t, \mathbf{u}_t, t)$ can be interpreted as the advantage function, i.e., $\tilde{A}_w(\mathbf{x}_t, \mathbf{u}_t, t) \approx Q_t^\pi(\mathbf{x}_t, \mathbf{u}_t) - V_t(\mathbf{x}_t)$. We can check that \tilde{A}_w is an advantage function by observing that $\mathbb{E}_{p_\theta(\tau)} [\tilde{A}_w(\mathbf{x}_t, \mathbf{u}_t, t)] = \mathbb{E}_{p_\theta(\tau)} [\nabla_\theta \log \pi_\theta(\mathbf{u}_t | \mathbf{x}_t, t)] \mathbf{w} = 0$. The advantage function \tilde{A}_w can be estimated by using temporal difference methods [80, 13]. However, in order to estimate the advantage function, such methods also require an estimate of the value function $V_t(\mathbf{x}_t)$ [61]. While the advantage function would be easy to learn as its basis functions are given by the compatible function approximation, appropriate basis functions for the value function are typically more difficult to specify. Hence, we typically want to find algorithms which avoid estimating a value function.

Episodic Natural Actor Critic. One such algorithm is the episodic Natural Actor Critic (eNAC) algorithm [60]. In the episodic policy search setup, i.e., with a limited time-horizon T , the estimation of the value function V_t can be avoided by considering whole sample paths. To see this, we first rewrite the Bellman equation for the advantage function

$$\tilde{A}_w(\mathbf{x}, \mathbf{u}, t) + V_t(\mathbf{x}) = r_t(\mathbf{x}, \mathbf{u}) + \int p(\mathbf{x}' | \mathbf{x}, \mathbf{u}) V_t(\mathbf{x}') d\mathbf{x}', \quad (2.32)$$

where $V_T(\mathbf{x}) = r_T(\mathbf{x})$ is the reward for the final state. Equation (2.32) can be rewritten as

$$\tilde{A}_w(\mathbf{x}_t, \mathbf{u}_t, t) + V_t(\mathbf{x}_t) = r_t(\mathbf{x}_t, \mathbf{u}_t) + V_t(\mathbf{x}_{t+1}) + \epsilon, \quad (2.33)$$

for a single transition from \mathbf{x}_t to \mathbf{x}_{t+1} , where ϵ is a zero-mean noise term. We now sum up Equation (2.33) along a sample path and get the following condition

$$\sum_{t=0}^{T-1} \nabla_\theta \log \pi_\theta(\mathbf{u}_t | \mathbf{x}_t, t) \mathbf{w} + V_0(\mathbf{x}_0) = \sum_{t=0}^{T-1} r_t(\mathbf{x}_t, \mathbf{u}_t) + r_T(\mathbf{x}_T) \quad (2.34)$$

Algorithm 8 Episodic Natural Actor Critic

Input: Policy parametrization θ ,
 data-set $\mathcal{D} = \left\{ \mathbf{x}_{1:T}^{[i]}, \mathbf{u}_{1:T-1}^{[i]}, r_{1:T}^{[i]} \right\}_{i=1 \dots N}$
for each sample $i = 1 \dots N$ **do**
 Compute returns: $R^{[i]} = \sum_{t=0}^T r_t^{[i]}$
 Compute features: $\psi^{[i]} = \begin{bmatrix} \sum_{t=0}^{T-1} \nabla_{\theta} \log \pi_{\theta} \left(\mathbf{u}_t^{[i]} | \mathbf{x}_t^{[i]}, t \right) \\ \varphi(\mathbf{x}_0^{[i]}) \end{bmatrix}$
end for
 Fit advantage function and initial value function

$$\mathbf{R} = \left[R^{[1]}, \dots, R^{[N]} \right]^T, \quad \mathbf{\Psi} = \left[\psi^{[1]}, \dots, \psi^{[N]} \right]^T$$

$$\begin{bmatrix} \mathbf{w} \\ \mathbf{v} \end{bmatrix} = (\mathbf{\Psi}^T \mathbf{\Psi})^{-1} \mathbf{\Psi}^T \mathbf{R}$$

return $\nabla_{\theta}^{\text{eNAC}} J_{\theta} = \mathbf{w}$

with $\nabla_{\theta} \log \pi_{\theta}(\mathbf{u}_t | \mathbf{x}_t, t) \mathbf{w} = \tilde{A}_{\mathbf{w}}(\mathbf{x}_t, \mathbf{u}_t, t)$. Now, the value function V_t needs to be estimated only for the first time step. For a single start state \mathbf{x}_0 , estimating $V_0(\mathbf{x}_0) = v_0$ reduces to estimating a constant v_0 . For multiple start states \mathbf{x}_0 , V_t needs to be parametrized $V_0(\mathbf{x}_0) \approx \tilde{V}_{\mathbf{v}}(\mathbf{x}_0) = \varphi(\mathbf{x})^T \mathbf{v}$. By using multiple sample paths $\tau^{[i]}$, we get a regression problem whose solution is given by

$$\begin{bmatrix} \mathbf{w} \\ \mathbf{v} \end{bmatrix} = (\mathbf{\Psi}^T \mathbf{\Psi})^{-1} \mathbf{\Psi}^T \mathbf{R}, \quad (2.35)$$

where the matrix $\mathbf{\Psi}$ contains the policy and value function features of the sample paths, i.e.,

$$\psi^{[i]} = \begin{bmatrix} \sum_{t=0}^{T-1} \nabla_{\theta} \log \pi_{\theta} \left(\mathbf{u}_t^{[i]} | \mathbf{x}_t^{[i]}, t \right) \\ \varphi(\mathbf{x}_0^{[i]})^T \end{bmatrix}, \quad \mathbf{\Psi} = \left[\psi^{[1]}, \dots, \psi^{[N]} \right]^T$$

and \mathbf{R} contains the returns of the sample paths. The eNAC algorithm is illustrated in Algorithm 8.

Natural Actor Critic. While the eNAC algorithm is efficient for the episodic reinforcement learning formulation, it uses the returns

$R^{[i]}$ for evaluating the policy, and consequently, gets less accurate for large time-horizons due to the large variance of the returns. For learning problems with a large time horizon, especially, for infinite horizon tasks, the convergence speed can be improved by directly estimating the value function V_t . Such a strategy is implemented by the Natural Actor Critic Algorithm (NAC) algorithm [61]. The NAC algorithm estimates the advantage function and the value function by applying temporal difference methods [80, 13]. To do so, temporal difference methods have to first be adapted to learn the advantage function.

We start this derivation by first writing down the Bellman equation in terms of the advantage function in the infinite horizon formulation

$$Q^\pi(\mathbf{x}, \mathbf{u}) = A^\pi(\mathbf{x}, \mathbf{u}) + V^\pi(\mathbf{x}) = r(\mathbf{x}, \mathbf{u}) + \gamma \int p(\mathbf{x}'|\mathbf{x}, \mathbf{u})V^\pi(\mathbf{x}')d\mathbf{x}. \quad (2.36)$$

Note that we now discuss the infinite horizon case, i.e., all functions are time independent and we introduced the discount factor γ . By inserting $\tilde{A}(\mathbf{x}, \mathbf{u}) \approx \nabla_{\boldsymbol{\theta}} \log \pi_{\boldsymbol{\theta}}(\mathbf{u}|\mathbf{x})\mathbf{w}$ and $V^\pi(\mathbf{x}) \approx \boldsymbol{\varphi}(\mathbf{x})^T \mathbf{v}$, we can rewrite the Bellman equation as a set of linear equations

$$\nabla_{\boldsymbol{\theta}} \log \pi_{\boldsymbol{\theta}} \left(\mathbf{u}^{[i]} \middle| \mathbf{x}^{[i]} \right) \mathbf{w} + \boldsymbol{\varphi}(\mathbf{x}^{[i]})^T \mathbf{v} = r(\mathbf{x}^{[i]}, \mathbf{u}^{[i]}) + \gamma \boldsymbol{\varphi}(\mathbf{x}'^{[i]})^T \mathbf{v} + \epsilon. \quad (2.37)$$

One efficient method to estimate \mathbf{w} and \mathbf{v} is to use the LSTD-Q(λ) [13, 61] algorithm. A simplified version of the NAC algorithm which uses LSTD-Q(0) to estimate \mathbf{w} and \mathbf{v} is given in Algorithm 9. For a more detailed discussion of the LSTD-Q algorithm we refer to the corresponding papers [13, 61, 42].

Episode-Based Natural Policy Gradients

The beneficial properties of the natural gradient can also be exploited for episode-based algorithms [88, 78]. While such methods come from the area of evolutionary algorithms, as they always maintain a ‘population’ of parameter-samples, they perform gradient ascent on a fitness function which is in the reinforcement learning context the expected long-term reward $J_{\boldsymbol{\omega}}$ of the upper-level policy. Hence, we categorize these methods as Policy Gradient methods.

Algorithm 9 Natural Actor Critic

Input: Policy parametrization θ ,
 data-set $\mathcal{D} = \{\mathbf{x}^{[i]}, \mathbf{u}^{[i]}, r^{[i]}, \mathbf{x}'^{[i]}\}_{i=1\dots N}$

for each sample $i = 1 \dots N$ **do**

 Compute features for current and successor state:

$$\boldsymbol{\psi}^{[i]} = \begin{bmatrix} \nabla_{\theta} \log \pi_{\theta}(\mathbf{u}^{[i]} | \mathbf{x}^{[i]}) \\ \boldsymbol{\varphi}(\mathbf{x}^{[i]}) \end{bmatrix}, \quad \boldsymbol{\psi}'^{[i]} = \begin{bmatrix} \mathbf{0} \\ \boldsymbol{\varphi}(\mathbf{x}'^{[i]}) \end{bmatrix}$$

end for

 Compute LSTD-Q solution

$$\mathbf{b} = \sum_{i=1}^N \boldsymbol{\psi}^{[i]} r^{[i]}, \quad \mathbf{A} = \sum_{i=1}^N \boldsymbol{\psi}^{[i]} (\boldsymbol{\psi}^{[i]} - \gamma \boldsymbol{\psi}'^{[i]})^T$$

$$\begin{bmatrix} \mathbf{w} \\ \mathbf{v} \end{bmatrix} = \mathbf{A}^{-1} \mathbf{b}$$

 return $\nabla_{\theta}^{\text{NAC}} J_{\theta} = \mathbf{w}$

Existing natural gradient methods in parameter space do not use a compatible function approximation to estimate the natural gradient but directly try to estimate the Fisher information matrix which is subsequently multiplied with the likelihood gradient $\nabla_{\omega}^{\text{PE}} J_{\omega}$ given in Equation (2.19) in parameter space, i.e.,

$$\nabla_{\omega}^{\text{NES}} J_{\omega} = \mathbf{F}_{\omega}^{-1} \nabla_{\omega}^{\text{PE}} J_{\omega}. \quad (2.38)$$

The natural gradient for parameter space was first used in the Natural Evolution Strategy [88] where the Fisher information matrix was determined empirically. However, the empirical estimation of the Fisher information matrix is problematic as the matrix may not be invertible due to redundant parameters or sampling errors. In [78], the authors compute the Fisher information matrix in closed form for Gaussian policies in parameter space. The authors also give derivations of an optimal baseline for their method. As these derivations are rather complex we refer to the corresponding paper for both derivations. The NES strategy has also been compared with the PEPG algorithm [68], indi-

cating that NES is more efficient for low-dimensional problems while PEPG has advantages for high-dimensional parameter spaces as the second-order type update of the natural gradient gets more difficult.

2.4.2 Expectation Maximization Policy Search Approaches

Policy gradient methods require the user to specify the learning rate. Setting the learning rate can be problematic and often results in an unstable learning process or slow convergence [38]. This problem can be avoided by formulating policy search as an inference problem with latent variables and, subsequently, using the Expectation-Maximization (EM) algorithm to infer a new policy. As in the standard Expectation-Maximization algorithm, the parameter update is computed as a weighted maximum likelihood estimate which has a closed form solution for most of the used policies. Hence, no learning rate is required.

We will first review the standard EM-algorithm and, subsequently, reformulate policy search as an inference problem by treating the reward as improper probability distribution. Finally, we will explain the resulting EM-based policy search algorithms.

2.4.2.1 The Standard Expectation Maximization Algorithm

The EM-algorithm [46, 12] is a well known algorithm for determining the maximum likelihood solution of a probabilistic latent variable model. Let us assume that \mathbf{y} defines an observed random variable, \mathbf{z} an unobserved random variable and $p_{\boldsymbol{\theta}}(\mathbf{y}, \mathbf{z})$ is the parametrized joint distribution of observed and unobserved variables with parameters $\boldsymbol{\theta}$. As \mathbf{z} is unobserved, it needs to be marginalized out to compute the likelihood of the parameters, i.e., $p_{\boldsymbol{\theta}}(\mathbf{y}) = \int p_{\boldsymbol{\theta}}(\mathbf{y}, \mathbf{z}) d\mathbf{z}$. Given a data set $\mathbf{Y} = [\mathbf{y}^{[1]}, \dots, \mathbf{y}^{[N]}]^T$, we now want to maximize the log-marginal-likelihood

$$\log p_{\boldsymbol{\theta}}(\mathbf{Y}) = \sum_{i=1}^N \log p_{\boldsymbol{\theta}}(\mathbf{y}^{[i]}) = \sum_{i=1}^N \log \int p_{\boldsymbol{\theta}}(\mathbf{y}^{[i]}, \mathbf{z}) d\mathbf{z} \quad (2.39)$$

with respect to the parameters $\boldsymbol{\theta}$, where we assumed i.i.d. data-points, i.e., $p_{\boldsymbol{\theta}}(\mathbf{Y}) = \prod_i p_{\boldsymbol{\theta}}(\mathbf{y}^{[i]})$. Since the logarithm is acting on the marginal

distribution $\int p_{\theta}(\mathbf{y}, \mathbf{z})d\mathbf{z}$ instead of the joint distribution $p_{\theta}(\mathbf{y}, \mathbf{z})$, we can not obtain a closed form solution for the parameters θ of our probability model $p_{\theta}(\mathbf{y}, \mathbf{z})$.

The EM-algorithm is an iterative procedure for estimating the maximum likelihood solution of latent variable models where the parameter updates of every iteration can be obtained in closed form. We will closely follow the derivation of EM from [12] as it can directly be applied to the policy search setup.

The costly marginalization over the hidden variables can be avoided by introducing an auxiliary distribution $q(\mathbf{Z})$, which we will denote as variational distribution, that is used to decompose the marginal log-likelihood by using the identity $p_{\theta}(\mathbf{Y}) = p_{\theta}(\mathbf{Y}, \mathbf{Z})/p_{\theta}(\mathbf{Z}|\mathbf{Y})$, i.e.,

$$\begin{aligned} \log p_{\theta}(\mathbf{Y}) &= \int q(\mathbf{Z}) \log p_{\theta}(\mathbf{Y})d\mathbf{Z} = \int q(\mathbf{Z}) \log \frac{q(\mathbf{Z})p_{\theta}(\mathbf{Y}, \mathbf{Z})}{q(\mathbf{Z})p_{\theta}(\mathbf{Z}|\mathbf{Y})}d\mathbf{Z} \\ &= \int q(\mathbf{Z}) \log \frac{p_{\theta}(\mathbf{Y}, \mathbf{Z})}{q(\mathbf{Z})}d\mathbf{Z} - \int q(\mathbf{Z}) \log \frac{p_{\theta}(\mathbf{Z}|\mathbf{Y})}{q(\mathbf{Z})}d\mathbf{Z} \\ &= \mathcal{L}_{\theta}(q) + \text{KL}(q(\mathbf{Z})||p_{\theta}(\mathbf{Z}|\mathbf{Y})). \end{aligned} \quad (2.40)$$

Since the KL-divergence is always larger or equal to zero, the term $\mathcal{L}_{\theta}(q)$ is a lower bound of the log marginal-likelihood $\log p_{\theta}(\mathbf{Y})$. The two update steps in EM each correspond to maximizing the lower bound \mathcal{L} and minimizing the KL-divergence term.

Expectation Step. In the expectation step (E-step), we update the variational distribution $q(\mathbf{Z})$ by minimizing the KL-divergence $\text{KL}(q(\mathbf{Z})||p_{\theta}(\mathbf{Z}|\mathbf{Y}))$ which is equivalent to setting $q(\mathbf{Z}) = p_{\theta}(\mathbf{Z}|\mathbf{Y})$. Note that the lower bound $\mathcal{L}_{\theta}(q)$ is tight after each E-step, i.e., $\log p_{\theta}(\mathbf{Y}) = \mathcal{L}_{\theta}(q)$, as the KL-divergence $\text{KL}(q(\mathbf{Z})||p_{\theta}(\mathbf{Z}|\mathbf{Y}))$ has been set to zero by the E-step. As $\log p_{\theta}(\mathbf{Y})$ is unaffected by the change of $q(\mathbf{Z})$, we observe from Equation (2.40) that the lower bound $\mathcal{L}_{\theta}(q)$ has to increase if we decrease the KL-divergence.

Maximization Step. In the maximization step (M-step), we optimize the lower bound with respect to θ , i.e.,

$$\begin{aligned}\theta_{\text{new}} &= \operatorname{argmax}_{\theta} \mathcal{L}_{\theta}(q) = \operatorname{argmax}_{\theta} \int q(\mathbf{Z}) \log p_{\theta}(\mathbf{Y}, \mathbf{Z}) d\mathbf{Z} + H(q) \\ &= \operatorname{argmax}_{\theta} \mathbb{E}_{q(\mathbf{Z})} [\log p_{\theta}(\mathbf{Y}, \mathbf{Z})] = \operatorname{argmax}_{\theta} \mathcal{Q}_{\theta}(q),\end{aligned}\quad (2.41)$$

where the term $H(q)$ denotes the entropy of q and can be neglected for estimating θ_{new} . The term in Equation (2.41) is also denoted as the *expected complete data log-likelihood* $\mathcal{Q}_{\theta}(q)$. The log now directly acts on the joint distribution $p_{\theta}(\mathbf{Y}, \mathbf{Z})$, and, hence, the M-step can be obtained in closed form. By examining the expected complete data log-likelihood

$$\mathcal{Q}_{\theta}(q) = \int q(\mathbf{Z}) \log p_{\theta}(\mathbf{Y}, \mathbf{Z}) d\mathbf{Z} \quad (2.42)$$

$$= \sum_{i=1}^N \int q_i(\mathbf{z}) \log p_{\theta}(\mathbf{y}^{[i]}, \mathbf{z}) d\mathbf{z} \quad (2.43)$$

in more detail, we can see that the M-step is based on a *weighted maximum likelihood estimate* of θ using the complete data points $[\mathbf{y}^{[i]}, \mathbf{z}]$ weighted by $q_i(\mathbf{z})$.

Note that after the E-step, the KL-term of Equation (2.40) is set to zero and, hence, the KL-term can only increase in the M-step. Consequently, $\log p_{\theta}(\mathbf{Y})$ is increased even more than the lower bound \mathcal{L} . The EM-algorithm is guaranteed to converge to a local maximum of the marginal log-likelihood $\log p_{\theta}(\mathbf{Y})$ as the lower bound is increased in each E-step and M-step, and the lower bound is tight after each E-step.

2.4.2.2 Policy Search as an Inference Problem

We will first formulate policy search as a latent variable inference problem and then show how EM can be applied to solve this problem. To do so, we define a binary reward event R as our observed variable. As we want to maximize the reward, we always want to observe the reward event, i.e., $R = 1$. The probability of this reward event is given by $p(R = 1|\tau)$. The trajectories τ are the latent variables in our model. A graphical model of policy search formulated as an inference problem is given in Figure 2.3.

Maximizing the reward implies maximizing the probability of the reward event, and, hence, our trajectory distribution $p_{\theta}(\tau)$ needs to assign high probability to trajectories with high reward probability $p(R = 1|\tau)$. As we are only concerned with the case $R = 1$ for estimating the trajectory distribution $p_{\theta}(\tau)$, we will write $p(R|\tau)$ instead of $p(R = 1|\tau)$.

If the return $R(\tau)$ of a trajectory is bounded, it can be directly transformed into a non-normalized probability distribution, i.e., $p(R|\tau) \propto R(\tau) - \min_{\tau_j} R(\tau_j)$ [19]. Otherwise, an exponential transformation of the reward signal can be used [59, 84], i.e., $p(R|\tau) \propto \exp(\beta R(\tau))$. This exponential transformation implies a soft-max distribution for the trajectories conditioned on the observation of the reward event, i.e.,

$$p_{\theta}(\tau|R) = \frac{p_{\theta}(\tau) \exp(\beta R(\tau))}{\int p_{\theta}(\tau) \exp(\beta R(\tau)) d\tau}.$$

The parameter β denotes the inverse temperature of the soft-max distribution. The higher we choose β , the more greedy the policy update becomes. This parameter is either specified by the user [38] or can be set by heuristics, for example, setting β to a multiple of the standard deviation of the rewards [48, 49].

We want to find a parameter vector θ that maximizes the probability of the reward event. In other words, we want to find the maximum likelihood solution for the log marginal-likelihood

$$\log p_{\theta}(R) = \log \int_{\tau} p(R|\tau) p_{\theta}(\tau) d\tau. \quad (2.44)$$

As for the standard EM algorithm, a variational distribution $q(\tau)$ is used to decompose the log-marginal likelihood into two terms

$$\log p_{\theta}(R) = \mathcal{L}_{\theta}(q) + \text{KL}(q(\tau) || p_{\theta}(\tau|R)), \quad (2.45)$$

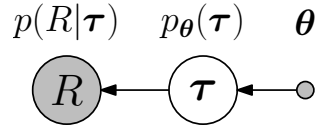


Fig. 2.3 Graphical Model for inference-based policy search. We introduce a binary reward event $R = 1$ as observation, the latent variables are given by the trajectories τ . We want to find the maximum likelihood solution for the parameters θ of observing the reward, i.e., $\theta_{\text{new}} = \text{argmax}_{\theta} \log p_{\theta}(R)$.

where $p(\boldsymbol{\tau}|R)$ is denoted as the *reward-weighted* trajectory distribution

$$p_{\boldsymbol{\theta}}(\boldsymbol{\tau}|R) = \frac{p(R|\boldsymbol{\tau})p_{\boldsymbol{\theta}}(\boldsymbol{\tau})}{\int p(R|\boldsymbol{\tau})p_{\boldsymbol{\theta}}(\boldsymbol{\tau})d\boldsymbol{\tau}} \propto p(R|\boldsymbol{\tau})p_{\boldsymbol{\theta}}(\boldsymbol{\tau}). \quad (2.46)$$

In the E-step, we minimize $\text{KL}(q(\boldsymbol{\tau})||p_{\boldsymbol{\theta}}(\boldsymbol{\tau}|R))$, and, hence, set $q(\boldsymbol{\tau})$ to the reward weighted model distribution $p_{\boldsymbol{\theta}}(\boldsymbol{\tau}|R)$. In the M-step, we maximize the expected complete data log-likelihood

$$\begin{aligned} \operatorname{argmax}_{\boldsymbol{\theta}} \mathcal{Q}_{\boldsymbol{\theta}}(q) &= \operatorname{argmax}_{\boldsymbol{\theta}} \int q(\boldsymbol{\tau}) \log(p(R|\boldsymbol{\tau})p_{\boldsymbol{\theta}}(\boldsymbol{\tau}))d\boldsymbol{\tau} \\ &= \operatorname{argmax}_{\boldsymbol{\theta}} \int q(\boldsymbol{\tau}) \log p_{\boldsymbol{\theta}}(\boldsymbol{\tau})d\boldsymbol{\tau} + f(q) \\ &= \operatorname{argmin}_{\boldsymbol{\theta}} \text{KL}(q(\boldsymbol{\tau})||p_{\boldsymbol{\theta}}(\boldsymbol{\tau})). \end{aligned} \quad (2.47)$$

Hence, the E- and the M-step use different KL-divergences for their iterative updates. We distinguish between two EM update procedures, called Monte-Carlo (MC-)EM approaches [38, 59] and Variational Inference for Policy Search [48], which we will discuss in the following sections. Both update procedures use different approximations to minimize the KL-divergences.

2.4.2.3 Monte-Carlo EM-based Policy Search.

Some of the most efficient policy search methods are Monte-Carlo Expectation-Maximization (MC-EM) methods [38, 59, 38, 85]. The MC-EM algorithm [46] is a variant of EM that uses a sample based approximation for the variational distribution q , i.e., in the E-step, MC-EM minimizes the KL-divergence $\text{KL}(q(\mathbf{Z})||p_{\boldsymbol{\theta}}(\mathbf{Z}|\mathbf{Y}))$ by using samples $Z_j \sim p_{\boldsymbol{\theta}}(\mathbf{Z}|\mathbf{Y})$. Subsequently, these samples \mathbf{Z}_j are used to estimate the expectation of the complete data log-likelihood

$$\mathcal{Q}_{\boldsymbol{\theta}}(q) = \sum_{j=1}^K \log p_{\boldsymbol{\theta}}(\mathbf{Y}, \mathbf{Z}_j). \quad (2.48)$$

In terms of policy search, MC-EM methods use samples $\boldsymbol{\tau}^{[i]}$ from the old trajectory distribution $p_{\boldsymbol{\theta}'}(\boldsymbol{\tau})$, where $\boldsymbol{\theta}'$ represents the old policy parameters, to represent the variational distribution $q(\boldsymbol{\tau}) \propto p(R|\boldsymbol{\tau})p_{\boldsymbol{\theta}'}(\boldsymbol{\tau})$ over trajectories. As $\boldsymbol{\tau}^{[i]}$ has already been sampled from $p_{\boldsymbol{\theta}'}(\boldsymbol{\tau})$, $p_{\boldsymbol{\theta}'}(\boldsymbol{\tau})$

Algorithm 10 Episode-Based MC-EM Policy Updates

Input: inverse temperature β data-set $\mathcal{D}_{\text{ep}} = \left\{ \mathbf{s}^{[i]}, \boldsymbol{\theta}^{[i]}, R^{[i]} \right\}_{i=1 \dots N}$ **Compute** weighting $d^{[i]} = f(R^{[i]})$ for each sample i e.g., $d^{[i]} \propto \exp(\beta R^{[i]})$ **Compute** weighted ML solution, see Equation (2.50) and (2.51)

$$\boldsymbol{\omega}_{\text{new}} = \operatorname{argmax}_{\boldsymbol{\omega}} \sum_{i=1}^N \sum_{t=0}^{T-1} d^{[i]} \log \pi_{\boldsymbol{\omega}}(\boldsymbol{\theta}^{[i]} | \mathbf{s}^{[i]})$$

return $\boldsymbol{\omega}_{\text{new}}$

cancels with the importance weights and can be skipped from the distribution $q(\boldsymbol{\tau}^{[i]})$ and, hence, $q(\boldsymbol{\tau}^{[i]}) \propto p(R | \boldsymbol{\tau}^{[i]})$. These samples are then used in the M-step for estimating the complete-data log-likelihood. Consequently, in the M-step, we have to maximize

$$\mathcal{Q}_{\boldsymbol{\theta}}(\boldsymbol{\theta}') = \sum_{\boldsymbol{\tau}^{[i]} \sim p_{\boldsymbol{\theta}'}(\boldsymbol{\tau})} p(R | \boldsymbol{\tau}^{[i]}) \log p_{\boldsymbol{\theta}'}(\boldsymbol{\tau}^{[i]}) \quad (2.49)$$

with respect to the new policy parameters $\boldsymbol{\theta}$. This maximization corresponds to a weighted maximum likelihood estimate of $\boldsymbol{\theta}$ where each sample $\boldsymbol{\tau}^{[i]}$ is weighted by $d^{[i]} = p(R | \boldsymbol{\tau}^{[i]})$.

Episode-based EM-Algorithms

The episode-based version of MC-EM policy search algorithms can straightforwardly be derived by replacing the trajectory distribution $p_{\boldsymbol{\theta}}(\boldsymbol{\tau}^{[i]})$ by the upper-level policy $\pi_{\boldsymbol{\omega}}(\boldsymbol{\theta})$ and has been used to generalize the upper-level policy to multiply contexts, i.e., learn $\pi_{\boldsymbol{\omega}}(\boldsymbol{\theta} | \mathbf{s})$ [37]. The policy update is given by the weighted maximum likelihood estimate of the parameters $\boldsymbol{\omega}$ of the upper level policy. The general setup for Episode-Based EM-updates is given in Algorithm 10.

Reward Weighted Regression. Reward Weighted Regression (RWR) uses a linear policy for $\pi_{\boldsymbol{\omega}}(\boldsymbol{\theta} | \mathbf{s})$, and, hence, the weighted maximum likelihood estimate performed by the EM-update is given by a

weighted linear regression. RWR was introduced in [59] to learn an inverse dynamics model for operational space control. However, the algorithm straightforwardly generalizes to episode-based policy search with multiple contexts. The policy $\pi_{\omega}(\boldsymbol{\theta}|\mathbf{s}) = \mathcal{N}(\boldsymbol{\theta}|\mathbf{W}^T\boldsymbol{\phi}(\mathbf{s}), \boldsymbol{\Sigma}_{\boldsymbol{\theta}})$ is represented as Gaussian linear model. Given the data-set \mathcal{D}_{ep} and the weightings $d^{[i]}$ for each sample $(\mathbf{s}^{[i]}, \boldsymbol{\theta}^{[i]})$ in \mathcal{D}_{ep} , the weighted maximum likelihood solution for \mathbf{W} is given by

$$\mathbf{W} = (\boldsymbol{\Phi}^T \mathbf{D} \boldsymbol{\Phi} + \lambda \mathbf{I})^{-1} \boldsymbol{\Phi}^T \mathbf{D} \boldsymbol{\Theta}, \quad (2.50)$$

where λ is a ridge factor, $\boldsymbol{\Phi} = [\boldsymbol{\phi}(\mathbf{s}^{[1]}), \dots, \boldsymbol{\phi}(\mathbf{s}^{[N]})]$ contains the feature vectors of the contexts, the diagonal matrix \mathbf{D} contains the weights $d^{[i]}$, and $\boldsymbol{\Theta} = [\boldsymbol{\theta}^{[1]}, \dots, \boldsymbol{\theta}^{[N]}]^T$ the parameter vectors $\boldsymbol{\theta}^{[i]}$. The covariance matrix $\boldsymbol{\Sigma}_{\boldsymbol{\theta}}$ can be updated according to a weighted maximum likelihood estimate. The update equations for $\boldsymbol{\Sigma}_{\boldsymbol{\theta}}$ are given in Appendix B.

Cost-Regularized Kernel Regression. Cost-Regularized Kernel Regression (CRKR) is the kernelized version of Reward-Weighted-Regression, and was one of the first algorithms used to learn upper-level policies for multiple contexts [38]. Similar to most kernel-regression methods, CRKR uses individual regressions for the individual output dimensions. The policy in CRKR $\pi_{\omega}(\boldsymbol{\theta}|\mathbf{s}) = \mathcal{N}(\boldsymbol{\omega}|\boldsymbol{\mu}_{\omega}(\mathbf{s}), \text{diag}(\boldsymbol{\sigma}_{\omega}(\mathbf{s})))$ is modeled as a Gaussian process. The mean and the variance for the h -th output dimension are therefore given by

$$\mu_h(\mathbf{s}) = \mathbf{k}(\mathbf{s})(\mathbf{K} + \lambda \mathbf{C})^{-1} \boldsymbol{\Theta}_h, \quad (2.51)$$

$$\sigma_h^2(\mathbf{s}) = k(\mathbf{s}, \mathbf{s}) + \lambda - \mathbf{k}(\mathbf{s})^T (\mathbf{K} + \lambda \mathbf{C})^{-1} \mathbf{k}(\mathbf{s}). \quad (2.52)$$

The term $\mathbf{K} = \boldsymbol{\Phi}^T \boldsymbol{\Phi}$ denotes the kernel matrix and $\mathbf{k}(\mathbf{s}) = \boldsymbol{\phi}(\mathbf{s})^T \boldsymbol{\Phi}$ represents the kernel vector for a new context query point \mathbf{s} , where the feature matrix $\boldsymbol{\Phi} = [\boldsymbol{\phi}(\mathbf{s}^{[1]}), \dots, \boldsymbol{\phi}(\mathbf{s}^{[N]})]$ contains the feature vectors of the contexts. The matrix \mathbf{C} is denoted as cost matrix because it is inversely related to the reward weighting used in RWR, i.e., $\mathbf{C} = \mathbf{D}^{-1}$. The cost matrix is treated as an input dependent noise prior for the Gaussian process [37].

As the standard kernel-regression formulation (for example, see Bishop [12], chapter 6.1), CRKR can be derived from linear regression

using the Woodbury Identity [37]. Instead of standard linear regression, reward weighted regression is used to derive CRKR.

As CRKR is a non-parametric method, it does not update a parameter vector. Instead, the policy is determined by the given data set. As CRKR is a kernel method, we do not need to specify a feature vector $\phi(\mathbf{s})$, but rather use a kernel. Kernels typically offer more flexibility in modeling a function than user-specified feature vectors. We refer to [65] for more details about kernel-methods for regression. The disadvantage of using a kernel-method is that the output dimensions of the policy $\pi_{\omega}(\boldsymbol{\theta}|\mathbf{s})$ are typically modeled as independent Gaussian distributions, and, hence, no correlations can be modeled. Such uncorrelated exploration strategies might result in a decreased performance of the algorithm as we will discuss in Section 2.1.

Step-based EM-Algorithms

Step-based EM-Algorithms decompose the complete data log-likelihood $\mathcal{Q}_{\boldsymbol{\theta}}(\boldsymbol{\theta}')$ into the single steps of the episode. We denote the parameter vector $\boldsymbol{\theta}'$ as the parameters of the old policy. We first show that $\mathcal{Q}_{\boldsymbol{\theta}}(\boldsymbol{\theta}')$ is a lower bound of the logarithmic expected return $\log J_{\boldsymbol{\theta}}$ where we will assume that no reward transformation has been used, i.e., $p(R|\boldsymbol{\tau}) \propto R(\boldsymbol{\tau})$,

$$\log J_{\boldsymbol{\theta}} = \log \int p_{\boldsymbol{\theta}}(\boldsymbol{\tau})R(\boldsymbol{\tau})d\boldsymbol{\tau} = \log p_{\boldsymbol{\theta}}(R). \quad (2.53)$$

As we know that $\mathcal{Q}_{\boldsymbol{\theta}}(\boldsymbol{\theta}')$ is a lower bound of $\log p_{\boldsymbol{\theta}}(R)$, we conclude that

$$\log J_{\boldsymbol{\theta}} \geq \mathcal{Q}_{\boldsymbol{\theta}}(\boldsymbol{\theta}') = \mathbb{E}_{p_{\boldsymbol{\theta}'}(\boldsymbol{\tau})} [R(\boldsymbol{\tau}) \log p_{\boldsymbol{\theta}}(\boldsymbol{\tau})], \quad (2.54)$$

where the *old* policy parameters $\boldsymbol{\theta}'$ have been used for generating the roll-outs. By differentiating $\mathcal{Q}_{\boldsymbol{\theta}}(\boldsymbol{\theta}')$ with respect to the new policy parameters $\boldsymbol{\theta}$, we get

$$\begin{aligned} \nabla_{\boldsymbol{\theta}} \mathcal{Q}_{\boldsymbol{\theta}}(\boldsymbol{\theta}') &= \mathbb{E}_{p_{\boldsymbol{\theta}'}(\boldsymbol{\tau})} [R(\boldsymbol{\tau}) \nabla_{\boldsymbol{\theta}} \log p_{\boldsymbol{\theta}}(\boldsymbol{\tau})], \\ &= \mathbb{E}_{p_{\boldsymbol{\theta}'}(\boldsymbol{\tau})} \left[R(\boldsymbol{\tau}) \sum_{t=0}^{T-1} \nabla_{\boldsymbol{\theta}} \log \pi_{\boldsymbol{\theta}}(\mathbf{u}_t|\mathbf{x}_t, t) \right]. \end{aligned} \quad (2.55)$$

Algorithm 11 Step-Based MC-EM Policy Updates

Input: Policy parametrization θ

$$\text{data-set } \mathcal{D} = \left\{ \mathbf{x}_{1:T}^{[i]}, \mathbf{u}_{1:T-1}^{[i]}, Q_t^{[i]} \right\}_{i=1 \dots N}$$

Compute weighting $d_t^{[i]} \propto Q_t^{[i]}$ or $d_t^{[i]} \propto \exp(\beta Q_t^{[i]})$

Compute weighted maximum ML estimate, see Equations (2.58) and (2.59)

$$\theta_{\text{new}} = \operatorname{argmax}_{\theta} \sum_{i=1}^N \sum_{t=0}^{T-1} d_t^{[i]} \log \pi_{\theta} \left(\mathbf{u}_t^{[i]} \mid \mathbf{x}_t^{[i]}, t \right)$$

Using the same insight as we used for the policy gradient theorem, i.e., past rewards are independent of future actions, we obtain

$$\nabla_{\theta} \mathcal{Q}_{\theta}(\theta') = \mathbb{E}_{p_{\theta'}(\tau)} \left[\sum_{t=0}^{T-1} Q_t^{\pi}(\mathbf{x}_t, \mathbf{u}_t) \nabla_{\theta} \log \pi_{\theta}(\mathbf{u}_t \mid \mathbf{x}_t, t) \right]. \quad (2.56)$$

Setting Equation (2.56) to zero corresponds to performing a weighted maximum likelihood estimate on the step-based data-set $\mathcal{D}_{\text{step}}$ for obtaining the new parameters θ of policy $\pi_{\theta}(\mathbf{u}_t \mid \mathbf{x}_t, t)$. The weighting is in the step-based case given by

$$Q_t^{\pi}(\mathbf{x}_t^{[i]}, \mathbf{u}_t^{[i]}) \approx \sum_{h=t}^{T-1} r(\mathbf{x}_h^{[i]}, \mathbf{u}_h^{[i]}).$$

From Equation (2.56) we can see that step-based EM algorithms reduce policy search to an iterative reward-weighted imitation learning procedure. This formulation is used to derive the widely used EM-based policy search algorithm, Policy learning by Weighting Exploration with Returns (PoWER), which was introduced in [38]. The general algorithm for step-based EM algorithms is illustrated in Algorithm 11.

Relation to Policy Gradients. There is a close connection between step-based gradient methods which were introduced in Section 2.4.1 and the step-based EM-based approach. In the limit, if the new parameters θ are close to the old parameters θ' , we obtain the policy

gradient theorem update rule from the lower bound $\mathcal{Q}_\theta(\theta')$,

$$\begin{aligned} \lim_{\theta \rightarrow \theta'} \nabla_\theta \mathcal{Q}_\theta(\theta') &= \lim_{\theta \rightarrow \theta'} \mathbb{E}_{p_{\theta'}(\tau)} \left[\sum_{t=0}^{T-1} Q_t^\pi(\mathbf{x}_t, \mathbf{u}_t) \nabla_\theta \log \pi_\theta(\mathbf{u}_t | \mathbf{x}_t) \right] \\ &= \mathbb{E}_{p_\theta(\tau)} \left[\sum_{t=0}^{T-1} Q_t^\pi(\mathbf{x}_t, \mathbf{u}_t) \nabla_\theta \log \pi_\theta(\mathbf{u}_t | \mathbf{x}_t) \right] \\ &= \nabla_\theta^{\text{PG}} J_\theta, \end{aligned} \quad (2.57)$$

From this comparison we conclude that, unlike policy gradient methods, the EM-based approach allows us to use a different parameter vector θ' for the expectation of the trajectories than for the estimation of the gradient $\nabla_\theta \log p_\theta(\tau)$. The EM-based approach aims at making actions with high future reward $Q_t^\pi(\mathbf{x}, \mathbf{u})$ more likely. However, in contrast to policy gradient methods, it neglects the influence of the policy update $\delta\theta$ on the trajectory distribution $p_{\theta+\delta\theta}(\tau)$.

However, such a relationship can only be obtained, if we can linearly transform the reward into an improper probability distribution. If we need to use an exponential transformation for the rewards, such a direct relationship with policy gradient methods cannot be established.

Episodic Reward Weighted Regression. Despite the name, Episodic Reward Weighted Regression (eRWR) is the step-based extension of RWR [38], which we presented in the previous section. Similar to RWR, episodic RWR assumes a linear model for the policy, which is in the step-based case given as $\pi_\theta(\mathbf{u}_t | \mathbf{x}_t, t) = \mathcal{N}(\mathbf{u}_t | \mathbf{W}^T \phi_t(\mathbf{x}), \Sigma_{\mathbf{u}})$. The weighted maximum likelihood estimate for $\theta = \{\mathbf{W}, \Sigma_{\mathbf{u}}\}$ is now performed on the step-based data set $\mathcal{D}_{\text{step}}$ with inputs $\mathbf{x}^{[i]}$, target vectors $\mathbf{u}^{[i]}$ and weightings $d^{[i]}$. As we have a Gaussian policy which is linear in the feature vectors ϕ_t , the weighted ML estimate of the weight vector \mathbf{W} is given by a weighted least squares linear regression,

$$\mathbf{W}_{\text{new}} = (\Phi^T \mathbf{D} \Phi)^{-1} \Phi^T \mathbf{D} \mathbf{U}, \quad (2.58)$$

where $\Phi = [\phi_0^{[1]}, \dots, \phi_{T-1}^{[1]}, \dots, \phi_0^{[N]}, \dots, \phi_{T-1}^{[N]}]$ contains the feature vectors for all time steps t of all trajectories $\tau^{[i]}$, \mathbf{D} is the diagonal weighting matrix containing the weightings $d_t^{[i]}$ of each sample and the

matrix \mathbf{U} contains the control vectors $\mathbf{u}_t^{[i]}$ for all t and i . For more details on Equation (2.58) please refer to Appendix B. The update of the covariance matrix $\Sigma_{\mathbf{u}}$ of $\pi_{\theta}(\mathbf{u}|\mathbf{x})$ can also be obtained by a weighted maximum likelihood estimate and is given in Appendix B, see Equation (4.4).

Policy learning by Weighting Exploration with Returns. The policy update of the PoWER algorithm [38] is similar to episodic RWR, however, PoWER uses a more structured exploration strategy which typically results in better performance of PoWER in comparison to episodic RWR. To simplify our discussion we will for now assume that the control action u_t is one dimensional. RWR directly perturbs the controls $u_t = \mathbf{w}^T \phi_t + \epsilon_t$ with zero-mean Gaussian noise. Instead, PoWER applies the perturbation to the parameters \mathbf{w} , i.e. $u_t = (\mathbf{w} + \epsilon_t)^T \phi_t$, where $\epsilon_t \sim \mathcal{N}(\mathbf{0}, \Sigma_{\mathbf{w}})$ is the noise term applied to the parameter vector at time step t . Such a policy can be written as $\pi_t(\mathbf{u}_t|\mathbf{x}_t) = \mathcal{N}(u_t|\mathbf{w}^T \phi_t, \phi_t^T \Sigma_{\mathbf{w}} \phi_t)$, i.e., as Gaussian policy where the variance also depends on the current features Φ_t .

If we assume that $\Sigma_{\mathbf{w}}$ is known, we can again determine the weighted maximum likelihood solution for the parameters \mathbf{w} . This solution is given in [38] as

$$\mathbf{w}_{\text{new}} = \mathbf{w}_{\text{old}} + \mathbb{E} \left[\sum_{t=0}^{T-1} \mathbf{L}_t(\mathbf{x}) Q_t^{\pi}(\mathbf{x}_t, \mathbf{u}_t) \right]^{-1} \mathbb{E} \left[\sum_{t=0}^{T-1} \mathbf{L}_t(\mathbf{x}) Q_t^{\pi}(\mathbf{x}_t, \mathbf{u}_t) \epsilon_t \right],$$

where $\mathbf{L}_t(\mathbf{x}) = \phi_t(\mathbf{x}) \phi_t(\mathbf{x})^T (\phi_t(\mathbf{x})^T \Sigma_{\mathbf{w}} \phi_t(\mathbf{x}))^{-1}$. As we can see, the exploration noise ϵ_t is weighted by the returns Q_t^{π} to obtain the new parameter vector \mathbf{w} . For the derivation of this equation we refer to [38]. However, the update rule of PoWER can also be written in terms of matrices, where we use the action vectors \mathbf{u}_t instead of using the noise terms ϵ_t as target values, i.e.,

$$\mathbf{w}_{\text{new}} = (\Phi^T \tilde{\mathbf{D}} \Phi)^{-1} \Phi^T \tilde{\mathbf{D}} \mathbf{U}, \quad (2.59)$$

where \mathbf{U} contains the actions of all time steps and all trajectories, Φ^T is defined as in Equation (2.50) and $\tilde{\mathbf{D}} \in \mathbb{R}^{NT \times NT}$ is a diagonal weighting matrix with the entries $\tilde{d}_t^{[i]} = \left(\phi_t(\mathbf{x}^{[i]})^T \Sigma_{\mathbf{w}} \phi_t(\mathbf{x}^{[i]}) \right)^{-1} 1 Q_t^{[i]}$ for each

sample i and time step t . For further details please consult Appendix B. We can now see, as the only difference between the eRWR policy and the policy used in PoWER is the state-dependent variance term, the only difference in the policy update is that in PoWER the data-points are additionally weighted by the precision $\left(\phi_t(\mathbf{x}^{[i]})^T \Sigma_{\mathbf{w}} \phi_t(\mathbf{x}^{[i]})\right)^{-1}$ of the policy for state $\mathbf{x}^{[i]}$. Consequently, data-points with less variance have a higher influence on the result of the regression. As this notation does not contain the noise terms $\epsilon_t^{[i]}$, it is also compatible with Algorithm 11.

2.4.2.4 Variational Inference-based Methods

As we have seen, the MC-EM approach uses a weighted maximum likelihood estimate to obtain the new parameters θ of the policy. While a weighted maximum likelihood estimate can be computed efficiently, such an approach might also suffer from a caveat: it averages over several modes of the reward function. Such a behavior might result in slow convergence to good policies as the average of several modes might be in an area with low reward [48].

Moment Projection and Information Projection. We observe that the maximization used for the MC-EM approach as defined in Equation (2.49) is equivalent to minimizing

$$\text{KL}(p(R|\tau)p_{\theta'}(\tau)||p_{\theta}(\tau)) = \int p(R|\tau^{[i]})p_{\theta'}(\tau^{[i]}) \log \frac{p(R|\tau)p_{\theta'}(\tau^{[i]})}{p_{\theta}(\tau^{[i]})}$$

with respect to the new policy parameters θ . This minimization is also called the *Moment-Projection* of the reward weighted trajectory distribution as it matches the moments of $p_{\theta}(\tau)$ with the moments of $p(R|\tau)p_{\theta'}(\tau)$. It forces $p_{\theta}(\tau)$ to have probability mass everywhere where $p(R|\tau)p_{\theta'}(\tau)$ has non-negligible probability mass. Consequently, if $p_{\theta}(\tau)$ is a Gaussian, the M-projection averages over all modes of the reward weighted trajectory distribution.

Alternatively, we can use the Information (I)-projection $\text{argmin}_{\theta} \text{KL}(p_{\theta}(\tau)||p(R|\tau)p_{\theta'}(\tau))$ to update the policy, as introduced in the Variational Inference for Policy Search [48] algorithm.

This projection forces the new trajectory distribution $p_{\theta}(\tau)$ to be zero everywhere where the reward weighted trajectory distribution is zero. When using a Gaussian distribution for $p_{\theta}(\tau)$, the I-projection will concentrate on a single mode of $p(R|\tau)p_{\theta'}(\tau)$ and lose information about the other modes contained in the samples. Unfortunately, it can not be determined in closed form for most distributions.

Variational Inference for Policy Search. In the Variational Inference for policy search approach [48], a parametric representation of the variational distribution $q_{\beta}(\tau)$ is used instead of a sample-based approximation as used in the MC-EM approach. It is convenient to choose $q_{\beta}(\tau)$ from the same family of distributions as $p_{\theta}(\tau)$. Now, a sample-based approximation is used to replace the integral in the KL-divergence $\text{KL}(q_{\beta}(\tau) \| p(R|\tau)p_{\theta'}(\tau))$ needed for the E-step, i.e.,

$$\begin{aligned} \beta &\in \operatorname{argmin}_{\tilde{\beta}} \text{KL}(q_{\tilde{\beta}}(\tau) \| p(R|\tau)p_{\theta'}(\tau)) \\ &\approx \operatorname{argmin}_{\tilde{\beta}} \sum_{\tau^{[i]}} q_{\tilde{\beta}}(\tau^{[i]}) \log \frac{q_{\tilde{\beta}}(\tau^{[i]})}{p(R|\tau^{[i]})p_{\theta'}(\tau^{[i]})}. \end{aligned} \quad (2.60)$$

The minimization of this KL-divergence is equivalent to the *I-projection* of the reward-weighted trajectory distribution $p(R|\tau)p_{\theta}(\tau)$. In the variational approach, the M-step now trivially reduces to setting the new parameter vector θ_{new} to θ' .

Hence, the MC-EM and the variational inference algorithm only differ in the employed projections of the reward-weighted trajectory distribution $p_{\theta}(\tau|R)$. As the projections are in general different, they each converge to a different (local) maximum of the lower bound $\mathcal{L}_{\theta}(q)$. The variational inference algorithm has been used in the episode-based formulation to learn an upper-level policy $\pi_{\omega}(\theta|\mathbf{s})$ for multiple contexts. If we use a Gaussian distribution for $\pi_{\omega}(\theta|\mathbf{s})$, the I-projection concentrates on a single mode, as shown in Figure 2.4. Such behavior can be beneficial if all modes are almost equally good. However, the I-projection might also choose a sub-optimal mode (which has a lower reward). The M-projection averages over all modes, and, therefore, might also include large areas of low reward in the distribution. The behavior of both approaches for a simple multi-modal toy problem

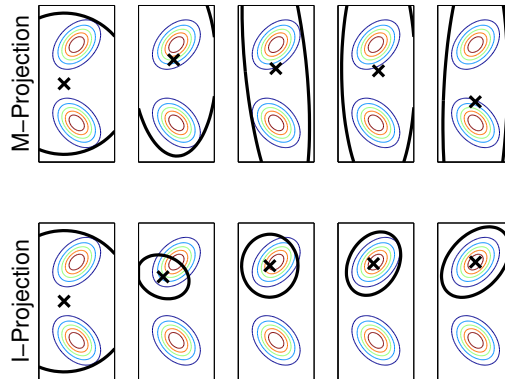


Fig. 2.4 Comparison of using the I-projection versus the M-projection for inference based policy search. While the M-projection averages over several modes of the reward function, the I-projection concentrates on a single mode, and, therefore, avoids including areas of low reward in the policy distribution.

is illustrated in Figure 2.4.

If the target distribution is uni-modal, both projections yield almost the same solutions. However, using the I-projection is computationally demanding, and, hence, the variational inference approach is generally not the method of choice. In addition, if we use a more complex distribution for modeling the policy, e.g., a mixture of Gaussians, the difference between the I- and the M-projection becomes less distinct.

2.4.3 Information-Theoretic Approaches

The main idea behind the information-theoretic approaches is to stay close to the ‘data’, i.e., the trajectory distribution after the policy update should not jump away from the trajectory distribution before the policy update. Information-theoretic approaches bound the distance between the old trajectory distribution $q(\boldsymbol{\tau})$ and the newly estimated trajectory distribution $p(\boldsymbol{\tau})$ at each update step. Such regularization of the policy update limits the information loss of the updates, and, hence, avoids that the new distribution $p(\boldsymbol{\tau})$ prematurely concentrates on local optima of the reward landscape. The first type of algorithms to implement this insight from information theory were the natural policy gradient algorithms [61], which have already been discussed in

the policy gradient section. Natural policy gradient algorithms always require a user-specified learning rate, an issue which was alleviated by EM-based methods. However, EM-based methods have other problems concerning premature convergence and stability of the learning process as they typically do not stay close to the data. The information theoretic insight was again taken up in [57] with the Relative Entropy Policy Search (REPS) algorithm to combine the advantages of both types of algorithms. REPS uses the same information theoretic bound as the NAC algorithm but simultaneously updates its policy by weighted maximum likelihood estimates, which do not require a learning rate.

REPS formulates the policy search problem as an optimization problem, wherein the optimization is done directly in the space of distributions p over trajectories, state-actions pairs or parameters without considering a direct or indirect parametrization of p . As we will see, the REPS optimization problem allows for a closed-form solution for computing p . The used distance measure $\text{KL}(p||q)$ forces p to be low everywhere where the old ‘data’ distribution q is also low. Intuitively, bounding $\text{KL}(p||q)$ prevents p to ‘move outside’ the ‘old data’ distribution q as such behavior is potentially dangerous for the robot. The use of the opposite KL-divergence $\text{KL}(q||p)$ would not exhibit this favorable property and also does not allow for a closed form solution for p .

2.4.3.1 Episode-based Relative Entropy Policy Search

We start our discussion with the episode-based formulation [17] of relative entropy policy search [57] as it is the simplest formulation. In the episode-based formulation, we need to learn an upper-level policy $\pi_{\omega}(\theta)$ for selecting the parameters of the lower-level policy $\pi_{\theta}(\mathbf{u}_t|\mathbf{x}_t)$ in order to maximize the average return J_{ω} as defined in Equation (2.2). At the same time, we want to bound the Kullback-Leibler divergence between the newly estimated policy $\pi_{\omega}(\theta)$ and the old policy $q(\omega)$.

Resulting Optimization Program. To do so, we can solve the following constrained optimization problem

$$\begin{aligned} \max_{\pi} \quad & \int \pi(\boldsymbol{\theta}) R(\boldsymbol{\theta}) d\boldsymbol{\theta}, \\ \text{s. t.} \quad & \epsilon \geq \int \pi(\boldsymbol{\theta}) \log \frac{\pi(\boldsymbol{\theta})}{q(\boldsymbol{\theta})} d\boldsymbol{\theta}, \\ & 1 = \int \pi(\boldsymbol{\theta}) d\boldsymbol{\theta}, \end{aligned} \quad (2.61)$$

This constrained optimization problem can be solved efficiently by the method of Lagrangian multipliers. Please refer to Appendix C for more details. From the Lagrangian, we can also obtain a closed-form solution for the new policy

$$\pi(\boldsymbol{\theta}) \propto q(\boldsymbol{\theta}) \exp\left(\frac{R(\boldsymbol{\theta})}{\eta}\right), \quad (2.62)$$

where η is the Lagrangian multiplier connected to the KL-bound constrained. It specifies a scaling factor for the reward and can be interpreted as the temperature of the soft-max distribution given in Equation (2.62).

The Dual Function. The parameter η is obtained by minimizing the dual-function $g(\eta)$ of the original optimization problem,

$$g(\eta) = \eta\epsilon + \eta \log \int q(\boldsymbol{\theta}) \exp\left(\frac{R(\boldsymbol{\theta})}{\eta}\right) d\boldsymbol{\theta}. \quad (2.63)$$

The derivation of the dual function is given in Appendix C. In practice, the integral in the dual function is approximated by samples, i.e.,

$$g(\eta) = \eta\epsilon + \eta \log \sum_i \frac{1}{N} \exp\left(\frac{R(\boldsymbol{\theta}^{[i]})}{\eta}\right) d\boldsymbol{\theta}. \quad (2.64)$$

Estimating the New Policy. The new policy $\pi(\boldsymbol{\theta})$ is also only known for samples $\boldsymbol{\theta}^{[i]}$ where we have evaluated the reward $R(\boldsymbol{\theta}^{[i]})$. Consequently, we need to fit a parametric distribution $\pi_{\omega}(\boldsymbol{\theta})$ to our samples. This parametric distribution is obtained by a weighted maximum likelihood estimate on the samples $\boldsymbol{\theta}^{[i]}$ with the weightings

$d^{[i]} = \exp\left(R(\boldsymbol{\theta}^{[i]})/\eta\right)$. Note that the distribution $q(\boldsymbol{\theta}^{[i]})$ can be dropped from the weighting as we have already sampled from q . Typically, the parametric policy is Gaussian, and, hence the new parameters are given by the weighted mean and covariance. For the resulting updates of the weighted maximum likelihood estimates please refer to Appendix B.

In theory, the expected return $R(\boldsymbol{\theta}^{[i]}) = \int p_{\boldsymbol{\theta}^{[i]}}(\boldsymbol{\tau})R(\boldsymbol{\tau})d\boldsymbol{\tau}$ is given as expectation over all possible roll-outs. However, to increase the sample efficiency in practice, the return $R(\boldsymbol{\theta}^{[i]})$ is typically approximated with a single sample. Similar to EM-based approaches, such strategy introduces a bias into the optimization problem, as the expectation is not performed inside the exp function, and, consequently, the resulting policy is risk-seeking. However, for moderately stochastic system no performance loss was observed.

Although it seems natural to define $q(\mathbf{x}, \mathbf{u})$ as the old policy $\pi(\boldsymbol{\omega}_{k-1})$, we can use the last K policies as $q(\mathbf{x}, \mathbf{u})$ to reuse samples from previous iterations.

2.4.3.2 Episode-Based Extension to Multiple Contexts.

In episode-based REPS, we can extend the upper-level policy $\pi_{\boldsymbol{\omega}}(\boldsymbol{\theta}|\mathbf{s})$ to select the policy parameters $\boldsymbol{\theta}$ based on the context \mathbf{s} . Our aim is to maximize the expected reward $\mathcal{R}_{\mathbf{s}\boldsymbol{\theta}}$ while bounding the expected relative entropy between $\pi_{\boldsymbol{\omega}}(\boldsymbol{\theta}|\mathbf{s})$ and the old policy $q(\mathbf{u}|\mathbf{x})$, i.e.,

$$\begin{aligned} \max_{\pi} \int \mu(\mathbf{s}) \int \pi_{\boldsymbol{\omega}}(\boldsymbol{\theta}|\mathbf{s}) \mathcal{R}_{\mathbf{s}\boldsymbol{\theta}} d\boldsymbol{\theta} d\mathbf{s} \\ \text{s.t.: } \epsilon \geq \int \mu(\mathbf{s}) \text{KL}(\pi_{\boldsymbol{\omega}}(\boldsymbol{\theta}|\mathbf{s})||q(\boldsymbol{\theta}|\mathbf{s})) d\mathbf{s}, \quad \forall \mathbf{s} : 1 = \int \pi_{\boldsymbol{\omega}}(\boldsymbol{\theta}|\mathbf{s}) d\boldsymbol{\theta}. \end{aligned} \tag{2.65}$$

However, this formulation requires that we have access to many parameter vector samples $\boldsymbol{\theta}^{[i,j]}$ for a single context vector $\mathbf{s}^{[i]}$. In order to relax this assumption, contextual REPS optimizes for the joint probabilities $p(\mathbf{s}, \boldsymbol{\theta})$ and enforces that $p(\mathbf{s}) = \int p(\mathbf{s}, \boldsymbol{\theta}) d\boldsymbol{\theta}$ still reproduces the correct context distribution $\mu(\mathbf{s})$ by using the constraints $\forall \mathbf{s} : p(\mathbf{s}) = \mu(\mathbf{s})$.

Matching Feature Averages. Yet, in continuous spaces, this formulation results in an infinite amount of constraints. Therefore, we need to resort to matching feature expectations instead of matching single probabilities, i.e., $\int p(\mathbf{s})\boldsymbol{\varphi}(\mathbf{s})d\mathbf{s} = \hat{\boldsymbol{\varphi}}\boldsymbol{\phi}$ where $\hat{\boldsymbol{\varphi}}\boldsymbol{\phi} = \int \mu(\mathbf{s})\boldsymbol{\varphi}(\mathbf{s})d\mathbf{s}$ is the observed average feature vector. For example, if $\boldsymbol{\varphi}$ contains all linear and quadratic terms of the context \mathbf{s} , we match the first and second order moments of both distribution, i.e., mean and variance.

Optimization Program for Contextual Policy Search. The resulting optimization program yields

$$\begin{aligned} & \max_p \iint p(\mathbf{s}, \boldsymbol{\theta}) \mathcal{R}_{\mathbf{s}\boldsymbol{\theta}} d\boldsymbol{\theta} d\mathbf{s} \\ \text{s.t: } & \epsilon \geq \iint p(\mathbf{s}, \boldsymbol{\theta}) \log \frac{p(\mathbf{s}, \boldsymbol{\theta})}{\mu(\mathbf{s})q(\boldsymbol{\theta}|\mathbf{s})} d\boldsymbol{\theta} d\mathbf{s} \\ & \hat{\boldsymbol{\phi}} = \iint p(\mathbf{s}, \boldsymbol{\theta}) \boldsymbol{\phi}(\mathbf{s}) d\boldsymbol{\theta} d\mathbf{s}, \quad 1 = \iint p(\mathbf{s}, \boldsymbol{\theta}) d\boldsymbol{\theta} d\mathbf{s}. \end{aligned} \quad (2.66)$$

Note that, by replacing $p(\mathbf{s})$ with $\mu(\mathbf{s})$, we end up with the original optimization problem from Equation (2.65). This optimization problem can be solved by the method of Lagrangian multipliers and yields a closed-form solution for $p(\mathbf{s}, \boldsymbol{\theta})$ that is given by

$$p(\mathbf{s}, \boldsymbol{\theta}) \propto q(\boldsymbol{\theta}|\mathbf{s})\mu(\mathbf{s}) \exp\left(\frac{\mathcal{R}_{\mathbf{s}\boldsymbol{\theta}} - V(\mathbf{s})}{\eta}\right), \quad (2.67)$$

where $V(\mathbf{s}) = \boldsymbol{\phi}(\mathbf{s})^T \mathbf{v}$ is a context dependent baseline which is subtracted from the reward. The parameters η and \mathbf{v} are Lagrangian multipliers which are obtained by minimizing the dual function $g(\eta, \boldsymbol{\theta})$ of the optimization problem [17]. The dual function is given in the Appendix C.

The function $V(\mathbf{s})$ also has an interesting interpretation, which can be obtained when looking at the optimality condition for $V(\mathbf{s}_i) = v_i$ for nominal context variables³, $V_i = \eta \log \int q(\boldsymbol{\theta}|\mathbf{s}_i) \exp(\mathcal{R}_{\mathbf{s}_i\boldsymbol{\theta}}/\eta) d\boldsymbol{\theta}$, and, hence, $V(\mathbf{s}_i)$ is given by a soft-max operator over the expected rewards in context \mathbf{s}_i . Consequently, $V(\mathbf{s}_i)$ can be interpreted as value function [57].

³In this case, we do not have to use features.

Algorithm 12 Episode-Based REPS Updates for Multiple Contexts

Input: KL-bounding ϵ

$$\text{data-set } \mathcal{D}_{ep} = \left\{ \mathbf{s}^{[i]}, \boldsymbol{\theta}^{[i]}, R^{[i]} \right\}_{i=1 \dots N}$$

Optimize dual-function $[\eta, \mathbf{v}] = \operatorname{argmin}_{\eta', \mathbf{v}'} g(\eta', \mathbf{v}')$, s.t. $\eta > 0$

$$g(\eta, \mathbf{v}) = \eta\epsilon + \mathbf{v}^T \hat{\boldsymbol{\varphi}} + \eta \log \left(\sum_{i=1}^N \frac{1}{N} \exp \left(\frac{R^{[i]} - \mathbf{v}^T \boldsymbol{\varphi}(\mathbf{s}^{[i]})}{\eta} \right) \right)$$

Obtain parametric policy $\pi_{\boldsymbol{\omega}}(\boldsymbol{\theta}|\mathbf{s})$ by weighted ML estimate

$$\boldsymbol{\omega}_{\text{new}} = \operatorname{argmax}_{\boldsymbol{\omega}} \sum_{i=1}^N \exp \left(\frac{R^{[i]} - \mathbf{v}^T \boldsymbol{\varphi}(\mathbf{s}^{[i]})}{\eta} \right) \log \pi_{\boldsymbol{\omega}}(\boldsymbol{\theta}^{[i]}|\mathbf{s}^{[i]})$$

The dual function and the new policy $\pi(\boldsymbol{\theta}|\mathbf{s})$ is again computed based on samples. Subsequently, a new parametric distribution is obtained by performing a weighted maximum likelihood (ML) estimate. Typically, a linear Gaussian policy is used to represent the upper-level policy. The weighted ML updates for this policy are given in Appendix B. The episode-based REPS algorithm for generalizing the upper-level policy to multiple contexts is given in Algorithm 12.

2.4.3.3 Learning Multiple Solutions with REPS

Using maximum likelihood estimates for the parameter updates is also beneficial for learning multiple solutions of a motor task as we can represent these multiple solutions as mixture model [17]. REPS can be extended to learning multiple solutions by reformulating the problem as a latent variable estimation problem. The upper-level policy $\pi_{\boldsymbol{\omega}}(\boldsymbol{\theta}|\mathbf{s})$ is now extended with another layer of hierarchy that consists of a gating-policy $\pi(o|\mathbf{s})$ which selects the option o to execute given the current context \mathbf{s} . Subsequently, the option policy $\pi(\boldsymbol{\theta}|\mathbf{s}, o)$ selects the parameter vector $\boldsymbol{\theta}$ of the lower level policy which controls the robot. The upper-level policy can now be written as mixture model, i.e.,

$$\pi_{\boldsymbol{\omega}}(\boldsymbol{\theta}|\mathbf{s}) = \sum_o \pi(o|\mathbf{s})\pi(\boldsymbol{\theta}|\mathbf{s}, o). \quad (2.68)$$

With such a hierarchical approach, we can represent multiple solutions for the same motor task, such as, fore-hand and back-hand strokes in robot table tennis [38]. The gating policy allows inferring which options are feasible for context, allowing us to construct complex policies out of simpler ‘option policies’.

Hierarchical Policy Search as a Latent Variable Estimation

Problem. For efficient data-usage, we have to allow parameter vectors $\boldsymbol{\theta}$ from other options o' to be used to update the option policy $\pi(\boldsymbol{\theta}|\mathbf{s}, o)$ of option o . To achieve this goal, the options are treated as latent variables. Hence, only $q(\mathbf{s}, \boldsymbol{\theta})$ can be accessed and not $q(\mathbf{s}, \boldsymbol{\theta}, o)$. The bound on the KL can still be written in terms of the joint distribution $p(\mathbf{s}, \boldsymbol{\theta}, o) = p(\mathbf{s})\pi(\mathbf{s}|o)\pi(\boldsymbol{\theta}|\mathbf{s}, o)$ as

$$\epsilon \geq \sum_o \iint p(\mathbf{s}, \boldsymbol{\theta}, o) \log \left(\frac{p(\mathbf{s}, \boldsymbol{\theta}, o)}{q(\mathbf{s}, \boldsymbol{\theta})p(o|\mathbf{s}, \boldsymbol{\theta})} \right) d\mathbf{s}d\boldsymbol{\theta}, \quad (2.69)$$

where $p(o|\mathbf{s}, \boldsymbol{\theta})$ is obtained by Bayes theorem. Furthermore, options should not overlap as we want to learn distinct solutions for the motor task. As a measure for the overlap of the options, the expected entropy of $p(o|\mathbf{s}, \boldsymbol{\theta})$ is used, i.e.,

$$\mathbb{E}_{\mathbf{s}, \boldsymbol{\theta}} [H(p(o|\mathbf{s}, \boldsymbol{\theta}))] = - \sum_o \iint p(\mathbf{s}, \boldsymbol{\theta}, o) \log p(o|\mathbf{s}, \boldsymbol{\theta}) d\mathbf{s}d\boldsymbol{\theta}. \quad (2.70)$$

The overlap of the options should decrease by a certain percentage in each policy update step. Hence, the following constraint is introduced

$$\kappa \geq \mathbb{E}_{\mathbf{s}, \boldsymbol{\theta}} [H(p(o|\mathbf{s}, \boldsymbol{\theta}))]. \quad (2.71)$$

The upper bound κ is usually set as a percentage of the currently measured overlap \hat{H}_q , i.e., $\kappa = \hat{H}_q \tilde{\kappa}$, where $1 - \tilde{\kappa}$ denotes the desired decrease of the overlap. Figure 2.5 illustrates the resulting policy updates with and without bounding the overlap on a simple multi-modal reward function. Without bounding the overlap of the options, both options concentrate on both modes of the reward function. As a consequence, the quality of both options is rather poor. By introducing the overlap constraint, both options separate early in the optimization process, and, thus, concentrate on the individual modes.

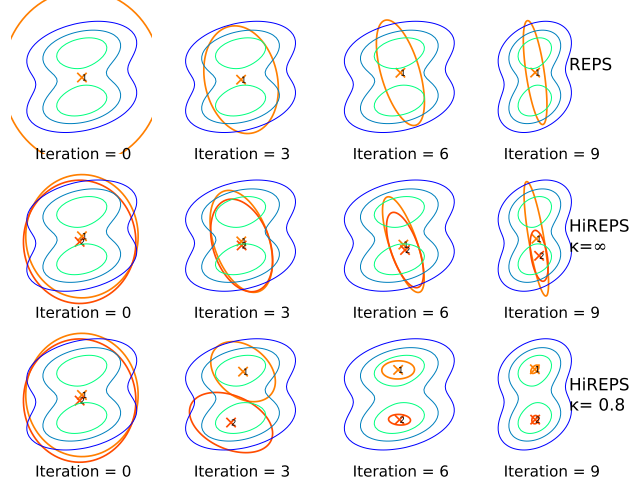


Fig. 2.5 Comparison REPS and HiREPS with and without bounding the overlap of the options on a simple bimodal reward function. The standard REPS approach only uses one option which averages over both modes. The HiREPS approach can use multiple options (two are shown). However, if we do not bound the overlap ($\kappa = \infty$), the options do not separate and concentrate on both modes. Only if we use the overlap constraint, we get a clear separation of the policies.

Lower Bound for the Optimization Problem with Latent Variables. Putting together the new constraints, HiREPS is defined as the following optimization problem:

$$\begin{aligned}
 & \max_p \sum_o \iint p(\mathbf{s}, \boldsymbol{\theta}, o) \mathcal{R}_{\mathbf{s}\boldsymbol{\theta}} d\mathbf{s} d\boldsymbol{\theta}, \\
 \text{s. t. } & \epsilon \geq \sum_o \iint p(\mathbf{s}, \boldsymbol{\theta}, o) \log \left(\frac{p(\mathbf{s}, \boldsymbol{\theta}, o)}{q(\mathbf{s}, \boldsymbol{\theta}) p(o|\mathbf{s}, \boldsymbol{\theta})} \right) d\mathbf{s} d\boldsymbol{\theta}, \quad (2.72) \\
 & \hat{\phi} = \iint p(\mathbf{s}, \boldsymbol{\theta}) \phi(\mathbf{s}) d\boldsymbol{\theta} d\mathbf{s}, \quad 1 = \iint p(\mathbf{s}, \boldsymbol{\theta}) d\boldsymbol{\theta} d\mathbf{s}, \\
 & \kappa \hat{H}_q \geq \mathbb{E}_{\mathbf{s}, \boldsymbol{\theta}} [H(p(o|\mathbf{s}, \boldsymbol{\theta}))], \\
 & 1 = \sum_o \iint p(\mathbf{s}, \boldsymbol{\theta}, o) d\mathbf{s} d\boldsymbol{\theta}. \quad (2.73)
 \end{aligned}$$

Unfortunately, this optimization problem can not be solved in closed form as it contains the conditional distribution $p(o|\mathbf{s}, \boldsymbol{\theta})$ inside the log-

arithm. However, a lower bound of this optimization problem can be obtained using an EM-like update procedure [17]. In the E-step, we estimate

$$\tilde{p}(o|\mathbf{s}, \boldsymbol{\theta}) = \frac{p(\mathbf{s}, \boldsymbol{\theta}, o)}{\sum_o p(\mathbf{s}, \boldsymbol{\theta}, o)}.$$

In the M-step, we use $\tilde{p}(o|\mathbf{s}, \boldsymbol{\theta})$ for $p(o|\mathbf{s}, \boldsymbol{\theta})$ in the optimization problem, and, therefore, neglect the relationship between $p(o|\mathbf{s}, \boldsymbol{\theta})$ and the joint distribution $p(\mathbf{s}, \boldsymbol{\theta}, o)$. Similar to EM, it can be shown that this iterative optimization procedure maximizes a lower bound of the original optimization problem that is tight after each E-step [17]. The resulting joint distribution has the following solution

$$p(\mathbf{s}, \boldsymbol{\theta}, o) \propto q(\mathbf{s}, \boldsymbol{\theta}) \tilde{p}(o|\mathbf{s}, \boldsymbol{\theta})^{1+\kappa/\eta} \exp\left(\frac{\mathcal{R}_{\mathbf{s}\boldsymbol{\theta}} - \varphi(\mathbf{s})^T \mathbf{v}}{\eta}\right), \quad (2.74)$$

where κ denotes the Lagrangian multiplier from bounding the overlap of the options, see Equation (2.71).

The dual function $g(\eta, \mathbf{v})$ that is needed to obtain the parameters η and \mathbf{v} , is given in the appendix C. As described in the previous section, the dual-function is approximated with samples and the probabilities $p(\mathbf{s}, \boldsymbol{\theta}, o)$ are only known for the given set of sample. Hence, we need to fit parametric models to the gating policy as well as to the option policies. Simple linear Gaussian models were used to represent the option policies $\pi(\boldsymbol{\theta}|o, \mathbf{s})$ and the gating policy was also given by a Gaussian gating. The episode-based HiREPS algorithm is summarized in Algorithm 13.

The advantage of using weighted maximum likelihood policy updates for determining hierarchical policies has yet still to be explored for more complex hierarchies. Exploiting structures such as hierarchies might well be the missing key to scale robot learning to more complex real world environments.

2.4.3.4 Step-based REPS for Infinite Horizon Problems

The original REPS formulation [57] is step-based and uses an infinite horizon formulation. The step-based algorithm uses the KL divergence $\text{KL}(p(\mathbf{x}, \mathbf{u})||q(\mathbf{x}, \mathbf{u}))$ on the resulting distribution over the state-action

Algorithm 13 Episode-Based HiREPS for Multiple Contexts

Input: KL-bounding ϵ , overlap-bounding κ

$$\text{data-set } \mathcal{D}_{\text{ep}} = \left\{ \mathbf{s}^{[i]}, \boldsymbol{\theta}^{[i]}, R^{[i]} \right\}_{j=1 \dots N}$$

old gating: $q(o|\mathbf{s})$, old option-policies $q(\boldsymbol{\theta}|\mathbf{s}, o)$ **Compute** $p(o|\mathbf{s}, \boldsymbol{\theta})$ for all options and samples i

$$\tilde{p}(o|i) = \frac{q(\boldsymbol{\theta}^{[i]}|\mathbf{s}^{[i]}, o)q(o|\mathbf{s}^{[i]})}{\sum_{o'} q(\boldsymbol{\theta}^{[i]}|\mathbf{s}^{[i]}, o')q(o'|\mathbf{s}^{[i]})}$$

Optimize dual-function, see Equation (4.20)

$$[\eta, \mathbf{v}] = \operatorname{argmin}_{\eta', \mathbf{v}'} g(\eta', \mathbf{v}'), \quad \text{s.t. } \eta > 0$$

Obtain option policies $\pi_{\boldsymbol{\omega}}^O(\boldsymbol{\theta}|\mathbf{s}, o)$ for all options o

$$\boldsymbol{\omega}_{\text{new}}^o = \operatorname{argmax}_{\boldsymbol{\omega}} \sum_{i=1}^N d_o^{[i]} \log \pi_{\boldsymbol{\omega}}^O(\boldsymbol{\theta}^{[i]}|\mathbf{s}^{[i]}, o)$$

Obtain gating policy $\pi_{\boldsymbol{\omega}}^G(o|\mathbf{s})$ by weighted ML estimate

$$\boldsymbol{\omega}_{\text{new}}^G = \operatorname{argmax}_{\boldsymbol{\omega}} \sum_{i=1}^N \sum_o d_o^{[i]} \log \pi_{\boldsymbol{\omega}}^G(o|\mathbf{s}^{[i]})$$

$$\text{with } d_o^{[i]} = \tilde{p}(o|i)^{1+\xi/\eta} \exp\left(\frac{R^{[i]} - \mathbf{v}^T \boldsymbol{\varphi}(\mathbf{s}^{[i]})}{\eta}\right)$$

pairs $p(\mathbf{x}, \mathbf{u})$, as a similarity measure of the new trajectory distribution $p(\boldsymbol{\tau})$ and the old trajectory distribution $q(\boldsymbol{\tau})$.

In the infinite horizon formulation, REPS maximizes the average reward per time step, given as

$$J_{\pi, \mu^\pi} = \mathbb{E}[r(\mathbf{x}, \mathbf{u})] = \iint \mu^\pi(\mathbf{x}) \pi(\mathbf{u}|\mathbf{x}) r(\mathbf{x}, \mathbf{u}) d\mathbf{x} d\mathbf{u}. \quad (2.75)$$

The distribution $\mu^\pi(\mathbf{x})$ denotes the stationary state distribution of the MDP with policy π .

Stationary State Distributions. The stationary state distribution $\mu^\pi(\mathbf{x})$ represents the probability of visiting state \mathbf{x} when following pol-

icy π . It can not be chosen freely but has to comply with the given state dynamics and the policy. Therefore, it has to fulfill the following constraint

$$\forall \mathbf{x}' : \mu^\pi(\mathbf{x}') = \iint_{\mathbf{x}, \mathbf{u}} \mu^\pi(\mathbf{x}) \pi(\mathbf{u}|\mathbf{x}) p(\mathbf{x}'|\mathbf{x}, \mathbf{u}) d\mathbf{x} d\mathbf{u}. \quad (2.76)$$

As these constraints are not feasible for continuous state spaces, we can again require that the distributions only match on their expected state-features $\varphi(\mathbf{x})$, i.e. the expected feature from the distribution $\mu^\pi(\mathbf{x}')$ need to match the expected features of the distribution

$$\tilde{\mu}^\pi(\mathbf{x}') = \iint \mu^\pi(\mathbf{x}) \pi(\mathbf{u}|\mathbf{x}) p(\mathbf{x}'|\mathbf{x}, \mathbf{u}) d\mathbf{x} d\mathbf{u},$$

which corresponds to $\mu^\pi(\mathbf{x}')$ after applying the policy and the system dynamics. Such a constraint can be formalized as

$$\int \mu^\pi(\mathbf{x}') \varphi(\mathbf{x}') d\mathbf{x}' = \iiint \mu^\pi(\mathbf{x}) \pi(\mathbf{u}|\mathbf{x}) p(\mathbf{x}'|\mathbf{x}, \mathbf{u}) \varphi(\mathbf{x}') d\mathbf{x} d\mathbf{u} d\mathbf{x}'. \quad (2.77)$$

Closeness-to-the-Data Constraint. To ensure the closeness to the old distribution, we bound the relative entropy between the old state-action distribution $q(\mathbf{x}, \mathbf{u})$ and the new state action distribution $\mu^\pi(\mathbf{x}) \pi_\theta(\mathbf{u}|\mathbf{x})$, i.e.,

$$\begin{aligned} \epsilon &\geq \text{KL}(\mu^\pi(\mathbf{x}) \pi(\mathbf{u}|\mathbf{x}) || q(\mathbf{x}, \mathbf{u})) \\ &= \iint \mu^\pi(\mathbf{x}) \pi(\mathbf{u}|\mathbf{x}) \log \frac{\mu^\pi(\mathbf{x}) \pi(\mathbf{u}|\mathbf{x})}{q(\mathbf{x}, \mathbf{u})} d\mathbf{x} d\mathbf{u}. \end{aligned} \quad (2.78)$$

This bound again limits the loss of information and ensures a smooth learning progress.

Resulting Optimization Program. The resulting optimization program can be formalized as follows:

$$\max_{\pi, \mu^\pi} \iint \mu^\pi(\mathbf{x}) \pi(\mathbf{u}|\mathbf{x}) r(\mathbf{x}, \mathbf{u}) d\mathbf{x} d\mathbf{u}, \quad (2.79)$$

$$\begin{aligned}
\text{s. t. : } \quad & \epsilon \geq \sum_{\mathbf{x}, \mathbf{u}} \mu^\pi(\mathbf{x}) \pi(\mathbf{u}|\mathbf{x}) \log \frac{\mu^\pi(\mathbf{x}) \pi(\mathbf{u}|\mathbf{x})}{q(\mathbf{x}, \mathbf{u})} d\mathbf{x} d\mathbf{u}, \\
& \int \mu^\pi(\mathbf{x}') \boldsymbol{\varphi}(\mathbf{x}') d\mathbf{x}' = \iiint \mu^\pi(\mathbf{x}) \pi(\mathbf{u}|\mathbf{x}) p(\mathbf{x}'|\mathbf{x}, \mathbf{u}) \boldsymbol{\varphi}(\mathbf{x}') d\mathbf{x} d\mathbf{u} d\mathbf{x}', \\
& 1 = \iint \mu^\pi(\mathbf{x}) \pi(\mathbf{u}|\mathbf{x}) d\mathbf{x} d\mathbf{u},
\end{aligned}$$

where the last constraint ensures that $\mu^\pi(\mathbf{x})\pi(\mathbf{u}|\mathbf{x})$ is a normalized probability distribution. This constrained optimization problem can again be solved efficiently by the method of Lagrangian multipliers. Please refer to Appendix C for more details. From the Lagrangian, we can also obtain a closed-form solution for the state-action distribution $\mu^\pi(\mathbf{x})\pi(\mathbf{u}|\mathbf{x})$ which is given as

$$\mu^\pi(\mathbf{x})\pi(\mathbf{u}|\mathbf{x}) \propto q(\mathbf{x}, \mathbf{u}) \exp\left(\frac{r(\mathbf{x}, \mathbf{u}) + \mathbb{E}_{\mathbf{x}'}[V(\mathbf{x}')] - V(\mathbf{x})}{\eta}\right) \quad (2.80)$$

for the joint distribution, and as

$$\pi(\mathbf{u}|\mathbf{x}) \propto q(\mathbf{u}|\mathbf{x}) \exp\left(\frac{r(\mathbf{x}, \mathbf{u}) + \mathbb{E}_{\mathbf{x}'}[V(\mathbf{x}')] - V(\mathbf{x})}{\eta}\right) \quad (2.81)$$

for the policy $\pi(\mathbf{u}|\mathbf{x})$. The parameter η denotes the Lagrangian multiplier for the relative entropy bound and the the function $V(\mathbf{x}) = \boldsymbol{\varphi}^T(\mathbf{x})\mathbf{v}$ includes the Lagrangian multipliers \mathbf{v} for the stationary distribution constraint from Equation (2.77).

The Lagrangian parameters can be efficiently obtained by minimizing the dual function $g(\eta, \mathbf{v})$ of the optimization problem. Intuitively, $V(\mathbf{x})$ can be seen as a value-function. As we can see, the expected value $\mathbb{E}_{p(\mathbf{x}'|\mathbf{x}, \mathbf{u})}[V(\mathbf{x}')]$ of the next state is added to the reward while the current value $V(\mathbf{x})$ is subtracted. With this interpretation, we can also interpret the term $\delta_V(\mathbf{x}, \mathbf{u}) = r(\mathbf{x}, \mathbf{u}) + \mathbb{E}_{p(\mathbf{x}'|\mathbf{x}, \mathbf{u})}[V(\mathbf{x}')] - V(\mathbf{x})$ as advantage function of state-action pair (\mathbf{x}, \mathbf{u}) . Hence, we now use the advantage function to determine the exponential weighting of the state-action pairs. The function $V(\mathbf{x})$ is highly connected to the policy gradient baseline. However, the REPS baseline directly emerged out of the derivation of the algorithm while it has to be added afterwards for the policy gradient algorithms to decrease the variance of the gradient estimate.

The Dual Function. The dual function $g(\eta, \mathbf{V})$ of the step-based REPS optimization problem is given in *log-sum-exp* form

$$g(\eta, \mathbf{V}) = \eta\epsilon + \eta \log \iint q(\mathbf{x}, \mathbf{u}) \exp\left(\frac{\delta_V(\mathbf{x}, \mathbf{u})}{\eta}\right) d\mathbf{x}d\mathbf{u}. \quad (2.82)$$

Due to its *log-sum-exp* form, the dual-function is convex in \mathbf{V} . Furthermore, we can approximate the expectation over $q(x, u)$ with a sum over samples $(x^{[i]}, u^{[i]})$ from the distribution $q(x, u)$, i.e.,

$$g(\eta, \mathbf{V}) \approx \eta\epsilon + \eta \log \frac{1}{N} \sum_{\mathbf{x}^{[i]}, \mathbf{u}^{[i]}} \exp\left(\frac{\delta_V(\mathbf{x}^{[i]}, \mathbf{u}^{[i]})}{\eta}\right). \quad (2.83)$$

Consequently, we do not need to know the distribution $q(\mathbf{x}, \mathbf{u})$ as a function, but only need to be able to sample from it. The parameters η and \mathbf{v} are obtained by minimizing the dual function. As η results from an inequality constraint, η needs to be larger than zero [14]. Solving the dual optimization problem is therefore given by the following program $[\eta, \mathbf{v}] = \operatorname{argmin}_{\eta', \mathbf{v}'} g(\eta', \mathbf{v}')$, s.t.: $\eta' > 0$.

Estimating the New Policy. To deal with continuous actions, we need to use a parametric policy $\pi_{\theta}(\mathbf{u}|\mathbf{x})$. Similar as to the episode-based algorithm, we can compute the probabilities $\mu^{\pi}(\mathbf{x}^{[i]})\pi(\mathbf{u}^{[i]}|\mathbf{x}^{[i]})$ only for the given set of samples and subsequently fit a parametric distribution $\pi_{\theta}(\mathbf{u}|\mathbf{x})$ to these samples. Fitting the policy corresponds to a weighted maximum likelihood estimate where the weighting is given by $d^{[i]} = \exp(\delta_V(\mathbf{x}^{[i]}, \mathbf{u}^{[i]})/\eta)$.

From this result, we can observe the close relationship with step-based EM-based policy search algorithms. For general reward functions, EM-based algorithms use an exponential transformation of the expected future return, i.e., $d^{[i]} = \exp(\beta Q^{\pi}(\mathbf{x}^{[i]}, \mathbf{u}^{[i]}))$, where the inverse temperature β has to be chosen by the user. In contrast, REPS always returns the optimized temperature η of the exponential weighting which exactly corresponds to the desired KL-bound. Note that the scaling η will be different for each policy update depending on the distribution of the current reward samples. Furthermore, REPS uses the value-function as a baseline to account for different achievable values in different states.

Representing the Transition Dynamics. The step-based REPS formulation requires the estimation of $\mathbb{E}_{p(\mathbf{x}'|\mathbf{x},\mathbf{u})}[V(\mathbf{x}')]$, which would require the knowledge of the model $p(\mathbf{x}'|\mathbf{x},\mathbf{u})$. However, in the current implementation [57], a single sample $\mathbf{x}'^{[i]}$ from the observed transition is used to approximate $\mathbb{E}_{p(\mathbf{x}'|\mathbf{x},\mathbf{u})}[V(\mathbf{x}')]$, and hence, no model needs to be known. Such approximation causes a bias as the expectation is not done inside the exponential function which is used for computing $\pi(\mathbf{u}|\mathbf{x})$, and, hence, the policy does not optimize the average reward any more. However, in our experience, this bias has only a minor effect on the quality of the learned solution if the noise in the system is not too high. We summarize the REPS algorithm for the infinite horizon formulation in Algorithm 14. The algorithm computes the expected feature change for each state action pair. However, for continuous states and actions, each state action pair is typically only visited once, and, hence the expectation is approximated using a single sample estimate. Instead of using q as the state-action distribution of the old policy, we can also reuse samples by assuming that $q(\mathbf{x},\mathbf{u})$ is the state-action distribution of the last K policies.

2.4.4 Miscellaneous Important Methods

In this section, we will cover two types of algorithms, stochastic optimization and policy improvements with path integrals, which lead to promising results in robotics. For the stochastic optimization algorithm, we will discuss the Covariance Matrix Adaptation - Evolutionary Strategy (CMA-ES) algorithm while for the path integral approach, we will discuss the Policy Improvements with Path Integral (PI²) algorithm, both of which are well-known in the field of robot learning.

2.4.4.1 Stochastic Optimization

Stochastic optimizers are black-box optimizers, and, hence, can be straightforwardly applied for policy search in the episode-based formulation. As it is typically the case with episode-based algorithms, they model an upper-level policy $\pi_{\omega}(\theta)$ to create samples in the parameter-space, which are subsequently evaluated on the real system.

Algorithm 14 REPS for infinite horizon problems

Input: KL-bounding ϵ data-set $\mathcal{D} = \{\mathbf{x}^{[i]}, \mathbf{u}^{[i]}, r^{[i]}, \mathbf{x}'^{[i]}\}_{i=1\dots N}$ **for** $i = 1 \dots N$ **do**State-action visits: $n(\mathbf{x}^{[i]}, \mathbf{u}^{[i]}) = \sum_j I_{ij}$ Summed reward: $\tilde{r}(\mathbf{x}^{[i]}, \mathbf{u}^{[i]}) = \sum_j I_{ij} r^{[j]}$ Summed features: $\delta\tilde{\varphi}(\mathbf{x}^{[i]}, \mathbf{u}^{[i]}) = \sum_j I_{ij} (\varphi(\mathbf{x}'^{[j]}) - \varphi(\mathbf{x}^{[j]}))$ **end for** (I_{ij} is 1 if $\mathbf{x}^{[i]} = \mathbf{x}^{[j]}$ and $\mathbf{u}^{[i]} = \mathbf{u}^{[j]}$, 0 elsewhere)**Compute sample bellman error**

$$\delta_{\mathbf{v}}^{[i]} = \frac{\tilde{r}(\mathbf{x}^{[i]}, \mathbf{u}^{[i]}) + \mathbf{v}^T \delta\tilde{\varphi}(\mathbf{x}^{[i]}, \mathbf{u}^{[i]})}{n(\mathbf{x}^{[i]}, \mathbf{u}^{[i]})}$$

Optimize dual-function $[\eta, \mathbf{v}] = \operatorname{argmin}_{\eta', \mathbf{v}'} g(\eta', \mathbf{v}')$, s.t. $\eta > 0$

$$g(\eta, \mathbf{v}) = \eta\epsilon + \eta \log \left(\frac{1}{N} \sum_{i=1}^N \exp \left(\frac{\delta_{\mathbf{v}}^{[i]}}{\eta} \right) \right)$$

Obtain parametric policy $\pi_{\theta}(\mathbf{u}|\mathbf{x})$ by weighted ML estimate

$$\theta_{k+1} = \operatorname{argmax}_{\theta} \sum_{i=1}^N \exp \left(\frac{\delta_{\mathbf{v}}^{[i]}}{\eta} \right) \log \pi_{\theta}(\mathbf{u}^{[i]} | \mathbf{x}^{[i]})$$

The Covariance Matrix Adaptation - Evolutionary Strategy.

The Covariance Matrix Adaptation - Evolutionary Strategy (CMA-ES) is considered as the state of the art in stochastic optimization [28]. CMA-ES was applied to policy search in robotics [30] and yielded promising results on standard benchmark tasks such as a cart-pole balancing task with two poles. The procedure of CMA-ES is similar to many episode-based policy search methods such as episode-based EM or episode-based REPS. CMA-ES maintains a Gaussian distribution $\pi_{\omega}(\theta)$ over the parameter vector θ , and uses the data-set \mathcal{D}_{ep} for the policy updates. Similar to the EM-based approaches, CMA-ES also uses a weight $d^{[i]}$ for each sample. However, for estimating the weight $d^{[i]}$ and updating the distribution $\pi_{\omega}(\theta)$ CMA-ES uses heuristics, which

often work well in practice but are not founded on a theoretical basis. For estimating the weight $d^{[i]}$, CMA-ES first sorts the samples $\theta^{[i]}$ according to their return $R^{[i]}$, and, subsequently, computes the weight of the best l samples by $d^{[i]} = \log(l+1) - \log(i)$. All other samples are neglected, i.e., get zero weight. Similar to weighted ML updates, the new mean μ_k of policy $\pi_\omega(\theta)$ is computed by the weighted average of the data points. However, the update of the covariance matrix Σ is based on a combination of the weighted sample-covariance and information about the ‘evolution path’ $\{\mu_j\}_{j=0\dots k}$.

The advantage of such an approach is that the covariance matrix update depends only on the current set of samples and, hence, requires only a few samples $\theta^{[i]}$. The number of samples N to evaluate for CMA-ES is typically fixed to $\max(4 + 3 \log D, 5)$, where D is the dimensionality of θ and the number of samples l used for the weighting is typically set to $N/2$. While CMA-ES is a black-box optimizer and, therefore, simple to use, it also has severe disadvantages. It cannot be used for generalizing the upper-level policy to multiple contexts. Furthermore, several roll-outs have to be evaluated if the evaluation $R^{[i]}$ is noisy. The minimum number of required roll-outs for a given parameter $\theta^{[i]}$ can be computed by Hoeffding and Bernstein races [29], which can slightly alleviate this problem. For a more detailed discussion about the CMA-ES updates we refer to [28].

2.4.4.2 Policy Improvement by Path Integrals

The Policy Improvements by Path Integrals (PI²) algorithm [81] is based on the path integral approach to optimal control. The path integral approach is designed for obtaining optimal control laws for non-linear continuous time systems of the form

$$\dot{\mathbf{x}}_t = \mathbf{f}(\mathbf{x}_t) + \mathbf{G}(\mathbf{x}_t)(\mathbf{u}_t + \mathbf{C}_u \epsilon_t) = \mathbf{f}_t + \mathbf{G}_t(\mathbf{u}_t + \epsilon_t), \quad (2.84)$$

where \mathbf{f}_t denotes a drift term, \mathbf{G}_t the control matrix of the system, ϵ_t is zero mean Brownian motion and \mathbf{C}_u is the diffusion coefficient. Note that these assumptions do not limit the generality of the approach as all physical systems linearly depend on the control action \mathbf{u}_t . Furthermore, the path integral theory assumes squared control costs of the form

$-\mathbf{u}_t \mathbf{R} \mathbf{u}_t$. The state-dependent part of the reward $r_t(\mathbf{x}_t)$ can be an arbitrary function. As path integrals are based on stochastic optimal control theory [83], we will now briefly review the relevant concepts. The stochastic optimal control (SOC) problem is now defined as finding the controls $\mathbf{u}_{1:T}$ which maximize the expected return

$$J(\mathbf{x}_1, \mathbf{u}_{1:T}) = r_T(\mathbf{x}_T) + \int_{t=0}^T r_t(\mathbf{x}_t, \mathbf{u}_t) dt. \quad (2.85)$$

The discrete time-formulation of the system for a fixed time step dt is given as

$$\mathbf{x}_{t+dt} = \mathbf{x}_t + \mathbf{f}_t dt + \mathbf{G}_t \left(\mathbf{u}_t dt + \mathbf{C}_u \boldsymbol{\epsilon}_t \sqrt{dt} \right), \quad (2.86)$$

where the term \sqrt{dt} appears because the variance of Brownian motion grows linearly with time, and, thus, the standard deviation grows with \sqrt{dt} . The probability of the next state given the action and the previous state can also be written down as Gaussian distribution

$$p(\mathbf{x}_{t+dt} | \mathbf{x}_t, \mathbf{u}_t) = \mathcal{N}(\mathbf{x}_t + \mathbf{f}_t dt + \mathbf{G}_t \mathbf{u}_t dt, \mathbf{G}_t \boldsymbol{\Sigma}_u \mathbf{G}_t^T dt) \quad (2.87)$$

with $\boldsymbol{\Sigma}_u = \mathbf{C}_u \mathbf{C}_u^T$. The expected return for the discrete time formulation is given as

$$J(\mathbf{x}_1, \mathbf{u}_{1:T}) = r_T(\mathbf{x}_T) + dt \left(\sum_{t=0}^T r_t(\mathbf{x}_t) - \mathbf{u}_t \mathbf{R}_t \mathbf{u}_t \right). \quad (2.88)$$

The discrete time formulations are needed for the derivation of the Hamilton-Jacobi Bellman (HJB) equation.

Hamilton-Jacobi Bellman Equation. The HJB Equation states the optimality conditions for a value function for continuous time systems as the one given in Equation (2.84). We start the derivation of the continuous time Bellman equation by with the discrete time system, and, after stating the optimality conditions for the discrete time system we will take the limit of $dt \rightarrow 0$ to get the continuous time formulation. The Bellman Equation for the discrete-time problem is given by

$$V(\mathbf{x}, t) = r_t(\mathbf{x}_t) dt + \max_{\mathbf{u}} \left(-\mathbf{u}_t \mathbf{R}_t \mathbf{u}_t dt + \mathbb{E}_{\mathbf{x}_{t+dt}} [V(\mathbf{x}_{t+dt}, t + dt)] \right),$$

where the expectation is done with respect to system dynamics $p(\mathbf{x}_{t+dt}|\mathbf{x}_t, \mathbf{u}_t)$. The HJB equation is now derived by using a second order Taylor approximation of the value function for time step $t + dt$,

$$V(\mathbf{x} + \delta\mathbf{x}, t + dt) \approx V(\mathbf{x}, t + dt) + \delta\mathbf{x}^T \mathbf{v}_x + \frac{1}{2} \delta\mathbf{x}^T \mathbf{V}_{xx} \delta\mathbf{x}, \quad (2.89)$$

with $\mathbf{v}_x = \partial V(\mathbf{x}, t + dt)/\partial \mathbf{x}$ and $\mathbf{V}_{xx} = \partial^2 V(\mathbf{x}, t + dt)/\partial^2 \mathbf{x}$. As $V(\mathbf{x}, t + dt)$ is now represented as quadratic function and $p(\mathbf{x}_{t+dt}|\mathbf{x}_t, \mathbf{u}_t)$ is Gaussian, we can solve the expectation over the next state analytically

$$\begin{aligned} \mathbb{E}_{\mathbf{x}_{t+dt}}[V(\mathbf{x}_{t+dt}, t + dt)] &= V(\mathbf{x}, t + dt) + (\mathbf{f}_t + \mathbf{G}_t \mathbf{u}_t)^T \mathbf{v}_x dt \\ &\quad + \frac{1}{2} \text{tr}(\boldsymbol{\Sigma}_x \mathbf{V}_{xx} dt) + O(dt^2) \end{aligned} \quad (2.90)$$

with $\boldsymbol{\Sigma}_x = \mathbf{G}_t \boldsymbol{\Sigma}_u \mathbf{G}_t^T$. Note that, the second order derivative \mathbf{V}_{xx} also appears in the first order approximation of $\mathbb{E}_{\mathbf{x}_{t+dt}}[V_{t+dt}(\mathbf{x}_{t+dt})]$. This term is a recurrent theme in stochastic calculus and directly relates to the Ito lemma [83]. Setting Equation (2.90) back into Equation (2.4.4.2), we get the following relationship

$$\begin{aligned} V(\mathbf{x}, t) &= V(\mathbf{x}, t + dt) + r_t(\mathbf{x}_t) dt \\ &\quad + \max_{\mathbf{u}} \left(-\mathbf{u}_t^T \mathbf{R}_t \mathbf{u}_t dt + (\mathbf{f}_t + \mathbf{G}_t \mathbf{u}_t)^T \mathbf{v}_x dt + \frac{1}{2} \text{tr}(\boldsymbol{\Sigma}_x \mathbf{V}_{xx} dt) \right). \end{aligned}$$

We can now move the $V(\mathbf{x}, t + dt)$ term to the other side, divide by dt and take the limit of $dt \rightarrow 0$ to end up with the continuous time optimality condition for the value function, also called Hamilton-Jacobi Bellman (HJB) equation,

$$\begin{aligned} -\dot{V}(\mathbf{x}, t) &= \lim_{dt \rightarrow 0} \frac{V(\mathbf{x}, t) - V(\mathbf{x}, t + dt)}{dt} \\ &= r_t(\mathbf{x}_t) + \max_{\mathbf{u}} \left(-\mathbf{u}_t^T \mathbf{R}_t \mathbf{u}_t + (\mathbf{f}_t + \mathbf{G}_t \mathbf{u}_t)^T \mathbf{v}_x + \frac{1}{2} \text{tr}(\boldsymbol{\Sigma}_x \mathbf{V}_{xx}) \right). \end{aligned}$$

As the HJB contains only quadratic terms of \mathbf{u}_t , we can obtain the optimal controls by setting the derivative with respect to \mathbf{u}_t to zero, i.e.,

$$\mathbf{u}_t = \mathbf{R}^{-1} \mathbf{G}_t^T \mathbf{v}_x. \quad (2.91)$$

Setting the optimal controls back into the HJB equation yields

$$-\dot{V}(\mathbf{x}, t) = r_t(\mathbf{x}) + \mathbf{v}_x^T \mathbf{f}_t + \mathbf{v}_x^T \mathbf{G}_t \mathbf{R}^{-1} \mathbf{G}_t^T \mathbf{v}_x + \frac{1}{2} \text{tr}(V_{\mathbf{x}\mathbf{x}} \Sigma_x), \quad (2.92)$$

the partial differential equation, which has to be satisfied by the optimal value function for the system in Equation (2.84).

Path Integrals. The HJB equation is now transformed to a system of linear partial differential equations by performing an exponential transformation of the optimal value function $\Psi(\mathbf{x}, t) = \exp(V(\mathbf{x}, t)/\lambda)$, where λ can be seen as the temperature of the exponential transformation. Under the assumption that the quadratic control cost matrix \mathbf{R} is given by the system noise, i.e., $\mathbf{R} = \lambda \Sigma_u^{-1}$, the solution for Ψ can be obtained by applying the Feynman-Kac theorem [55] for solving partial differential equations. The solution is given by

$$\Psi(\mathbf{x}, t) = \int p_{\text{uc}}(\boldsymbol{\tau}|\mathbf{x}, t) \exp\left(\frac{\sum_{h=t}^T r_h(\mathbf{x}_h) dt + r_T(\mathbf{x}_T)}{\lambda}\right) d\boldsymbol{\tau}, \quad (2.93)$$

where $p_{\text{uc}}(\boldsymbol{\tau}|\mathbf{x}, t)$ is the probability of the process using the uncontrolled dynamics $\dot{\mathbf{x}}_t = \mathbf{f}_t + \mathbf{G}_t(\boldsymbol{\epsilon}_t)$, i.e., $\mathbf{u}_t = 0$, when starting in state \mathbf{x} at time step t [81]. Note that Equation (2.93) is given in its discrete-time form, wherein the discretization time step is again denoted as dt . For more details on the exponential transformation of the value function and the Feynman-Kac theorem we refer to [81]. From the assumption $\mathbf{R} = \lambda \Sigma_u^{-1}$, we can also conclude that λ specifies the control costs. The higher we choose the temperature λ of the exponential transformation, the less greedy the exponential transformation will become. This intuition is also reflected in the increasing control costs with increasing λ , and, consequently, the resulting solution will become more similar to the uncontrolled process. In practice, the assumption for the control cost matrix \mathbf{R} is rather limiting, as we are not free in our choice of the control costs.

We will further define $S(\boldsymbol{\tau}|\mathbf{x}_t, t)$ to be the path integral of trajectory $\boldsymbol{\tau}$ starting at time step t in state \mathbf{x} which is given as

$$S(\boldsymbol{\tau}|\mathbf{x}_t, t) = r_T(\mathbf{x}_T) + \sum_{h=t}^T r_h(\mathbf{x}_h) dt + \log p_{\text{uc}}(\boldsymbol{\tau}|\mathbf{x}, t). \quad (2.94)$$

Hence, the path integral of a trajectory for time step t is given by the reward to come plus a logarithmic punishment term which renders less likely trajectories less attractive. Due to the assumption of $\mathbf{R} = \lambda \Sigma_u^{-1}$, the return and the probability of a trajectory can be treated in a unified way. The optimal controls \mathbf{u}_t can be obtained by

$$\mathbf{u}_t = -\mathbf{R}^{-1} \mathbf{G}_t^T \mathbf{v}_x = \lambda \mathbf{R}^{-1} \mathbf{G}_t^T \frac{\partial / \partial \mathbf{x} \Psi(\mathbf{x}, t)}{\Psi(\mathbf{x}, t)}. \quad (2.95)$$

Determining the term $(\partial / \partial \mathbf{x} \Psi(\mathbf{x}, t)) / \Psi(\mathbf{x}, t)$ and setting this term into Equation (2.95), yields [81]

$$\mathbf{u}_t = \int p_{\text{co}}(\boldsymbol{\tau} | \mathbf{x}_t, t) \mathbf{u}_L(\boldsymbol{\tau}, t) d\boldsymbol{\tau}, \quad (2.96)$$

where

$$p_{\text{co}}(\boldsymbol{\tau} | \mathbf{x}_t, t) = \frac{\exp(S(\boldsymbol{\tau}_t | \mathbf{x}_t, t) / \lambda)}{\int \exp(S(\boldsymbol{\tau}_t | \mathbf{x}_t, t) / \lambda) d\boldsymbol{\tau}} \quad (2.97)$$

$$\propto p_{\text{uc}}(\boldsymbol{\tau}_t | \mathbf{x}_t, t) \exp\left(\sum_{h=t}^T r_h(\mathbf{x}_h) dt + r_T(\mathbf{x}_T)\right) \quad (2.98)$$

and

$$\mathbf{u}_L(\boldsymbol{\tau}_t) = \mathbf{R}^{-1} \mathbf{G}_t (\mathbf{G}_t \mathbf{R}^{-1} \mathbf{G}_t^T)^{-1} \mathbf{G}_t \boldsymbol{\epsilon}_t. \quad (2.99)$$

We will denote the distribution $p_{\text{co}}(\boldsymbol{\tau} | \mathbf{x}_t, t)$ as controlled process distribution, as it denotes the trajectory distribution connected to the optimal (transformed) value function Ψ_t . We observe that the controlled process distribution is represented as a soft-max distribution which has the path integrals $S(\boldsymbol{\tau} | \mathbf{x}, t)$ of the trajectory in its exponent. Alternatively, $p_{\text{co}}(\boldsymbol{\tau} | \mathbf{x}_t, t)$ can also be written as the uncontrolled process distribution $p_{\text{uc}}(\boldsymbol{\tau} | \mathbf{x}_t, t)$ that is weighted by the exponentially transformed reward to come.

The action $\mathbf{u}_L(\boldsymbol{\tau}, t)$ is denoted as correction action of trajectory $\boldsymbol{\tau}$. It is defined as the action which follows the trajectory $\boldsymbol{\tau}$ while minimizing the immediate control costs [81] at time step t . The term $\boldsymbol{\epsilon}_t$ is the noise term applied at time step t for trajectory $\boldsymbol{\tau}$. The matrix $\mathbf{R}^{-1} \mathbf{G}_t (\mathbf{G}_t \mathbf{R}^{-1} \mathbf{G}_t^T)^{-1} \mathbf{G}_t$ projects the applied noise term into the null-space of the control matrix \mathbf{G}_t , and, hence, eliminates the unnecessary part of the control action \mathbf{u} .

By combining Equations (2.99) and (2.96), we can summarize that the optimal control law is given in the form of a soft-max distribution $p_{\text{co}}(\boldsymbol{\tau}|\boldsymbol{x}, t)$, which weights relevant parts of the noise term $\boldsymbol{\epsilon}_t$ according to their path integrals $S(\boldsymbol{\tau}, t)$. The main advantage of path integrals is that the optimal action can be obtained by performing Monte-Carlo roll-outs instead of applying dynamic programming. As in the REPS approach, the maximum-operation, which is typically needed to obtain the optimality of the action, is replaced by a soft-max operator, which is easier to perform. In the following, the condition on the start state \boldsymbol{x} of the trajectories will be dropped in order to make the computations feasible. However, this simplification might again add a bias to the resulting path integral approach. The effect of this bias still has to be examined.

Iterative Path Integral Control. In practice, we have to sample from the uncontrolled process to solve the integral in Equation (2.95) over the trajectory space. However, this sampling process can be inefficient for high-dimensional systems, wherein high number of samples is needed to obtain an accurate estimate. The number of required samples can be reduced by using an iterative approach. At each iteration k , we only compute the optimal change $\delta_k \boldsymbol{u}_t$ in the action for time step t . Subsequently, the mean action $\boldsymbol{u}_{k,t}$ for time step t is updated by $\boldsymbol{u}_{k,t} = \boldsymbol{u}_{k-1,t} + \delta_k \boldsymbol{u}_t$. Such procedure allows for the use of a smaller system noise $\boldsymbol{\Sigma}_{\boldsymbol{u}}$ for exploration, and hence, we can search for a locally optimal improvement of the options. As we now search for an optimal change of the action $\delta \boldsymbol{u}_t$, the current estimate $\boldsymbol{u}_{k,t}$ of the action is subsumed in the drift-term of the dynamics, i.e., $\tilde{\boldsymbol{f}}_{k,t} = \boldsymbol{f}_t + \boldsymbol{G}_t \boldsymbol{u}_{k,t}$. However, the explicit dependence of $\boldsymbol{u}_{k,t}$ from the current state is ignored in this iterative formulation. Note that the path integral approach only estimates the mean control action of the policy and not the variance. Exploration is solely performed by the user-specified uncontrolled dynamics.

Policy Improvements by Path Integrals. The Policy Improvement by Path Integrals (PI²) algorithm [81] applies the path inte-

gral theory to the problem of learning Dynamic Movement Primitives (DMPs) as policy representations. In this section, we restrict ourselves to learning a DMP for a single joint, for multiple joints, learning can be performed by applying the discussed update rules for each joint separately. The path integral theory can be applied directly to DMPs by treating the DMP parameter vector \mathbf{w} as the control action for the dynamical system

$$\ddot{y} = \underbrace{\tau^2 \alpha_y (\beta_y (g - y) - \dot{y})}_{f_t} + \underbrace{\tau^2 \phi_t^T}_{G_t} (\mathbf{w}_t + \boldsymbol{\epsilon}_t), \quad (2.100)$$

which defines the DMP trajectory. Note that, as in the PoWER approach [38], we have assumed to use a different parameter vector \mathbf{w}_t in each time step. The drift term f_t is given by the spring and damper system of the DMP and the control matrix G_t is given by the basis functions ϕ_t of the DMP. Since we apply the exploration noise $\boldsymbol{\epsilon}_t$ at each time step to the parameter vector \mathbf{w}_t , the exploration policy is given by $\mathcal{N}(u_t | f_t, H_t)$ with $H_t = \phi_t^T \boldsymbol{\Sigma}_w \phi_t$.

In order to define the path-integral update rule, we need to compute the path integral $S(\boldsymbol{\tau} | t)$ for a given trajectory. This step requires knowledge of the uncontrolled trajectory distribution $p_{\text{uc}}(\boldsymbol{\tau} | t)$, and, hence, knowledge of the system model. However, if we assume the rewards r_t to depend only on the state of the DMP, and not on the real state of the robot or its environment, i.e., $r_t(\mathbf{x}_t) = r_t(y_t, \dot{y}_t)$, the uncontrolled dynamics $p_{\text{uc}}(\boldsymbol{\tau} | \mathbf{x}_t) = \prod_{l=0}^T p(y_l, \dot{y}_l | y_{l-1}, \dot{y}_{l-1})$ are straightforward to compute. The path integral can be rewritten as

$$S(\boldsymbol{\tau}, t) = r_T(\mathbf{x}_T) + \sum_{l=t}^{T-1} r_l(\mathbf{x}_l) + (\mathbf{w}_l + \mathbf{M}_l \boldsymbol{\epsilon}_l)^T \mathbf{R}(\mathbf{w}_l + \mathbf{M}_l \boldsymbol{\epsilon}_l), \quad (2.101)$$

where $\mathbf{M}_l = \phi_l H_l^{-1} \phi_l^T$ [81]. The new parameter vector $\mathbf{w}_{\text{new},t}$ for time step t is computed by the iterative path integral update rule. Using Equation (2.95) for the optimal control action yields

$$\mathbf{w}_{\text{new},t} = \mathbf{w}_t + \int p_{\text{co}}(\boldsymbol{\tau}_t) \mathbf{M}_t \boldsymbol{\epsilon}_t d\boldsymbol{\tau}, \quad (2.102)$$

where $p_{\text{co}}(\boldsymbol{\tau}_t)$ is given by the soft-max distribution in Equation (2.98). So far, we used a different parameter vector w_t for each time step.

However, in practice, we can use only a single parameter vector \mathbf{w} for one episode, and, thus, we need to average over the parameter updates for all time-steps. The average update is computed by weighting each parameter vector $\mathbf{w}_{\text{new},t}$ with the number of time steps to go. Additionally, the update for the j -th dimension of the parameter vector \mathbf{w} for time step t is weighted by the activation of the j -th feature function $[\phi_t]_j$,

$$\begin{aligned} [\mathbf{w}_{\text{new}}]_j &= \frac{\sum_{t=0}^{T-1} (T-t) [\phi_t]_j [\mathbf{w}_{\text{new},t}]_j}{\sum_{t=0}^{T-1} (T-t) [\phi_t]_j} \\ &\approx [\mathbf{w}_{\text{old}}]_j + \sum_{t=0}^{T-1} \sum_{i=1}^N p_{\text{co}}(\tau_t^{[i]}) d_{t,j} [\boldsymbol{\epsilon}_t^{[i]}]_j \end{aligned}$$

with

$$d_{t,j} = (T-t) [\phi_t]_j / \left(\sum_{t'=0}^{T-1} (T-t') [\phi_{t'}]_j \right). \quad (2.103)$$

Equation (2.103) defines the update rule of the original PI^2 algorithm given in [81]. We observe that PI^2 fits into our categorization of using a step-based policy evaluation strategy, which uses the future reward in the current trajectory plus a punishment term for unlikely trajectories as the evaluation. Similar to the PoWER [38] algorithm, this evaluation is used by the soft-max policy to obtain a weighting for each of the samples.

Episode-Based PI^2 . The PI^2 algorithm has also been used in the episode-based formulation [77], which also has been used for updating the exploration strategy. The basic PI^2 algorithm does not update its exploration strategy and has to rely on the uncontrolled process dynamics, which is typically set by the user. However, the noise-variance of the uncontrolled process dynamics has a large impact on the learning performance and, hence needs to be chosen appropriately. To automatically estimate the exploration strategy, the PI^2 algorithm was reformulated in the episode-based policy search formulation and the condition that the noise covariance $\boldsymbol{\Sigma}_{\mathbf{u}}$ needs to match the control costs matrix \mathbf{R} was explicitly ignored. Due to the episode-based formulation,

the policy can be directly estimated in parameter space. Instead of using the uncontrolled process for exploration, the previously estimated policy is used for exploration. Furthermore, the log-term for the uncontrolled dynamics $p_{\text{uc}}(\boldsymbol{\tau}|\mathbf{x}_t)$ in the path integral $S(\boldsymbol{\tau})$ is neglected as this term does not seem to improve the performance. Consequently, in the episode-based formulation the returns $R^{[i]}$ are directly used as path integrals. The covariance matrix was updated by a weighted maximum likelihood estimate, where the soft-max distribution $p_{\text{co}}(\boldsymbol{\tau}_t)$ was used as weighting. The resulting episode-based PI² is a simplified version of the original PI² algorithm, but shows an improved learning performance.

Relation to Information-Theoretic and EM Approaches. Despite that information theoretic and EM approaches were developed from different principles than the path integral approach, all these approaches share similar characteristics.

The episode-based formulation of PI² shares many similarities with the episode-based REPS formulation. By pulling the $\log p_{\text{uc}}(\boldsymbol{\tau}_t|\mathbf{x}_t)$ term outside the exponent, we realize that $p_{\text{co}}(\boldsymbol{\tau}|\mathbf{x}_t)$ shares a similar soft-max structure as the closed form solution in REPS for $\pi(\mathbf{u}|\mathbf{x})$. In REPS, the trajectory distribution $q(\boldsymbol{\tau})$ of the old policy is used instead of the uncontrolled process distribution $p_{\text{uc}}(\boldsymbol{\tau}_t|\mathbf{x}_t)$. A similar strategy is emulated by using the iterative path integrals update, where the mean of the exploration policy is updated, but the exploration noise is always determined by the uncontrolled process. While REPS is designed to be an iterative algorithm, the iterative sampling process of PI² needs to be motivated from a heuristic view point.

We also observe that the temperature parameter λ in path integral control corresponds to the Lagrangian parameter η in REPS. This parameter is automatically set in REPS according to the relative entropy bound, while for path integral control, λ is set by heuristics. The path integral approach also relies on the assumption that the control cost matrix \mathbf{R} is predefined as $\lambda\mathbf{R}^{-1} = \boldsymbol{\Sigma}_u$. Dropping this assumption results in a more efficient algorithm [77], but the theoretical justifications for such an algorithm are also lost. Similar update rules emerge naturally for the REPS algorithm without the need for heuristics. However, REPS has so far only been introduced for the episode-based policy

search learning formulation, and, hence, makes inefficient usage of the available trajectory samples $\tau^{[i]}$. PI^2 has been derived from a step-based formulation, and, hence, might have advantages over REPS in some applications.

The original step-based version of the PI^2 algorithm is also closely related to the PoWER algorithm. If we also use an exponential transformation of the reward for PoWER, the policy updates are essentially the same. While PoWER uses a weighted maximum likelihood update to obtain the new policy, PI^2 averages the update rules for the single time steps.

2.5 Real Robot Applications with Model-Free Policy Search

In this section, we present selected results for model-free policy search in the area of robotics. These experiments include Baseball, Ball-In-The-Cup, Dart-Throwing, Pan-Cake Flipping and Tetherball. All experiments have been conducted with dynamic movement primitives as policy representation. In all applications, learning with DMPs takes place in two phases [38]. In the first phase, imitation learning is used to reproduce recorded trajectories. Subsequently, reinforcement learning is used to improve upon the imitation. The use of imitation learning to initialize the learning process allows for the incorporation of experts knowledge and can considerably speed up the learning process. Most experiments have been performed with a specific algorithm in mind that was at the time of the experiment state of the art. However, more recent approaches such as REPS or PI^2 would also have worked in most setups⁴.

2.5.1 Learning Baseball with eNAC

In the baseball experiment, a Sarcos Master Arm was used to hit a soft baseball placed on a T-stick such that it flies as far as possible [61]. This game is also called T-Ball and used to teach children how to hit

⁴While REPS also works for contextual policy search, PI^2 was so far not extended to the multi-task setup.

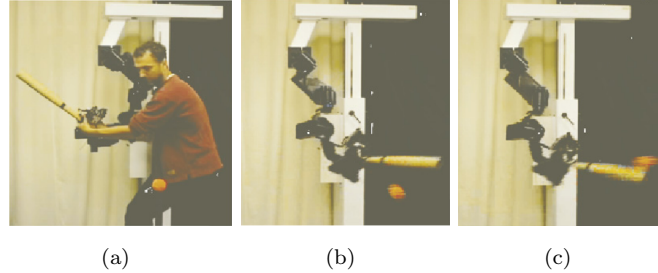


Fig. 2.6 Learning a baseball swinging movement to hit a ball placed on a T-stick [61]. (a) Obtaining an initial solution by imitation learning. (b) Initial solution replayed by the robot. The robot misses the ball. (c) Movement learned with the eNAC algorithm after 300 roll-outs. The robot learned to hit the ball robustly.

a baseball. The robot had seven degrees of freedom (DoF) and ten basis functions were used for each DoF. Only the shape parameters of the DMP were adapted during learning, resulting in a 70-dimensional weight vector \mathbf{w} . The reward function was given by

$$R(\boldsymbol{\tau}) = c_1 p_x - c_2 \sum_{t=0}^{T-1} \ddot{\mathbf{q}}_t^T \ddot{\mathbf{q}}_t,$$

where p_x is the distance the ball traveled and c_1 and c_2 are constants to weight the objective of hitting the ball versus minimizing the energy consumption. An initial solution was obtained by imitation learning, see Figure 2.6(a). However, the robot failed to exactly reproduce the behavior and missed the ball, Figure 2.6(b). The behavior could be improved by employing the episodic Natural Actor Critic algorithm. After 200 to 300 trials of learning, the robot was able to reliably hit the ball, see Figure 2.6(c).

2.5.2 Learning Ball-in-the-Cup with PoWER

The PoWER algorithm was used in [38] to learn the game ‘Ball-in-the-Cup’. The used robot platform was a Barrett WAM robot arm with seven degrees of freedom. In total, the authors selected 31 basis functions to learn the task. The shape parameters of the DMP were learned as well as the variance σ_i^2 for each basis function. All degrees of freedom are perturbed separately but share the reward, which is

zero except for the time step t_c , where the ball passes the cup rim in a downward direction. A stereo vision system was used to track the position $\mathbf{b} = [x_b, y_b, z_b]^T$ of the ball. This ball position was used for determining the reward, but not for feedback during the motion. The reward at time-step t_c was given by

$$r_{t_c} = \exp(-\alpha(x_c - x_b) - \alpha(y_c - y_b)),$$

where x_c and y_c are the x and y -coordinates of the cup and x_b and y_b the coordinates of the ball at this time step. The parameter α is a scaling parameter which is set to 100. The exp function is used to transform the squared distance of the ball to the cup into an improper probability distribution, as PoWER requires this type of reward function. The policy was initialized with imitation learning. After 75 trials, the robot could reliably catch the ball with the cup. The resulting learning progress is illustrated in Figure 2.7, which shows the position of the ball closest to the cup for a different number of roll-outs. This experiment was extended in [35] to include feedback for the position of the ball. Here, a DMP was used to learn the expected trajectory of the ball. For each degree of freedom, an additional PD-controller was learned which is added to the forcing function of the DMP. The PD controller takes the deviation of the current ball position and velocity to the expected ball position and velocity as input. This experiment was only conducted in simulation, the initial position and velocity of the ball were perturbed randomly. The DMP, including the additional feedback terms, contained 91 parameters which were optimized by the PoWER approach. Due to the high variance in the initial state of the ball, the PoWER algorithm now converged after 500 to 600 episodes.

2.5.3 Learning Pan-Cake Flipping with PoWER/RWR

In [40], the PoWER algorithm was used to learn to flip a pan-cake with a Barrett WAM. As underlying policy representation, an extension of the DMP approach was used where the spring-damper system of the DMP was also modeled as time dependent, and, additionally, the

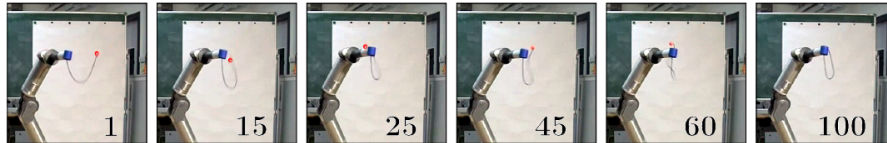


Fig. 2.7 Learning the game ‘Ball-in-the-cup’ [38]. The robot has to learn how to catch the ball which is connected with a string to the cup. The initial policy is learned by imitation and shown on the left. The figures show the minimum distance of the ball to the cup with an increasing number of roll-outs. After 75 episodes learning with the PoWER algorithm, the robot was able to catch the ball reliably.

spring-damper system for each dimension were coupled⁵. The DMP were defined in task space, i.e., in the 3-D Cartesian coordinates of the end-effector position. The coupling terms of the spring damper system as well as the goal attractor were learned. The reward function was composed of three terms: The rotation error of the pan-cake before landing in the pan, the distance of the pan-cake to the center of the pan when catching the pan-cake, and the maximum height the pan-cake reached. The authors state that they used the PoWER algorithm, however, we would rather classify the used algorithm as Reward Weighted Regression as it uses an episode-based evaluation and exploration strategy. The intermediate steps of the episode, a key characteristics of the PoWER algorithm, are neglected. The robot was able to learn the task after 50 roll-outs. A successful pan-cake-flipping episode is indicated in Figure 2.8. The robot also learned how to adjust the stiffness of the joints due to the time-varying spring damper system of the DMP. When catching the pan-cake, the robot increased the compliance of the arm such that the arm can freely move down and prevents the pan-cake from bouncing out of the pan.

2.5.4 Learning Dart Throwing with CRKR

In [37], CRKR was used to learn dart throwing with a robot. The robot had to hit different locations on a dart board. The dart is placed

⁵Coupling is typically achieved with the phase variable in the DMP framework. However, using coupling also for the feedback terms of the spring damper system allows for the use of correlated feedback.

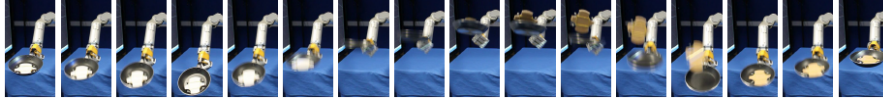


Fig. 2.8 Learning the pan-cake-flipping task [40]. The robot has to flip a pan-cake such that it rotates 180 degrees in the air, and, subsequently, the robot has to catch it again. The figure illustrates the learned movement after 50 roll-outs. The authors applied the RWR algorithm to learn the movement.

on a launcher attached to the end-effector and held there by stiction. CRKR learns an upper-level policy $\pi_{\omega}(\boldsymbol{\theta}|\mathbf{s})$ where the context \mathbf{s} is given by the two-dimensional location of the target on the dart board. The parameters $\boldsymbol{\theta}$ contained the the goal-attractor \mathbf{g} , the final velocity $\dot{\mathbf{g}}$ and the time scale constant τ of the DMP. The shape parameters \mathbf{w} of the DMP were obtained via imitation learning and fixed during learning. The reward function was given by the negative distance of the impact position \mathbf{d} of the dart to the desired position \mathbf{x} and an additional punishment term for fast movements, i.e.,

$$R(\boldsymbol{\theta}, \mathbf{x}) = -10\|\mathbf{d} - \mathbf{x}\| - \tau,$$

where τ is the time-scale constant of the DMP, which controls the velocity of the movement. The experiment was conducted on three real-robot platforms, the Barrett WAM, the humanoid robot CBi and the Kuka CR 6. For all robot platforms, CRKR was able to learn a successful throwing movement and hit the desired target within the range of $\pm 10\text{cm}$ after 200 to 300 trials, which was within the reproduction accuracy of the robots. The learned movement for the CBi is illustrated in Figure 2.9.

2.5.5 Learning Table Tennis with CRKR

CRKR was also used to learn to return table-tennis balls [37]. The authors used a Barrett WAM with seven DoF, which was mounted on the ceiling. The hitting movement for table-tennis was decomposed into three dynamic movement primitives. In the first phase, the robot swings back. The second movement primitive is then used to hit the ball while the third primitive is used to smoothly go back the initial position. The

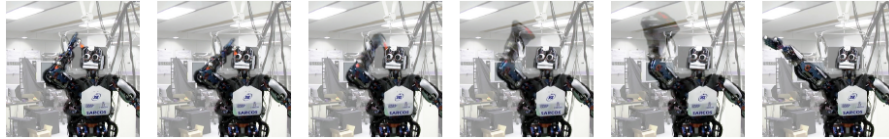


Fig. 2.9 Learning the dart throwing task [37]. CRKR was used to generalize a dart throwing movement to different target locations. After 300 episodes, the humanoid robot CBi was able to reliably throw the dart within a range of $\pm 10\text{cm}$ distance to the target. CRKR adapted the meta-parameters of the DMP which included the goal position \mathbf{g} and goal velocity $\dot{\mathbf{g}}$ for each joint and the time scaling constant of the DMP.

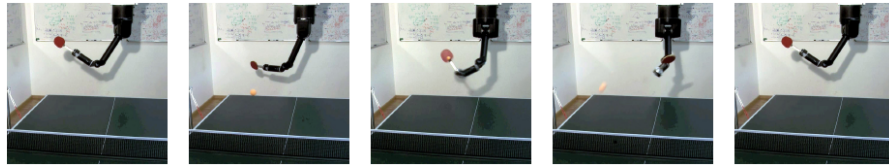


Fig. 2.10 Learning to hit a table tennis ball in the air [37]. CRKR was used to learn an upper-level policy $\pi_{\omega}(\theta|\mathbf{s})$, where the context \mathbf{s} was given by the position and velocity of the ball when it passes the net. The policy parameters θ included the final positions \mathbf{g} and velocities $\dot{\mathbf{g}}$ of the joint as well as a waiting time when to start the hitting movement. The robot was able to hit the ball for different configurations of the incoming ball and improved its success rate for hitting the ball from 30% to 80%.

meta-parameters of the second primitive, including the final position \mathbf{g} and final velocities $\dot{\mathbf{g}}$ of the movement are learned. Additionally, a timing parameter t_{hit} is learned that controls the the transition from swing-back to hitting primitive. Hence, the resulting parameter vector θ included 15 parameters. The reward function was given by the negative distance of the racket to the ball at the estimated hitting time point t_{hit} . The upper-level policy was initialized with five successful examples obtained from another player. The robot was able to increase its success rate to hit the ball from an initial 30% to 80% after 100 roll-outs. A successful hitting movement is illustrated in Figure 2.10.

2.5.6 Learning Tetherball with HiREPS

The HiREPS algorithm was used to learn several solutions for the game of Tetherball [18]. A ball was attached to a string which hung from the

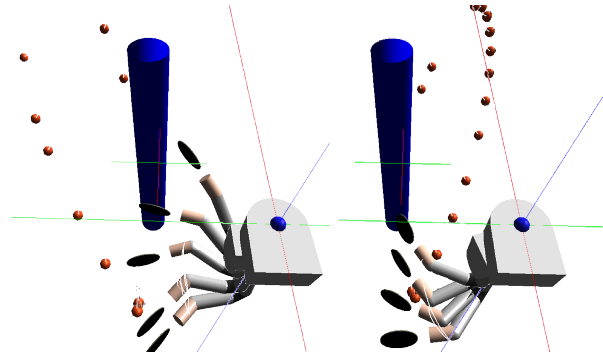


Fig. 2.11 Learning the Tetherball task [18]. The task is to hit a ball such that it winds around the pole. There are two solutions, hitting the ball to the left and hitting the ball to the right. HiREPS was able to learn both solutions in 300 trials. The figure shows both learned movement in simulation.

ceiling. A pole is placed in front of the robot. The task of the robot is to wind the ball around the pole. To achieve the task, the robot first needed to hit the ball once to displace it from its resting pose, and, subsequently, hit it again to arc it around the pole. This task has two solutions: wind the ball around the pole clockwise or counter-clockwise.

The movement was decomposed into a swing-in motion and a hitting motion. For both motions, the shape parameters \mathbf{w} were extracted by kinesthetic teach-in. Both motions were represented by a single set of parameters and the parameters for the two DMPs were jointly learned. For both movements, the final positions \mathbf{g} and velocities $\dot{\mathbf{g}}$ of all seven joints were learned. Additionally, the waiting time between both movements was included in the parameters. This task setup results in a 29-dimensional parameter space.

The reward was determined by the speed of the ball when the ball winds around the pole, where winding around the pole was defined as the ball passing the pole on the opposite side from the initial position. After 300 roll-outs HiREPS was able to learn both solutions within one learning trial as shown in Figure 2.11.

3

Model-based Policy Search

Model-free policy search methods as described in Section 2 are inherently based on sampling trajectories $\tau^{[i]}$ using the robot to find good policies π^* . Sampling trajectories is relatively straightforward in computer simulation. However, when working with mechanical systems, such as robots, each sample corresponds to interacting directly with the robot, which often requires substantial experimental time and causes wear and tear in non-industrial robots. Depending on the task, it can either be easier to learn a model or to learn a policy directly. Model-based policy search methods attempt to address the problem of sample inefficiency by using observed data to learn a forward model of the robot’s dynamics. Subsequently, this forward model is used for *internal* simulations of the robot’s dynamics, based on which the policy is learned.

Model-based policy search algorithms typically assume the following set-up: The state \mathbf{x} evolves according to the Markovian dynamics

$$\mathbf{x}_{t+1} = f(\mathbf{x}_t, \mathbf{u}_t) + \mathbf{w}, \quad \mathbf{x}_0 \sim p(\mathbf{x}_0), \quad (3.1)$$

where f is a nonlinear function, \mathbf{u} is a control signal (action), and \mathbf{w} is additive noise, often chosen to be i.i.d. Gaussian. Moreover, an episodic

set-up is considered where the robot is reset to an initial state \mathbf{x}_0 after executing a policy. The initial state distribution $p(\mathbf{x}_0)$ is often given by a Gaussian distribution $\mathcal{N}(\mathbf{x}_0 | \boldsymbol{\mu}_0^x, \boldsymbol{\Sigma}_0^x)$. Furthermore, we consider finite horizon problems, i.e., the policy-search objective is to find a parametrized policy $\pi_{\boldsymbol{\theta}}^*$ that maximizes the expected long-term reward

$$\pi_{\boldsymbol{\theta}}^* \in \arg \max_{\pi_{\boldsymbol{\theta}}} J_{\boldsymbol{\theta}} = \arg \max_{\pi} \sum_{t=1}^T \gamma^t \mathbb{E}[r(\mathbf{x}_t, \mathbf{u}_t) | \pi_{\boldsymbol{\theta}}], \quad \gamma \in [0, 1], \quad (3.2)$$

where r is an immediate reward signal, γ a discount factor, and the policy $\pi_{\boldsymbol{\theta}}$ is parametrized by parameters $\boldsymbol{\theta}$. Therefore, finding $\pi_{\boldsymbol{\theta}}^*$ in Equation (3.2) is equivalent to finding the corresponding optimal policy parameters $\boldsymbol{\theta}^*$.

For some problems, model-based RL methods have the promise of requiring fewer interactions with the robot than model-free RL by learning a model of the transition mapping in Equation (3.1), while efficiently generalizing to unforeseen situations using a model learned from observed data [6].

The general idea of model-based RL is depicted in Figure 3.1. The learned model is used for internal simulations, i.e., predictions about how the real robot and its environment would behave if it followed the current policy. Based on these internal simulations, the quality of the policy is evaluated using Equation (3.2) and improved accordingly. Subsequently, the updated policy is again evaluated using Equation (3.2) and improved. This policy evaluation/improvement loop terminates when the policy is learned, i.e., it no longer changes and a (local) optimum is attained. Once a policy is learned, it is applied to the robot and a new data set is recorded. Combined with previously collected data, the data set is used to update and refine the learned dynamics model. In theory, this loop continues forever. Note that, only the *application* of the policy requires interacting with the robot; internal simulations and policy learning only use the learned computer model of the robot dynamics.

While the idea of using models in the context of robot learning is well-known since Aboaf's work in the 1980s [2], it has been limited by its strong dependency on the quality of the learned model, which becomes also clear from Figure 3.1: The learned policy is inherently

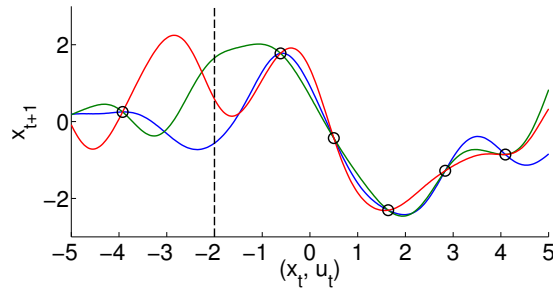


Fig. 3.2 Model errors. In this example, six function values have been observed (black circles). Three functions are shown that could have generated these observations. Any single function (e.g., the maximum likelihood function) produces more or less arbitrary predictions in regions with sparse training data, one example location of which is shown by the dashed line. Instead of selecting a single function approximator, we describe our uncertainty about the underlying function by a probability distribution to be robust to such model errors.

based on internal simulations using the learned model. When the model exactly corresponds to the true dynamics of the robot, sampling from the learned model is equivalent to sampling from the real robot.

However, in practice, the learned model is *not* exact, but only a more or less accurate approximation to the real dynamics. For example, in regions where the training data is sparse, the quality of the learned model can be insufficient as illustrated in Figure 3.2. There are multiple plausible functions that could have generated the observed function values (black circles). In regions with sparse training data, the models and their predictions differ significantly. Any single model

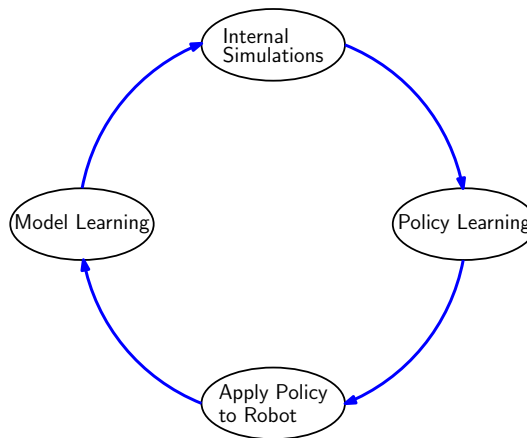


Fig. 3.1 General loop in model-based RL: The learned model is used for internal simulations (mental rehearsal) based on which the policy is improved. The learned policy is then applied to the robot. The data from this interaction is used to refine the model, and the loop continues until the model and the learned policy converge.

leads to overconfident predictions that, in turn, can result in control strategies that are not robust to model errors. This behavior can have a drastic effect in robotics, for example, resulting in the estimation of negative masses or negative friction coefficients. These implausible effects are often exploited by the learning system since they insert energy into the system, causing the system to believe in “perpetuum mobiles”. Therefore, instead of selecting a single model (e.g., the maximum likelihood model), we should describe our uncertainty about the latent underlying function f by a probability distribution $p(f)$ to be robust to such model errors [72, 7, 20]. By taking model uncertainty into account, the perpetuum mobile-effect is substantially less likely.

Since the learned policy inherently relies on the quality of the learned *forward model*, which essentially serves as a simulator of the robot, model errors can not only cause degraded policies, but they also often drastically bias the learning process. Hence, the literature on model-based policy search largely focuses on model building, i.e., explaining what kind of model is used for the forward dynamics and how it is trained.

Approaches to Dealing with Uncertain Models. Dealing with inaccurate dynamics models is one of the biggest challenges in model-based RL since small errors in the model can lead to large errors in the policy [6]. Inaccuracies might stem from an overly restrictive model class or from the lack of sufficiently rich data sets used for training the models, which can lead to under-modeling the true forward dynamics. Moreover, system noise typically adds another source of uncertainty.

Typically, a certainty-equivalence assumption is made¹, and the maximum likelihood model is chosen for planning [72, 7]. However, this certainty-equivalence assumption is violated in most interesting cases and can lead to large policy errors. Moreover, as mentioned already in [2], many approaches obtain derivatives of the expected return by back-propagating derivatives through learned forward models of the system. This step is particularly prone to model errors since gradient-based

¹It is assumed that the optimal policy for the learned model corresponds to the optimal policy for the true dynamics. Uncertainty about the learned model is neglected.

optimizers improve the policy parameters along their gradients. The learned policy needs to be robust to compensate for model errors such that it results in good performance when applied to the real system. Learning faithful dynamics models is crucial for building robust policies and remains one of the biggest challenges in model-based policy search.

In [45], the authors model unknown error dynamics using receptive-field weighted regression [70]. Explicit modeling of unknown disturbances leads to increased robustness of the learned controllers. The idea of designing controllers in the face of inaccurate (idealized) forward models is closely related to robust control in classical robotics. Robust control aims to achieve guaranteed performance or stability in the presence of (typically bounded) modeling errors. For example, \mathcal{H}_∞ loop-shaping [44] guarantees that the system remains close to its expected behavior even if (bounded) disturbances enter the system. In adaptive control, parameter uncertainties are usually described by unbounded probability distributions [5]. Model parameter uncertainty is typically not used in designing adaptive control algorithms. Instead, the estimates of the parameters are treated as the true ones [90]. An approach to designing adaptive controllers that do take uncertainty about the model parameters into account is stochastic adaptive control [5]. When reducing parameter uncertainty by probing, stochastic adaptive control leads to the principle of dual control [26]. Adaptive dual control has been investigated mainly for linear systems [90]. An extension of dual adaptive control to the case of nonlinear systems with affine controls was proposed in [25]. A minimum-variance control law is obtained, and uncertainty about the model parameters is penalized to improve their estimation, eliminating the need for prior system identification.

RL approaches that explicitly address the problem of inaccurate models in robotics have only been introduced recently [72, 7, 45, 51, 53, 1, 34, 21, 20]. For instance, in [1], the key idea is to use a real-life trial to evaluate a policy, but then use a crude model of the system to estimate the derivative of the evaluation with respect to the policy parameters (and suggest local improvements). In particular, the suggested algorithm iteratively updates the model f according to $f^{[i+1]}(\mathbf{x}_t, \mathbf{u}_t) = f^{[i]}(\mathbf{x}_t, \mathbf{u}_t) + \mathbf{x}_{t+1}^{[i]} - f^{[i]}(\mathbf{x}_t^{[i]}, \mathbf{u}_t^{[i]})$, where t indexes time.

Here, $\mathbf{x}_t^{[i]}, \mathbf{u}_t^{[i]}$ are taken from an observed trajectory. It follows that the updated model $f^{[i+1]}$ predicts the observed trajectory exactly, i.e., $f^{[i+1]}(\mathbf{x}_t^{[i]}, \mathbf{u}_t^{[i]}) = \mathbf{x}_{t+1}^{[i]}$. The algorithm evaluates the policy gradients along the trajectory of states and controls in the real system. In contrast, a typical model-based approach evaluates the derivatives along the trajectory predicted by the model, which does not correspond to the trajectory of the real system when the model is inexact. Note that the approach in [1] does not directly generalize to stochastic transition dynamics or systems with hidden states. Moreover, an approximate parametric model of the underlying dynamics need to be known in advance.

Major Challenges in Model-based Policy Search. There are three general challenges that need to be addressed in model-based policy search methods: what model to learn, how to use the model for long-term predictions, and how to update the policy based on the long-term predictions. These three challenges correspond to the three components in Figure 3.1 that do not require physical interaction with the robot: model learning, internal simulations, and policy learning.

In Section 3.1, we give a brief overview of two models that are frequently used in model-based policy search [7, 51, 53, 34, 21, 20], locally weighted (Bayesian) regression (LWBR) and Gaussian processes (GPs) [65]. In Section 3.2, we discuss two general methods of how to use these models for long-term predictions: stochastic inference, i.e., sampling, and deterministic approximate inference. In Section 3.3, we briefly discuss multiple options of updating the policy.

After having introduced these general concepts, in Section 3.4, we discuss model-based policy search algorithms that combine the learned models and inference algorithms shown in Table 3.1. We focus on the episodic case, wherein the start state distribution $p(\mathbf{x}_0)$ is known. In particular, we detail four model-based policy search methods. First, we start with PEGASUS, a general concept for efficient trajectory sampling in stochastic MDPs for a given model [52]. Second, we present two ways of combining PEGASUS with LWBR for learning robust controllers to fly helicopters [7, 51, 53]. Third, we present an approach

Table 3.1 Model-based policy search approaches grouped by the learned model and the method of generating trajectories. Due to the simplicity of sampling trajectories, most policy search methods follow this approach. Deterministic trajectory predictions can only be performed in special cases where closed-form approximate inference is possible. Although they are mathematically more involved than trajectory sampling, they do not suffer from large variances of the samples. Moreover, they can allow for analytic gradient computations, which is crucial in the case of hundreds of policy parameters.

<i>Learned Forward Model</i>	<i>Trajectory Prediction</i>	
	stochastic	deterministic
(Locally) linear models	[7, 51, 53]	—
Gaussian Processes	[34]	[20, 21]

for using the PEGASUS algorithm for sampling trajectories from GP forward models [34] in the context of learning a blimp controller. Finally, as a fourth approach, we outline the ideas of the PILCO policy search framework [20, 21] that combines efficient deterministic approximate inference for long-term predictions with GP dynamics models for learning to control mechanical systems and robot manipulators.

3.1 Probabilistic Forward Models

To reduce the effect of model errors, probabilistic models that express uncertainty about the underlying transition dynamics are preferable to deterministic models that imply a certainty-equivalence assumption, e.g., maximum-likelihood models of the transition dynamics.

In the following, we briefly introduce two promising nonparametric probabilistic models that are frequently used for learning the forward dynamics in the context of reinforcement learning and robotics: Locally weighted Bayesian regression (LWBR), [7, 51, 53] and Gaussian processes [34, 20, 21].

3.1.1 Locally Weighted Bayesian Regression

Let us start with the linear regression model, where the transition dynamics are given as

$$\mathbf{x}_{t+1} = [\mathbf{x}_t, \mathbf{u}_t]^T \boldsymbol{\psi} + \mathbf{w}, \quad \mathbf{w} \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}_w). \quad (3.3)$$

Here, $\boldsymbol{\psi}$ are the parameters of the Bayesian linear regression model, and $\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}_w)$ is i.i.d. Gaussian system noise. The model is linear

in the unknown parameters $\boldsymbol{\psi}$ that weight the input $(\boldsymbol{x}_t, \boldsymbol{u}_t)$.

In Bayesian linear regression, we place prior distributions on the parameters $\boldsymbol{\psi}$ and on the noise variance $\boldsymbol{\Sigma}_w$. Typically, the prior distribution on $\boldsymbol{\psi}$ is Gaussian with mean \boldsymbol{m} and covariance \boldsymbol{S} , and the prior on the diagonal entries $1/\sigma_i^2$ of $\boldsymbol{\Sigma}_w^{-1}$ is a Gamma distribution with scale and shape parameters η and $\boldsymbol{\Xi}$, respectively, such that the overall model is conjugate, see Figure 3.3, where we denote the training inputs by $[\boldsymbol{X}, \boldsymbol{U}]$ and the training targets by \boldsymbol{y} , respectively. In the (Bayesian) linear regression model Equation (3.3), it is fairly straightforward to find maximum likelihood estimates or posterior distributions of the parameters $\boldsymbol{\psi}$. However, the model itself is not very expressive since it assumes an overall linear relationship between the inputs $(\boldsymbol{x}_t, \boldsymbol{u}_t)$ and the successor state \boldsymbol{x}_{t+1} .

The idea of *locally weighted linear regression* (LWR) is to exploit the good properties of the linear regression model but to allow for a more general class of functions: locally linear functions. LWR finds a locally linear approximation of the underlying function [15]. For this purpose, every test input $(\boldsymbol{x}_t, \boldsymbol{u}_t)$ is equipped with a weighting factor b_i that determines how close training point $(\boldsymbol{x}_i, \boldsymbol{u}_i)$ is to $(\boldsymbol{x}_t, \boldsymbol{u}_t)$. An example for such a weight is the Gaussian-shaped weighting $b_i = \exp(-\|(\boldsymbol{x}_i, \boldsymbol{u}_i) - (\boldsymbol{x}_t, \boldsymbol{u}_t)\|^2/\kappa^2)$. If the distance between $(\boldsymbol{x}_i, \boldsymbol{u}_i)$ and $(\boldsymbol{x}_*, \boldsymbol{u}_*)$ is much larger than κ , the corresponding weight b_i declines to 0. Since these weights have to be computed for each query point, it is insufficient to store only the parameters $\boldsymbol{\psi}$, but the entire training data set $[\boldsymbol{X}, \boldsymbol{U}]$ is required, resulting in a nonparametric approach.

As in Bayesian linear regression, we can place priors on the parameters and the noise covariance. For simplicity, let us assume a known noise covariance matrix $\boldsymbol{\Sigma}_w$ and a zero-mean prior Gaussian distribu-

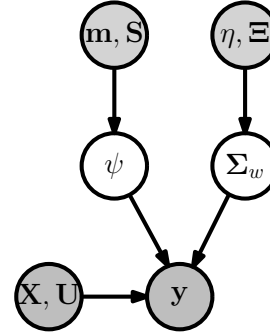


Fig. 3.3 Graphical model for Bayesian linear regression: A Gaussian prior with parameters $\boldsymbol{m}, \boldsymbol{S}$ is placed on the model parameters $\boldsymbol{\psi}$ and Gamma priors with parameters $\eta, \boldsymbol{\Xi}$ are placed on the diagonal entries of the precision matrix $\boldsymbol{\Sigma}_w^{-1}$. Training inputs and targets are denoted by $[\boldsymbol{X}, \boldsymbol{U}]$ and \boldsymbol{y} , respectively.

tion $\mathcal{N}(\boldsymbol{\psi} | \mathbf{0}, \mathbf{S})$ on the parameters $\boldsymbol{\psi}$. For each query point $(\mathbf{x}_t, \mathbf{u}_t)$, a posterior distribution over the parameters $\boldsymbol{\psi}$ is computed according to Bayes' theorem as

$$p(\boldsymbol{\psi} | \mathbf{X}, \mathbf{U}, \mathbf{y}) = \frac{p(\boldsymbol{\psi})p(\mathbf{y} | \mathbf{X}, \mathbf{U}, \boldsymbol{\psi})}{p(\mathbf{y} | \mathbf{X}, \mathbf{U})} \propto p(\boldsymbol{\psi})p(\mathbf{y} | \mathbf{X}, \mathbf{U}, \boldsymbol{\psi}). \quad (3.4)$$

For notational convenience, we define $\tilde{\mathbf{X}} := [\mathbf{X}, \mathbf{U}]$. The posterior mean and covariance of $\boldsymbol{\psi}$ at $(\mathbf{x}_t, \mathbf{u}_t)$ are given as

$$\mathbb{E}[\boldsymbol{\psi} | \tilde{\mathbf{X}}, \mathbf{y}] = \mathbf{S} \tilde{\mathbf{X}} \mathbf{B} \underbrace{(\mathbf{B} \tilde{\mathbf{X}}^T \mathbf{S} \tilde{\mathbf{X}} \mathbf{B} + \boldsymbol{\Sigma}_w)^{-1}}_{=\boldsymbol{\Omega}^{-1}} \mathbf{y} = \mathbf{S} \tilde{\mathbf{X}} \mathbf{B} \boldsymbol{\Omega}^{-1} \mathbf{y} \quad (3.5)$$

$$\text{cov}[\boldsymbol{\psi} | \tilde{\mathbf{X}}, \mathbf{y}] = \mathbf{S} - \mathbf{S}^T \tilde{\mathbf{X}} \mathbf{B} \boldsymbol{\Omega}^{-1} \mathbf{B} \tilde{\mathbf{X}}^T \mathbf{S}, \quad (3.6)$$

$$b_i = \exp(-\|(\mathbf{x}_i, \mathbf{u}_i) - (\mathbf{x}_t, \mathbf{u}_t)\|^2 / \kappa^2), \quad (3.7)$$

respectively, where $\mathbf{B} = \text{diag}(b_1, \dots, b_n)$, and \mathbf{y} are the training targets.

Predictive Distribution. The mean and covariance of the predictive distribution $p(\mathbf{x}_{t+1})$ for a given state-control pair $(\mathbf{x}_t, \mathbf{u}_t)$ are

$$\boldsymbol{\mu}_{t+1}^x = [\mathbf{x}_t, \mathbf{u}_t]^T \mathbb{E}[\boldsymbol{\psi} | \tilde{\mathbf{X}}, \mathbf{y}] = [\mathbf{x}_t, \mathbf{u}_t]^T \mathbf{S} \tilde{\mathbf{X}} \mathbf{B} \boldsymbol{\Omega}^{-1} \mathbf{y}, \quad (3.8)$$

$$\boldsymbol{\Sigma}_{t+1}^x = [\mathbf{x}_t, \mathbf{u}_t]^T \text{cov}[\boldsymbol{\psi} | \tilde{\mathbf{X}}, \mathbf{y}] [\mathbf{x}_t, \mathbf{u}_t] + \boldsymbol{\Sigma}_w, \quad (3.9)$$

respectively. In practice, the posterior mean and covariance over the parameters $\boldsymbol{\psi}$ can be computed more efficiently by applying matrix inversion lemmas [15] and exploiting sparsity.

Let us have a look at the predictive covariance when $[\mathbf{x}_t, \mathbf{u}_t]$ is far away from the training set $[\mathbf{X}, \mathbf{U}]$: The weight matrix \mathbf{B} is almost zero, which leads to a posterior variance over the model parameters $\boldsymbol{\psi}$ that is equal to the prior uncertainty \mathbf{S} , see Equation (3.6). Hence, the predictive variance at $[\mathbf{x}_t, \mathbf{u}_t]$ is non-zero, unlike the non-Bayesian locally weighted regression case.

3.1.2 Gaussian Process Regression

A Gaussian process is a distribution $p(f)$ over functions f . Formally, a GP is a collection of random variables f_1, f_2, \dots any finite number of which is Gaussian distributed [65]. In the context of this section, a GP

is placed over transition functions. Since the GP is a nonparametric model, it suffices to specify high-level assumptions, such as differentiability or periodicity, on the underlying function. These high-level properties are typically easier to specify than an explicit parametric model.

A GP is completely specified by a mean function $m(\cdot)$ and a positive semidefinite covariance function/kernel $k(\cdot, \cdot)$. Standard assumptions in GP models are a prior mean function $m \equiv 0$ and the covariance function

$$k(\tilde{\mathbf{x}}_p, \tilde{\mathbf{x}}_q) = \sigma_f^2 \exp\left(-\frac{1}{2}(\tilde{\mathbf{x}}_p - \tilde{\mathbf{x}}_q)^T \mathbf{\Lambda}^{-1}(\tilde{\mathbf{x}}_p - \tilde{\mathbf{x}}_q)\right) + \delta_{pq} \sigma_w^2 \quad (3.10)$$

with $\tilde{\mathbf{x}} := [\mathbf{x}^T \mathbf{u}^T]^T$. In Equation (3.10), we defined $\mathbf{\Lambda} := \text{diag}([\ell_1^2, \dots, \ell_D^2])$, which depends on the characteristic length-scales ℓ_i , and σ_f^2 is the prior variance of the latent function f . Given n training inputs $\tilde{\mathbf{X}} = [\tilde{\mathbf{x}}_1, \dots, \tilde{\mathbf{x}}_n]$ and corresponding training targets $\mathbf{y} = [y_1, \dots, y_n]^T$, the posterior GP hyper-parameters (length-scales ℓ_i , signal variance σ_f^2 , and noise variance σ_w^2) are learned using evidence maximization [43, 65].

Predictive Distribution. The posterior GP is a one-step prediction model, and the predicted successor state \mathbf{x}_{t+1} is Gaussian distributed

$$p(\mathbf{x}_{t+1} | \mathbf{x}_t, \mathbf{u}_t) = \mathcal{N}(\mathbf{x}_{t+1} | \boldsymbol{\mu}_{t+1}^x, \boldsymbol{\Sigma}_{t+1}^x), \quad (3.11)$$

$$\boldsymbol{\mu}_{t+1}^x = \mathbb{E}_f[f(\mathbf{x}_t, \mathbf{u}_t)], \quad \boldsymbol{\Sigma}_{t+1}^x = \text{var}_f[f(\mathbf{x}_t, \mathbf{u}_t)], \quad (3.12)$$

where the mean and variance of the GP prediction are

$$\boldsymbol{\mu}_{t+1}^x = \mathbf{k}_*^T \mathbf{K}^{-1} \mathbf{y} = \mathbf{k}_*^T \boldsymbol{\beta}, \quad (3.13)$$

$$\boldsymbol{\Sigma}_{t+1}^x = k_{**} - \mathbf{k}_*^T \mathbf{K}^{-1} \mathbf{k}_*, \quad (3.14)$$

respectively, with $\mathbf{k}_* := k(\tilde{\mathbf{X}}, \tilde{\mathbf{x}}_t)$, $k_{**} := k(\tilde{\mathbf{x}}_t, \tilde{\mathbf{x}}_t)$, $\boldsymbol{\beta} := \mathbf{K}^{-1} \mathbf{y}$, and where \mathbf{K} is the kernel matrix with entries $K_{ij} = k(\tilde{\mathbf{x}}_i, \tilde{\mathbf{x}}_j)$.

Note that, far away from the training data, the predictive uncertainty in Equation (3.14) corresponds to the prior uncertainty about the function, i.e., even for deterministic systems where $\boldsymbol{\Sigma}_w = \mathbf{0}$, we obtain $k_{**} > 0$. Therefore, the GP is a nonparametric, non-degenerate model.

3.2 Long-Term Predictions with a Given Model

In the following, we assume that a model for the transition dynamics is known. Conditioned on this model, we distinguish between two approaches for generating long-term predictions: approaches based on Monte-Carlo sampling or deterministic approximate inference.

3.2.1 Sampling-based Trajectory Prediction: PEGASUS

PEGASUS (Policy Evaluation-of-Goodness And Search Using Scenarios) is a conceptual framework for trajectory sampling in stochastic MDPs [52]. The key idea is to transform the stochastic MDP into an augmented deterministic MDP. For this purpose, PEGASUS assumes access to a simulator with no internal random number generator. When sampling from this model, PEGASUS provides the random numbers externally in advance. In this way, PEGASUS reduces the sampling variance drastically. Therefore, sampling following the PEGASUS approach is also commonly used in model-free policy search, see Section 2.

3.2.1.1 Trajectory Sampling and Policy Evaluation

Assume that a forward model of the system at hand is given that can be used for sampling trajectories. If the state transitions are stochastic, computing the expected long-term reward

$$J_{\theta} = \sum_{t=0}^T \gamma^t \mathbb{E}[r(\mathbf{x}_t) | \pi_{\theta}], \quad \mathbf{x}_0 \sim p(\mathbf{x}_0), \quad \gamma \in [0, 1], \quad (3.15)$$

will require many sample trajectories for computing the approximation \tilde{J} to J , where

$$\tilde{J}_{\theta} = \frac{1}{m} \sum_{i=1}^m J_{\theta}(\mathbf{x}_0^{[i]}) \quad (3.16)$$

with samples $\mathbf{x}_0^{[i]}$ from $p(\mathbf{x}_0)$. Computing reliable policy gradients will require even more samples for robust derivatives. However, as the limit of an infinite number of samples, we obtain $\lim_{m \rightarrow \infty} \tilde{J}_{\theta} = J_{\theta}$.

For more efficient computations, PEGASUS augments the state \mathbf{x} by an externally given sequence of random values $\mathbf{w}_0, \mathbf{w}_1, \dots$. To draw

Algorithm 15 PEGASUS algorithm for sampling trajectories

Init: $g(\mathbf{x}, \mathbf{u}, \mathbf{w})$, reward r , random numbers $\mathbf{w}_0, \mathbf{w}_1, \dots$, initial state distribution $p(\mathbf{x}_0)$, policy π_θ

for $i = 1, \dots, m$ **do**

$\mathbf{x}_0^{[i]} \sim p(\mathbf{x}_0)$ \triangleright Sample “scenario” from initial state distribution

for $t = 0, \dots, T - 1$ **do**

$\mathbf{x}_{t+1}^{[i]} = g(\mathbf{x}_t^{[i]}, \pi_\theta(\mathbf{x}_t), \mathbf{w}_t)$ \triangleright Succ. state in augmented MDP

end for

end for

$\tilde{J}_\theta = \frac{1}{m} \sum_{i=1}^m \sum_{t=1}^T \gamma^t r(\mathbf{x}_t^{[i]})$ \triangleright Estimate expected long-term reward

a sample from $p(\mathbf{x}_{t+1} | \mathbf{x}_t, \mathbf{u}_t)$, PEGASUS uses a-priori given noise values \mathbf{w}_t to compute the sampled state \mathbf{x}_{t+1} , such that $\mathbf{x}_{t+1} = g(\mathbf{x}_t, \mathbf{u}_t, \mathbf{w}_t)$. Since the sequence of random numbers is fixed, repeating the same experiment results in the identical sample trajectory.

PEGASUS can be described as generating m Monte Carlo trajectories and taking their average reward, but the randomization is determined in advance. It can be shown that solving the augmented deterministic MDP is equivalent to solving the original stochastic MDP [52]. The PEGASUS algorithm is summarized in Algorithm 15.

3.2.1.2 Practical Considerations

While sampling using the PEGASUS algorithm can be performed relatively efficiently, robust low-variance estimates of the expected long-term cost require a lot of samples. Therefore, policy search methods based on trajectory sampling are practically limited by a relatively small number of a few tens of policy parameters they can manage [50].² Although the policy gradients can also be computed (e.g., with finite difference approximations), the sample-based nature of PEGASUS leads to derivatives with high variance, which renders them largely useless. Thus, policy updates in the context of PEGASUS do usually not

²“Typically, PEGASUS policy search algorithms have been using [...] maybe on the order of ten parameters or tens of parameters; so, 30, 40 parameters, but not thousands of parameters [...]” [50]

rely on gradients [53, 7, 34]. Of course all model-free approaches from Section 2 for estimating the policy gradients can be used in conjunction with the PEGASUS idea of sampling trajectories from a given model.

An alternative to sampling trajectories are deterministic approximate inference methods for predicting trajectories, such as linearization [4], moment matching, or the unscented transformation [32].

3.2.2 Deterministic Long-Term Predictions

Instead of performing stochastic sampling, a probability distribution $p(\boldsymbol{\tau})$ over trajectories $\boldsymbol{\tau} = (\mathbf{x}_0, \dots, \mathbf{x}_T)$ can also be computed using deterministic approximations, such as linearization [4], sigma-point methods (e.g., the unscented transformation [32]), or moment matching. These common inference methods approximate unwieldy predictive distributions by Gaussians.

Assuming a joint Gaussian probability distribution $p(\mathbf{x}_t, \mathbf{u}_t) = \mathcal{N}([\mathbf{x}_t, \mathbf{u}_t] | \boldsymbol{\mu}_t^{xu}, \boldsymbol{\Sigma}_t^{xu})$, the problem of computing the successor state distribution $p(\mathbf{x}_{t+1})$ corresponds to solving the integral

$$p(\mathbf{x}_{t+1}) = \iiint p(\mathbf{x}_{t+1} | \mathbf{x}_t, \mathbf{u}_t) p(\mathbf{x}_t, \mathbf{u}_t) d\mathbf{x}_t d\mathbf{u}_t d\mathbf{w}, \quad (3.17)$$

where $\mathbf{x}_{t+1} = f(\mathbf{x}_t, \mathbf{u}_t) + \mathbf{w}$. If the transition function f is nonlinear, $p(\mathbf{x}_{t+1})$ is non-Gaussian and we have to resort to approximate inference techniques. A convenient approximation of the unwieldy predictive distribution $p(\mathbf{x}_{t+1})$ is the Gaussian $\mathcal{N}(\mathbf{x}_{t+1} | \boldsymbol{\mu}_{t+1}^x, \boldsymbol{\Sigma}_{t+1}^x)$. The mean $\boldsymbol{\mu}_{t+1}^x$ and covariance $\boldsymbol{\Sigma}_{t+1}^x$ of this predictive distribution can be computed in various ways. In the following, we outline three commonly used approaches: linearization, the unscented transformation, and moment matching.

Linearization. One way of computing $\boldsymbol{\mu}_{t+1}^x$ and $\boldsymbol{\Sigma}_{t+1}^x$ is to linearize the transition function $f \approx \mathbf{F}$ locally around $(\boldsymbol{\mu}_t^x, \boldsymbol{\mu}_t^u)$ and, subsequently, estimate the predictive covariance by mapping the Gaussian input distribution through the linearized system. With linearization, we obtain the predictive mean and covariance given by $\boldsymbol{\mu}_{t+1}^x = f(\boldsymbol{\mu}_t^{xu})$ and $\boldsymbol{\Sigma}_{t+1}^x = \mathbf{F} \boldsymbol{\Sigma}_t^{xu} \mathbf{F}^T + \boldsymbol{\Sigma}_w$, respectively. Figure 3.4 illustrates the idea of linearization.

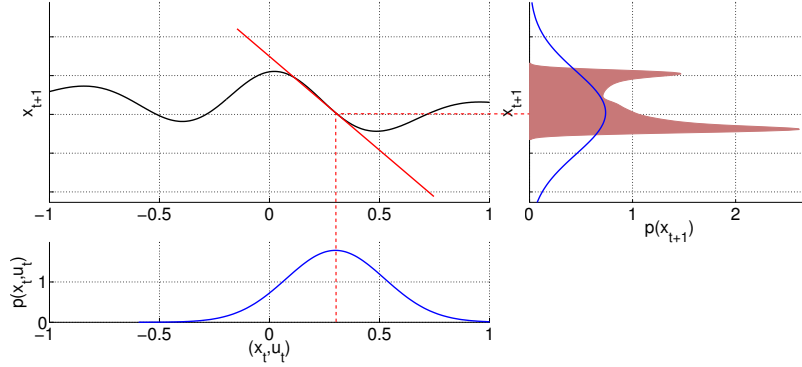


Fig. 3.4 Computing an approximate predicted distribution using linearization. A Gaussian distribution $p(\mathbf{x}_t, \mathbf{u}_t)$ (lower-left panel) needs to be mapped through a nonlinear function (black, upper-left panel). The true predictive distribution is represented by the shaded area in the right panel. To obtain a Gaussian approximation of the unwieldy shaded distribution, the nonlinear function is linearized (red line, upper-left panel) at the mean of the input distribution. Subsequently, the Gaussian is mapped through this linear approximation and yields the blue Gaussian approximate predictive distribution $p(\mathbf{x}_{t+1})$ shown in the right panel.

Linearization is conceptually straightforward and computationally efficient. Note that this approach leaves the Gaussian input distribution $p(\mathbf{x}_t, \mathbf{u}_t)$ untouched but approximates the transition function f . A potential disadvantage is that the transition function f needs to be differentiable to perform the linearization.³ Moreover, linearization can easily underestimate predictive variances, which can cause policies to be too aggressive, causing damage on real robot systems.

Unscented Transformation. The key idea behind the unscented transformation [32] is to represent the distribution $p(\mathbf{x}_t, \mathbf{u}_t)$ by a set of deterministically chosen sigma points $(\mathcal{X}_t^{[i]}, \mathcal{U}_t^{[i]})$. For these sigma points, the corresponding exact function values are computed. The mean $\boldsymbol{\mu}_{t+1}^x$ and covariance $\boldsymbol{\Sigma}_{t+1}^x$ of the predictive distribution $p(\mathbf{x}_{t+1})$ are computed from the weighted mapped sigma points. In particular,

³Differentiability assumptions can be problematic in robotics. For instance, contacts in locomotion and manipulation can render this assumption invalid.

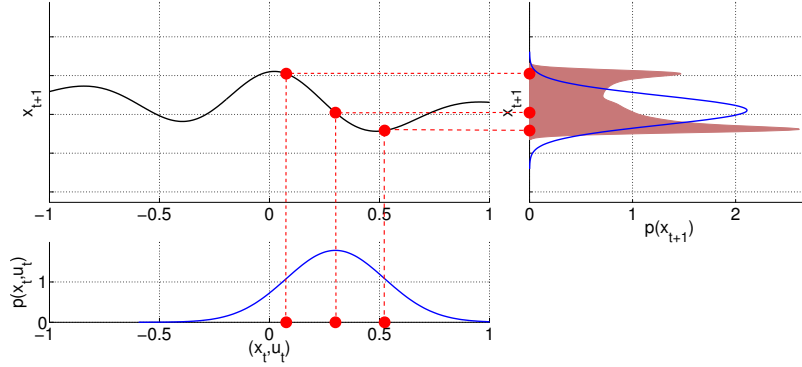


Fig. 3.5 Computing an approximate predicted distribution using the unscented transformation. A Gaussian distribution $p(\mathbf{x}_t, \mathbf{u}_t)$ (lower-left panel) needs to be mapped through a nonlinear function (black, upper-left panel). The true predictive distribution is represented by the shaded area in the right panel. To obtain a Gaussian approximation of the unwieldy shaded distribution, the input distribution is represented by three sigma points (red dots in lower-left panel). Subsequently, the sigma points are mapped through the nonlinear function (upper-left panel) and their sample mean and covariance yield the blue Gaussian approximate predictive distribution $p(\mathbf{x}_{t+1})$ shown in the right panel.

we obtain

$$\boldsymbol{\mu}_{t+1}^x = \sum_{i=0}^{2d} w_m^{[i]} f(\mathcal{X}_t^{[i]}, \mathcal{U}_t^{[i]}), \quad (3.18)$$

$$\boldsymbol{\Sigma}_{t+1}^x = \sum_{i=0}^{2d} w_c^{[i]} (f(\mathcal{X}_t^{[i]}, \mathcal{U}_t^{[i]}) - \boldsymbol{\mu}_{t+1}^x)(f(\mathcal{X}_t^{[i]}, \mathcal{U}_t^{[i]}) - \boldsymbol{\mu}_{t+1}^x)^T, \quad (3.19)$$

respectively, where d is the dimensionality of (\mathbf{x}, \mathbf{u}) , $(\mathcal{X}_t^{[i]}, \mathcal{U}_t^{[i]})$ are *sigma points*, i.e., deterministically chosen samples from the joint distribution $p(\mathbf{x}_t, \mathbf{u}_t)$, and $w_m^{[i]}$ and $w_c^{[i]}$ are weights. For further details on the unscented transformation, we refer to [32, 82]. Figure 3.5 illustrates the idea of the unscented transformation.

Note that the unscented transformation approximates the Gaussian distribution $p(\mathbf{x}_t, \mathbf{u}_t)$ by sigma points, which are subsequently mapped through the original transition function f . The unscented transformation does not require differentiability and is expected to yield more accurate approximations of the predictive distribution $p(\mathbf{x}_{t+1})$ than an explicit linearization [91].

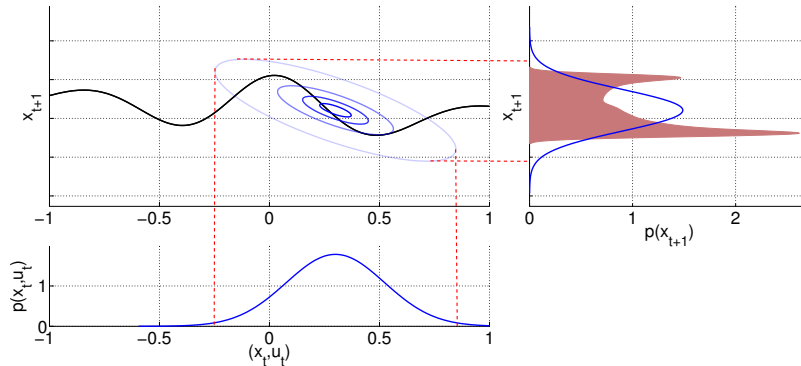


Fig. 3.6 Computing an approximate predicted distribution using moment matching. A Gaussian distribution $p(\mathbf{x}_t, \mathbf{u}_t)$ (lower-left panel) needs to be mapped through a nonlinear function (black, upper-left panel). The true predictive distribution is represented by the shaded area in the right panel. To obtain a Gaussian approximation of the unwieldy shaded distribution, the mean and covariance of the shaded distribution are computed analytically. These first and second-order moments fully determine the blue Gaussian approximate predictive distribution $p(\mathbf{x}_{t+1})$ shown in the right panel. The contour lines in the upper-left panel represent the joint distribution between inputs and prediction.

Moment Matching. The idea of moment matching is to compute the predictive mean and covariance of $p(\mathbf{x}_{t+1})$ exactly and approximate $p(\mathbf{x}_{t+1})$ by a Gaussian that possesses the exact mean and covariance. Here, neither the joint distribution $p(\mathbf{x}_t, \mathbf{u}_t)$ nor the transition function f are approximated. The moment-matching approximation is the best unimodal approximation of the predictive distribution in the sense that it minimizes the Kullback-Leibler divergence between the true predictive distribution and the unimodal approximation [12]. Figure 3.6 illustrates the idea of moment matching.

3.2.2.1 Practical Considerations

The exact moments can be computed only in special cases since the required integrals for computing the predictive mean and covariance might be intractable. Moreover, an exact moment-matching approximation is typically computationally more expensive than approximations by means of linearization or sigma points.

Unlike sampling-based approaches such as PEGASUS, deterministic approximate inference methods for long-term planning can be used

to learn several thousands of policy parameters [20]. We will see examples in Section 3.4.2.1. The reason why deterministic long-term predictions can learn policies with many parameters is that gradients can be computed analytically. Therefore, these gradient estimates do not suffer from high variances, a typical problem with sampling-based estimation. Nevertheless, deterministic inference often requires more implementation effort than sampling approaches.

3.3 Policy Updates

Having introduced two major model classes and two general ways of performing long-term predictions with these models, in the following, we will discuss ways of updating the policy. We distinguish between gradient-free and gradient-based policy updates.

3.3.1 Model-based Policy Updates without Gradient Information

Gradient-free methods are probably the easiest way of updating the policy since they do not require the computation or estimation of policy gradients. By definition they also have no differentiability constraints on the policy or the transition model. Standard gradient-free optimization methods are the Nelder-Mead method [47], a heuristic simplex method, or hill-climbing, a local search method that is closely related to simulated annealing [69]. Due to their simplicity and the small required computational effort, they are commonly used in the context of model-based policy search [7, 34, 51, 53], especially in combination with sampling-based trajectory generation.

A clear disadvantage of gradient-free optimization is their relatively slow convergence rate. For faster convergence, we can use gradient-based policy updates, which are introduced in the following sections.

3.3.2 Model-based Policy Updates with Gradient Information

Gradient-based policy updates are expected to yield faster convergence than gradient-free updates. We distinguish between two cases: a

sample-based estimation of the policy gradients and an analytic computation of the policy gradients $dJ_{\theta}(\boldsymbol{\theta})/d\boldsymbol{\theta}$.

3.3.2.1 Sampling-based Policy Gradients

When we use sample trajectories $\tau^{[i]}$ from the learned model to estimate the expected long-term reward J_{θ} in Equation (3.2), we can numerically approximate the gradient $dJ_{\theta}/d\boldsymbol{\theta}$.

The easiest way of estimating gradients is to use *finite difference methods*. However, finite difference methods require $O(F)$ many evaluations of the expected long-term reward J_{θ} , where F is the number of policy parameters $\boldsymbol{\theta}$. Since each of these evaluations is based on the average of m sample roll-outs, the required number of sample trajectories quickly becomes excessive. In the model-based set-up, this is just a computational problem but not a problem of wearing the robot out since the samples are generated from the model and not the robot itself.

There are several ways of making model-based gradient estimation more efficient: First, for a more robust estimate of J_{θ} , i.e., an estimate with smaller variance, the PEGASUS approach [52] can be used. Second, for more efficient gradient estimation any of the model-free methods presented in Section 2 for gradient estimation can be used in the model-based context. The only difference is that instead of the robot, the learned model is used to generate trajectories. To the best of our knowledge, there are currently not many approaches for model-based policy search based on sampling-based gradients using the methods from Section 2.

3.3.2.2 Analytic Policy Gradients

Computing the gradients $dJ_{\theta}/d\boldsymbol{\theta}$ analytically requires the policy, the (expected) reward function, and the learned transition model to be differentiable. Despite this constraint, analytic gradient computations are a viable alternative to sampling-based gradients for two major reasons: First, they do not suffer from the sampling variance, which is especially pronounced when computing gradients. Second, the computational effort scales very favorably with the number of policy parameters, allowing for learning policies with thousands of parameters. However, due

to the repeated application of the chain-rule, the computation of the gradient itself is often mathematically more involved than a sampling-based estimate.

Let us consider an example where the immediate reward r only depends on the state (generalizations to control-dependent rewards are straightforward) and the system dynamics are deterministic, such that $\mathbf{x}_{t+1} = f(\mathbf{x}_t, \mathbf{u}_t) = f(\mathbf{x}_t, \pi_{\boldsymbol{\theta}}(\mathbf{x}_t, \boldsymbol{\theta}))$, where f is a (nonlinear) transition function, $\pi_{\boldsymbol{\theta}}$ is the (deterministic) policy, and $\boldsymbol{\theta}$ are the policy parameters. The gradient of the long-term reward $J_{\boldsymbol{\theta}} = \sum_t \gamma^t r(\mathbf{x}_t)$ with respect to the policy parameters is obtained by applying the chain-rule repeatedly:

$$\frac{dJ_{\boldsymbol{\theta}}}{d\boldsymbol{\theta}} = \sum_t \gamma^t \frac{dr(\mathbf{x}_t)}{d\boldsymbol{\theta}} = \sum_t \gamma^t \frac{\partial r(\mathbf{x}_t)}{\partial \mathbf{x}_t} \frac{d\mathbf{x}_t}{d\boldsymbol{\theta}} \quad (3.20)$$

$$= \sum_t \gamma^t \frac{\partial r(\mathbf{x}_t)}{\partial \mathbf{x}_t} \left(\frac{\partial \mathbf{x}_t}{\partial \mathbf{x}_{t-1}} \frac{d\mathbf{x}_{t-1}}{d\boldsymbol{\theta}} + \frac{\partial \mathbf{x}_t}{\partial \mathbf{u}_{t-1}} \frac{d\mathbf{u}_{t-1}}{d\boldsymbol{\theta}} \right). \quad (3.21)$$

From these equations we observe that the total derivative $d\mathbf{x}_t/d\boldsymbol{\theta}$ depends on the total derivative $d\mathbf{x}_{t-1}/d\boldsymbol{\theta}$ at the previous time step. Therefore, the derivative $dJ_{\boldsymbol{\theta}}/d\boldsymbol{\theta}$ can be computed iteratively.

Extension to Probabilistic Models and Stochastic MDPs. For the extension to derivatives in stochastic MDPs and/or probabilistic models, we have to make a few adaptations to the gradients in Equation (3.20)–(3.21): When the state \mathbf{x}_t is represented by a probability distribution $p(\mathbf{x}_t)$, we have to compute the *expected* reward $\mathbb{E}[r(\mathbf{x}_t)] = \int r(\mathbf{x}_t)p(\mathbf{x}_t) d\mathbf{x}_t$. Moreover, we need to compute the derivatives with respect to the parameters of the state distribution, assuming that $p(\mathbf{x}_t)$ has a parametric representation.

For example, if $p(\mathbf{x}_t) = \mathcal{N}(\mathbf{x}_t | \boldsymbol{\mu}_t^x, \boldsymbol{\Sigma}_t^x)$, we compute the derivatives of $\mathbb{E}[r(\mathbf{x}_t)]$ with respect to the mean $\boldsymbol{\mu}_t^x$ and covariance $\boldsymbol{\Sigma}_t^x$ of the state distribution and continue applying the chain-rule similarly to Equation (3.20)–(3.21): With the definition $\mathcal{E}_t := \mathbb{E}_{\mathbf{x}_t}[r(\mathbf{x}_t)]$, we obtain the

gradient

$$\begin{aligned} \frac{dJ_{\boldsymbol{\theta}}}{d\boldsymbol{\theta}} &= \sum_t \gamma^t \frac{d\mathcal{E}_t}{d\boldsymbol{\theta}}, \\ \frac{d\mathcal{E}_t}{d\boldsymbol{\theta}} &= \frac{\partial\mathcal{E}_t}{\partial p(\mathbf{x}_t)} \frac{dp(\mathbf{x}_t)}{d\boldsymbol{\theta}} := \frac{\partial\mathcal{E}_t}{\partial \boldsymbol{\mu}_t^x} \frac{d\boldsymbol{\mu}_t^x}{d\boldsymbol{\theta}} + \frac{\partial\mathcal{E}_t}{\partial \boldsymbol{\Sigma}_t^x} \frac{d\boldsymbol{\Sigma}_t^x}{d\boldsymbol{\theta}}, \end{aligned} \quad (3.22)$$

where we used the shorthand notation $\partial\mathcal{E}_t/\partial p(\mathbf{x}_t) = \{\partial\mathcal{E}_t/\partial \boldsymbol{\mu}_t^x, \partial\mathcal{E}_t/\partial \boldsymbol{\Sigma}_t^x\}$ for taking the (total) derivative of \mathcal{E}_t with respect to the parameters of $p(\mathbf{x}_t) = \mathcal{N}(\mathbf{x}_t | \boldsymbol{\mu}_t^x, \boldsymbol{\Sigma}_t^x)$, i.e., the mean and covariance. The mean $\boldsymbol{\mu}_t^x$ and the covariance $\boldsymbol{\Sigma}_t^x$ are functionally dependent on the moments $\boldsymbol{\mu}_{t-1}^x$ and $\boldsymbol{\Sigma}_{t-1}^x$ of $p(\mathbf{x}_{t-1})$ and the controller parameters $\boldsymbol{\theta}$. By applying the chain-rule to Equation (3.22), we obtain

$$\frac{d\boldsymbol{\mu}_t^x}{d\boldsymbol{\theta}} = \frac{\partial\boldsymbol{\mu}_t^x}{\partial \boldsymbol{\mu}_{t-1}^x} \frac{d\boldsymbol{\mu}_{t-1}^x}{d\boldsymbol{\theta}} + \frac{\partial\boldsymbol{\mu}_t^x}{\partial \boldsymbol{\Sigma}_{t-1}^x} \frac{d\boldsymbol{\Sigma}_{t-1}^x}{d\boldsymbol{\theta}} + \frac{\partial\boldsymbol{\mu}_t^x}{\partial \boldsymbol{\theta}}, \quad (3.23)$$

$$\frac{d\boldsymbol{\Sigma}_t^x}{d\boldsymbol{\theta}} = \frac{\partial\boldsymbol{\Sigma}_t^x}{\partial \boldsymbol{\mu}_{t-1}^x} \frac{d\boldsymbol{\mu}_{t-1}^x}{d\boldsymbol{\theta}} + \frac{\partial\boldsymbol{\Sigma}_t^x}{\partial \boldsymbol{\Sigma}_{t-1}^x} \frac{d\boldsymbol{\Sigma}_{t-1}^x}{d\boldsymbol{\theta}} + \frac{\partial\boldsymbol{\Sigma}_t^x}{\partial \boldsymbol{\theta}}. \quad (3.24)$$

Note that the total derivatives $d\boldsymbol{\mu}_{t-1}^x/d\boldsymbol{\theta}$ and $d\boldsymbol{\Sigma}_{t-1}^x/d\boldsymbol{\theta}$ are known from time step $t-1$.

If all these computations can be performed in closed form, the policy gradients $d\tilde{J}_{\boldsymbol{\theta}}/d\boldsymbol{\theta}$ can be computed analytically by repeated application of the chain-rule without the need for sampling. Therefore, standard optimization techniques (e.g., BFGS or CG) can be used to learn policies with thousands of parameters [20].

3.3.3 Discussion

Using the gradients of the expected long-term reward $J_{\boldsymbol{\theta}}$ with respect to the policy parameters $\boldsymbol{\theta}$ often leads to faster learning than gradient-free policy updates. Moreover, gradient-free methods are typically limited to a few tens of policy parameters [50]. Computing the gradients can be unwieldy and requires additional computational resources. When computing gradients, exact analytic gradients are preferable over sampling-based gradients since the latter ones often suffer from large variance. These variances can even lead to slower convergence than gradient-free policy updates [7, 34]. For analytic gradients, we impose assumptions on

the differentiability of the reward function r and the transition function f .⁴ Moreover, for analytic gradients, we rely on deterministic approximate inference methods, e.g., moment matching or linearization, such that only an approximation \tilde{J}_θ to J_θ can be computed; but with the exact gradients $d\tilde{J}_\theta/d\theta$.

For updating the policy we recommend using gradient information to exploit better convergence properties. Ideally, the gradients are determined analytically without any approximations. Since this aim can only be achieved for linear systems, we have to resort to approximations, either by using sampling-based approaches or analytic approximate gradients. Sampling-based approaches are practically limited to fairly low-dimensional policy parameters $\theta \in \mathbb{R}^k$, $k \leq 50$. For high-dimensional policy parameters with $k > 50$, we recommend using analytic policy gradients if they are available.

3.4 Model-based Policy Search Algorithms with Robot Applications

In this section, we briefly describe policy search methods that have been successfully applied to learning policies for robots. We distinguish between approaches that evaluate the expected long-term reward J_θ using either sampling methods as described in Section 3.2.1 or deterministic approximate inference methods as described in Section 3.2.2.

3.4.1 Sampling-based Trajectory Prediction

Sampling directly from the learned simulator has been dealt with by a series of researchers for maneuvering helicopters [7, 51, 53] and for controlling blimps [34]. All approaches use the PEGASUS algorithm [52] to generate trajectories from the learned stochastic models.

Ng et al. [53, 51] learn models for hovering a helicopter based on locally weighted linear regression. To account for noise and model inaccuracies, this originally deterministic model was made stochastic by adding i.i.d. Gaussian (system) noise to the transition dynamics.

Unlike [53, 51], Bagnell and Schneider [7] explicitly describe uncer-

⁴In stochastic MDPs, this assumption is usually valid.

tainty about the learned model by means of a posterior distribution over a finite set of locally affine models. Trajectories are sampled from this mixture of models for learning the policy

Ko et al. [34] combine idealized parametric models with nonparametric Gaussian processes for modeling the dynamics of an autonomous blimp. The GPs are used to model the discrepancy between the nonlinear parametric blimp model and the data. Trajectories are sampled from this hybrid model when learning the policy. In the following, we discuss these policy search algorithms.

3.4.1.1 Locally Weighted Regression Forward Models and Sampling-based Trajectory Prediction

In [7], locally weighted Bayesian regression was used for learning forward models for hovering an autonomous helicopter. To account for model uncertainty, a posterior probability distribution over the model parameters ψ and, hence, the model itself was considered instead of a point estimate of the model parameters. Trajectories were sampled from this mixture of models for learning the policy.

Trajectories $\tau^{[i]}$ were generated using the PEGASUS approach [52]. At each time step, a model parameter set ψ_i was sampled from the posterior distribution $p(\psi|\mathbf{X}, \mathbf{U}, \mathbf{y}, \mathbf{x}_t, \mathbf{u}_t)$. After every transition, the dynamics model was updated with the observed (simulated) transition. After each generated trajectory $\tau^{[i]}$, the model was reset by deleting the simulated trajectory $\tau^{[i]}$ from the model [7]. For Bayes-optimal model estimators, this procedure is equivalent to sampling the model *and* sampling a full trajectory from it. Algorithm 16 summarizes how to sample trajectories from the learned model while incorporating the posterior uncertainty about the model itself. The model uncertainty is implicitly integrated out by averaging over the expected long-term rewards for all generated trajectories $\tau^{[i]}$.

In [7], a gradient-free optimization, the Nelder-Mead method, was used to update the policy parameters θ , which outperformed naive gradient-based optimization. The resulting approach learned a neural-network controller with ten parameters to hover a helicopter about a fixed point [7], see Figure 3.7(a). Extrapolation outside the range of

Algorithm 16 Policy evaluation and T -step predictions [7]

```

1: Input: transition model  $f$ , posterior distribution over model pa-
   parameters  $p(\psi|\mathbf{X}, \mathbf{U}, \mathbf{y})$ , policy parameters  $\theta$ 
2: for  $i = 1, \dots, m$  do
3:   for  $t = 0, \dots, T - 1$  do ▷ Sample trajectory  $\tau^{[i]}$ 
4:     Sample local model parameters  $\psi_i \sim p(\psi|\mathbf{X}, \mathbf{U}, \mathbf{y}, \mathbf{x}_t, \mathbf{u}_t)$ 
5:     Compute control  $\mathbf{u}_t = \pi_\theta(\mathbf{x}_t)$ 
6:     Generate a sample state transition  $\mathbf{x}_{t+1} \sim p(\mathbf{x}_{t+1}|\mathbf{x}_t, \mathbf{u}_t, \psi_i)$ 
7:     Update  $\mathbf{X}, \mathbf{U}$  with simulated transition  $(\mathbf{x}_t, \mathbf{u}_t, \mathbf{x}_{t+1})$ 
8:   end for
9:   Compute  $J_{\theta,i}$ 
10:  Reset the learned model to the original model  $f$ 
11: end for
12:  $\tilde{J}_\theta = \frac{1}{m} \sum_{i=1}^m J_{\theta,i}$ 

```



(a) Helicopter hovering [7].



(b) Inverted helicopter hovering [53].

Fig. 3.7 Model-based policy search methods with stochastic inference were used for learning to hover helicopters.

collected data was discouraged by large penalties on the corresponding states. The learned controller that was based on the probability distribution over models, was substantially less oscillatory than a controller learned by using the maximum likelihood model, i.e., a point estimate of the model parameters.

Ng et al. [53, 51] learn models for helicopter hovering based on locally-weighted linear regression, see Figure 3.7(b). Unlike in [7], a point estimate of the parameters ψ in Equation (3.3) was determined, for instance by maximum likelihood or maximum-a-posteriori estima-



Fig. 3.8 Combination of a parametric prior and GPs for modeling and learning to control an autonomous blimp [34].

tion. To account for noise and model inaccuracies, this originally deterministic model was made stochastic by adding i.i.d. Gaussian (system) noise to the transition dynamics. Angular velocities expressed in helicopter coordinates were integrated and subsequently transformed into angles in world coordinates, which made the model necessarily nonlinear. With this approach, models for helicopter hovering in a standard [53] or inverse [51] pose were determined using data collected from human pilots' trajectories.

For learning a controller with these stochastic nonlinear transition dynamics, the PEGASUS [52] sampling method was used to sample trajectories from the model. With these sampled trajectories, a Monte-Carlo estimate of the expected long-term reward was computed. A greedy hill-climbing method was used to learn the parameters θ of the policy π_{θ} , which was represented by a simplified neural network [51].

In the case of inverted helicopter hovering, a human pilot flipped the helicopter upside down. Then, the learned controller took over and stabilized the helicopter in the inverted position [51], an example of which is shown in Figure 3.7(b).

3.4.1.2 Gaussian Process Forward Models and Sampling-based Trajectory Prediction

In [34], GP forward models were learned to model the yaw-dynamics of an autonomous blimp, see Figure 3.8. The GP models were combined with an idealized parametric model of the blimp’s dynamics, i.e., the GP modeled the discrepancy between the parametric nonlinear model and the observed data. The model was trained on blimp trajectory data generated by a human flying the blimp using a remote control.

The PEGASUS approach [52] was used to sample long-term trajectories. Each new sample was incorporated into the model by updating the kernel matrix. The controller was learned using the gradient-free Nelder-Mead [47] optimization method. The four controller parameters θ were the drag coefficient, the right/left motor gains, and the slope of a policy-smoothing function. The learned controller was an open-loop controller, i.e., the controls were pre-computed offline and, subsequently, applied to the blimp. The controller based on the learned GP dynamics outperformed the optimal controller solely based on the underlying idealized parametric blimp dynamics model [34].

3.4.2 Deterministic Trajectory Predictions

In the following, we summarize policy search methods that perform deterministic trajectory predictions for policy evaluation.

3.4.2.1 Gaussian Process Forward Models and Deterministic Trajectory Prediction

The PILCO (probabilistic inference for learning control) policy search framework [20, 21] uses a GP forward model of the robot’s dynamics to consistently account for model errors. In combination with deterministic inference by means of moment matching for predicting state trajectories $p(\tau) = (p(\mathbf{x}_1), \dots, p(\mathbf{x}_T))$ PILCO computes an analytic approximation \tilde{J} to the expected long-term reward J_θ in Equation (3.2). Moreover, the gradients $d\tilde{J}/d\theta$ of the expected long-term reward with respect to the policy parameters are computed analytically. For policy learning, standard optimizers (e.g., BFGS or CG) can be used. The

Algorithm 17 PILCO policy search framework [20, 21]

Init: Sample controller parameters $\theta \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$. Apply random control signals and record data.

repeat

Learn probabilistic (GP) dynamics model using all data.

repeat

Compute $p(\mathbf{x}_1), \dots, p(\mathbf{x}_T)$ using moment matching and \tilde{J}_θ

Analytically compute policy gradients $d\tilde{J}_\theta/d\theta$

Update parameters θ (line-search in BFGS or CG).

until convergence; **return** θ^*

Apply π_{θ^*} to system (single trial/episode) and record data.

until task learned

PILCO algorithm is outlined in Algorithm 17. The algorithm is typically initialized uninformatively, i.e., the policy parameters are sampled randomly, and data is recorded from a short state trajectory generated by applying random actions. In the following, we briefly outline details about computing the long-term predictions, the policy gradients, and application of the PILCO framework to control and robotic systems.

Long-Term Predictions. For predicting a distribution $p(\tau|\pi_\theta)$ over trajectories for a given policy, PILCO iteratively computes the state distributions $p(\mathbf{x}_1), \dots, p(\mathbf{x}_T)$. For these predictions, the posterior uncertainty about the learned GP forward model is explicitly integrated out. Figure 3.9 illustrates this scenario: Let us assume, a joint Gaussian distribution $p(\mathbf{x}_t, \mathbf{u}_t)$ is given. For predicting the distribution $p(\mathbf{x}_{t+1})$ of the next state, the joint distribution $p(\mathbf{x}_t, \mathbf{u}_t)$ in the lower-left panel has to be mapped through the posterior GP distribution on the latent transition function, shown in the upper-left panel. Exact inference is intractable due to the nonlinear covariance function. Extensive Monte-Carlo sampling yields a close approximation to the predictive distribution, which is represented by the shaded bimodal distribution in the right panel. PILCO computes the mean and the variance of this shaded distribution exactly and approximates the shaded distribution by a Gaussian with the correct mean and variance as shown by the

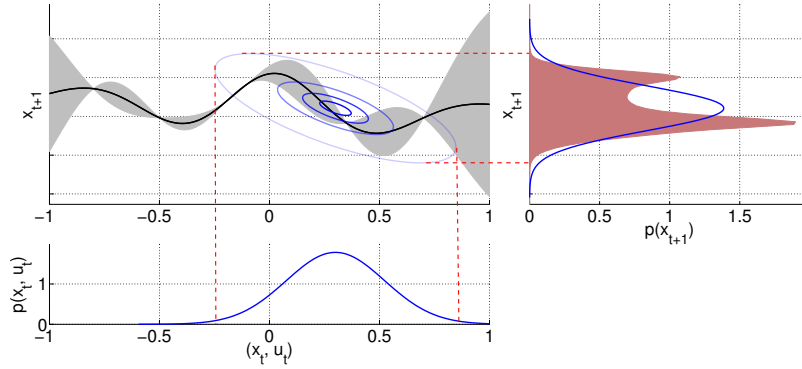


Fig. 3.9 Approximate predictions with Gaussian processes at uncertain inputs: In order to determine the predictive distribution $p(\mathbf{x}_{t+1})$, it is required to map the input distribution $p(\mathbf{x}_t, \mathbf{u}_t)$ (lower-left panel) through the posterior GP distribution (upper-left panel) while explicitly averaging out the model uncertainty (shaded area). Exhaustive Monte-Carlo sampling yields the exact distribution, represented by the red shaded distribution (right panel). The deterministic moment-matching approximation computes the mean and variance of the exact predictive distribution and fits a Gaussian (blue, right panel) with the exact first two moments to it. The contour lines in the upper-left panel represent the joint distribution between inputs and prediction.

blue distribution in the upper-right panel [20, 21].

Analytic Policy Gradients. The predicted states are not point estimates but represented by Gaussian probability distributions $p(\mathbf{x}_t)$, $t = 1, \dots, T$. When computing the policy gradients $dJ_{\theta}/d\theta$, PILCO explicitly accounts for the probabilistic formulation by analytically computing the policy gradients for probabilistic models presented in Section 3.3.2.2.⁵

Due to the explicit incorporation of model uncertainty into long-term predictions and gradient computation, PILCO typically does not suffer severely from model errors.

Robot Applications. As shown in Figure 3.11, the PILCO algorithm achieved an unprecedented speed of learning on a standard benchmark task, the under-actuated cart-pole swing-up and balancing task, see

⁵A software package implementing the PILCO learning framework is publicly available at <http://mlg.eng.cam.ac.uk/pilco>.

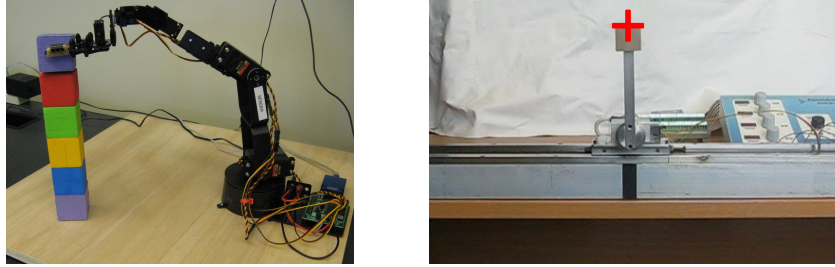


Fig. 3.10 PILCO learning successes. Left: Autonomous learning to stack a block of blocks using an off-the-shelf low-cost manipulator [21]. Right: Autonomous learning to swing up and balance a freely swinging pendulum attached to a cart [20].

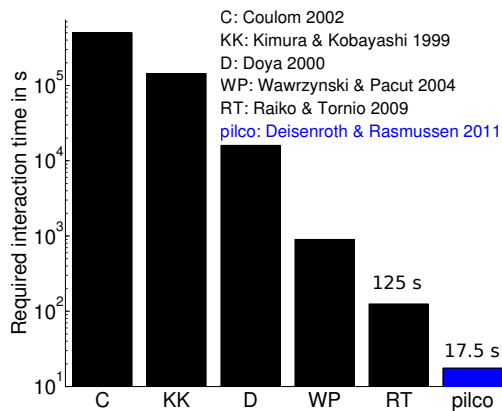


Fig. 3.11 PILCO achieves an unprecedented speed of learning on the cart-pole swing-up task. The horizontal axis gives references to RL approaches that solve the same task, the vertical axis shows the required interaction time in seconds on a logarithmic scale.

Figure 3.10. In particular, the cart-pole swing-up was learned requiring an order of magnitude less interaction time with the robot than any other RL method that also learns from scratch, i.e., without an informative initialization by demonstrations for instance. For the cart-pole swing-up problem, the learned nonlinear policy was a radial-basis function network with 50 axes-aligned Gaussian basis functions. The policy parameters θ were the weights, the locations, and the widths of the basis functions resulting in 305 policy parameters.

PILCO was also successfully applied to efficiently learning controllers from scratch in a block-stacking task with a low-cost five degrees of

Table 3.2 Overview of model-based policy search algorithms with robotic applications.

Algorithm	Predictions	Forward Model	Policy Update	Application
[7]	sampling (PEGASUS)	LWBR	gradient free	helicopter hovering
[51, 53]	sampling (PEGASUS)	LWR+noise	gradient free	helicopter hovering
[34]	sampling (PEGASUS)	GP	gradient free	blimp control
[20, 21]	moment matching	GP	gradient based	manipulator, cart pole

freedom robot arm [21], see also Figure 3.10. The state \mathbf{x} of the system was defined as the 3D coordinates of the block in the end-effector of the manipulator. For tracking these coordinates, an RGB-D camera was used. The learned policy π_{θ} was an affine function of the state, i.e., $\mathbf{u} = \mathbf{A}\mathbf{x} + \mathbf{b}$. Exactly following Algorithm 17, PILCO learned to stack a tower of six blocks in less than 20 trials. State-space constraints for obstacle avoidance were straightforwardly incorporated into the learning process as well [21].

3.4.3 Overview of Model-based Policy Search Algorithms

Table 3.2 summarizes the model-based policy search approaches that were presented in this section. Each algorithm is listed according to their prediction method (sampling/deterministic), their learned forward model (LWR/LWBR/GP), the policy updates (gradient free/gradient based), and the corresponding robotic applications.

Note that in [51, 53], model errors and local minima are dealt with by injecting additional noise to the system to reduce the danger of overfitting. Generally, noise in the system (either by artificially injecting it [51, 53] or by using probabilistic models [7, 34, 20, 21]) also smoothens out the objective function and, hence, local minima.

3.5 Important Properties of Model-based Methods

In the following, we briefly discuss three important topics related to model-based policy search. In particular, we first discuss advantages and disadvantages of stochastic inference versus deterministic inference. Then, we discuss how uncertainty about the learned model itself is treated in the literature. Finally, we shed some light on the requirements when a policy is learned from scratch, i.e., learning must happen without a good initialization. The latter point is important if neither

informed knowledge about the dynamics nor “good” data sets from demonstrations are available. Instead, the robot has to learn starting from, potentially sparse, data and uninformed prior knowledge.

3.5.1 Deterministic and Stochastic Long-Term Predictions

We discussed two general model-based approaches for computing distributions over trajectories and the corresponding long-term reward: Monte-Carlo sampling using the PEGASUS trick and deterministic predictions using linearization, unscented transformation, or moment matching. The advantage of stochastic sampling is that the sampler will return a correct estimate of the expected long-term reward J_{θ} in the limit of an infinite number of sampled trajectories. Exhaustive sampling can be computationally inefficient, but it can be straightforwardly parallelized. A more significant issue with sampling, even when using the PEGASUS approach [52], is that it is only practical for several tens of policy parameters.

As an alternative to stochastic sampling, deterministic predictions only compute an exact trajectory distribution for linear-Gaussian systems. Therefore, in nonlinear systems, only an approximation to the expected long-term reward is returned. The computations required for computing predictive distributions are non-trivial and can be computationally expensive. Unlike stochastic sampling, deterministic predictions are not straightforwardly parallelizable. On the other hand, deterministic predictions have several advantages that can outweigh its disadvantages: First, despite the fact that deterministic predictions are computationally more expensive than generating a single sample transition, the requirement of many samples quickly gets computationally even more expensive. A striking advantage of deterministic predictions is that gradients with respect to the policy parameters can be computed analytically. Therefore, policy search with deterministic prediction methods can learn policies with thousands of parameters [20].

Table 3.3 summarizes the properties of deterministic and stochastic trajectory predictions. The table lists whether the expected long-term reward J_{θ} and the corresponding gradients $dJ_{\theta}/d\theta$ can be evaluated exactly or only approximately. For stochastic trajectory predic-

Table 3.3 Properties of deterministic and stochastic trajectory predictions in model-based policy search.

	Stochastic	Deterministic
J_{θ}	exact in the limit	approximate
$dJ_{\theta}/d\theta$	exact in the limit	exact
Computations	simple	involved
# Policy parameters	$1 \leq \theta \leq 50$	$1 \leq \theta \leq ?$

tions, i.e., sampling, the required computations are relatively simple whereas the computations for deterministic predictions are mathematically more involved. Finally, we give practicable bounds on the number of policy parameters that can be learned using either of the prediction methods. For stochastic trajectory generation, J_{θ} can be evaluated exactly in the limit of infinitely many sample trajectories. The corresponding policy gradients converge even slower. In practice, where only a finite number of samples are available both J_{θ} and $dJ_{\theta}/d\theta$ cannot be evaluated exactly.

3.5.2 Treatment of Model Uncertainty

Expressing uncertainty about the learned model is important for model-based policy search to be robust to model errors. When predicting or generating trajectories, there are two general ways of treating model uncertainty.

In [72, 20, 21], model uncertainty is treated as *temporally uncorrelated* noise, i.e., model errors at each time step are considered independent. This approach is computationally relatively cheap and allows for the consideration of an infinite number of models during model averaging [20, 21]. Alternatively, sampling the model parameters initially and fixing the model parameters for the generated trajectory has the advantage that temporal correlation is automatically accounted for when the partially sampled trajectories are treated as training data until the model parameters are resampled [7]. Here, temporal correlation means that the state at one time step along a trajectory is correlated with the state at the previous time step. On the other hand, only a finite number of models can be sampled.

3.5.3 Extrapolation Properties of Models

In model-based policy search, it is assumed that models are known or have been trained in a pre-processing step [7, 51, 53, 34]. Here, humans were asked to maneuver the robot (e.g., a helicopter or a blimp) in order to collect data for model building. A crucial aspect of the collected data is that it covers the regions of the state space that are relevant for successfully learning the task at hand. Nevertheless, it is possible that it could be optimal (according to the reward function) to explore regions outside the training data of the current model. In this case, however, the learned model must be able to faithfully predict its confidence far away from the training data. Deterministic models (e.g., LWR or neural networks) cannot faithfully represent their confidence far away from the training data, which is why extrapolation is often discouraged by large penalty terms in the reward function [7]. Two models that possess credible error bars outside the training set are locally weighted Bayesian regression and Gaussian processes. Therefore, they can even be used for learning from scratch in a robotics context, i.e., without the need to ask a human expert to generate good data for model learning or a reasonably innate starting policy—if the robot is relatively robust to initial arbitrary exploration [20, 21].

3.5.4 Huge Data Sets

In robotics, it is not uncommon that huge data sets with millions of data points are available. For instance, recording 100 s of data at a frequency of 1 kHz leads to a data set with a million data points. For global models, such as the standard GP, these data set sizes lead to impractical computation time. For instance, the GP would need to repeatedly store and invert a $10^6 \times 10^6$ kernel matrix during training. A common way of reducing the size of the data set is subsampling, e.g., taking only every 10th or 100th data point. This is often possible because the dynamics are sufficiently smooth, and the state of the robot does not change much in 1/1000 s. Additionally, there are sparse approximations to the global GP [76, 63], which scale more favorably. However, even for these methods millions of data points are impractical. Therefore, local models, such as LWR or local GP models [54] should be employed if the

data sets are huge. The idea of local models is to train many models with local information, and combine predictions of these models.

4

Conclusion and Discussion

In this review, we have provided an overview of successful policy search methods in the context of robot learning, where high-dimensional and continuous state-action space challenge any RL algorithm. We distinguished between model-free and model-based policy search methods.

4.1 Conclusion

Model-free policy search is very flexible as it does not make any assumptions on the underlying dynamics. Instead of learning models, data from the robot is directly used to evaluate and update the policy. When prior knowledge in the form of demonstrations or low-level policies is available, model-free policy search often yields relatively fast convergence. In Section 2, we distinguished between the used policy evaluation strategy, policy update strategy, and exploration strategy. The policy evaluation strategies can be categorized in step-based and episode-based evaluation. While step-based exploration makes more efficient use of the sampled trajectory data, algorithms using episode-based policy evaluation strategies typically learn an upper-level policy $\pi(\omega)$ which can also capture the correlation of the parameters for an

efficient exploration. We presented four main policy update strategies, policy gradients, expectation-maximization, information-theoretic policy updates and policy updates based on path integrals.

Model-Free Methods. For the policy updates, we identified a main insight from information theory, i.e., the “distance” between the old trajectory distribution and the new trajectory distribution should be bounded, as an important key for a fast and stable learning process. This insight is used by the natural policy gradient algorithms [61, 62, 78] and by the REPS algorithm [57, 17]. While policy gradient methods require a user-specified learning rate which is not always easy to choose, REPS performs a weighted maximum likelihood (ML) estimate to determine the new policy, which can be obtained in closed form and does not require a learning rate.

EM-based methods such as PoWER [38] and methods based on path integrals [81] also employ a weighted maximum likelihood estimate. However, in contrast to REPS, those methods are also available in the step-based policy evaluation formulation, and, thus, might use data more efficiently. Hence, PoWER and PI² might show a better performance as the REPS approach in scenarios where the step-based information can be effectively exploited.

All three methods which are based on weighted ML, REPS, PoWER and PI² use a soft-max distribution to determine the weighting of the data-points. While in PoWER the temperature of the soft-max distribution is set by hand, PI² uses a heuristic which works well in practice. For the information-theoretic REPS approach, this temperature is determined by the relative entropy bound used in the algorithm and automatically recomputed for each policy update. Furthermore, episode-based REPS uses a baseline $V(\mathbf{s})$ to remove the state-dependent reward from the reward samples $R^{[i]}$. This baseline emerges naturally from the additional constraint to reproduce the given context distribution $p(\mathbf{s})$. The usage of this baseline still needs to be explored for the alternative EM and path integral approaches. Based on the beneficial properties of the information theoretic approaches, our recommendation as policy update strategy is REPS [57].

The exploration strategy creates new trajectory samples which are

subsequently used to determine the policy update. Here we identified a clear trend to use an exploration strategy in parameter space which chooses the exploration only at the beginning of the episode. Furthermore, correlated exploration is preferable to uncorrelated exploration strategies as long as the number of parameters allow for an accurate estimation of the full covariance matrix used for correlated exploration strategies.

Model-Based Methods. Model-free policy search imposes only general assumptions on the entire learning process, but the number of policy parameters we can manage is limited by the number of samples that can be generated. While tens to hundred parameters are still feasible, learning several hundreds or thousands of policy parameters seems impractical due to an excessive need for real-robot experiments.

The objective of model-based policy search is to increase the data efficiency compared to model-free methods. For this purpose, an internal model of the robot is learned that, subsequently, is used for long-term predictions and policy improvement. The learned policy is, therefore, inherently limited by the quality of the model. Thus, it is crucial to account for potential model errors during policy learning by expressing uncertainty about the learned model itself. This idea has been successfully implemented by all model-based algorithms presented in Section 3. However, model learning imposes assumptions, such as differentiability, on the robot’s forward dynamics, effectively reducing the generality of model-free policy search.

We distinguished between stochastic and deterministic approaches for trajectory predictions. While sampling-based inference is conceptually simple and can be easily parallelized, it is currently limited to successfully learn policies with several tens of parameters, similarly to model-free policy search. In cases with hundreds or thousands of policy parameters, we have to resort to deterministic approximate inference (e.g., linearization or moment matching), ideally in combination with an analytic computation of policy gradients. An example of a method with these characteristics is PILCO.

4.2 Current State of the Art

In the following, we qualitatively compare model-free and model-based approaches and discuss future trends in policy search for robotics.

Characteristics of Model-Free and Model-Based Policy Search Applications. Model-free policy search applications typically rely on a compact policy representation, which does not use more than 100 parameters. Typically, time-dependent representations are used, e.g., the Dynamic Movement Primitives [31, 71] approach, since such representations can encode movements for high-dimensional systems with a relatively small number of parameters. Most applications of model-free policy search rely on imitation learning to initialize the learning process.

When good models can be learned, model-based policy search is a promising alternative to model-free methods. Good models can often be learned when no abrupt changes in the dynamics occur, such as contacts in locomotion and manipulation. The presented model-based algorithms were applied to learning models for flying helicopters, blimps, and robot arms—in all cases, the underlying dynamics were relatively smooth.

Advantages of Model-Free and Model-Based Policy Search. Learning a policy is often easier than learning accurate forward models of the robot and its environment. In the literature, model-free policy search is the predominant approach in comparison to model-based methods since the difficulty of learning a model is avoided. Model-free methods do not place assumptions on the underlying process, e.g., the system dynamics do not have to be smooth. Hence, model-free policy search can also be applied to environments, which include discrete events such as hitting a table tennis ball. Episode-based policy search algorithms can also be used when the reward is not composed of the rewards of the intermediate steps.

Model-based policy search uses the learned model as a simulator of the robot. Therefore, model-based policy search is the predominant approach for fast and data-efficient learning: Once a model is learned,

no interaction with the real robot is required to update the policy. Moreover, policies can be tested using the model without the risk of damaging the robot. When a model is available, the time required on the robot for running experiments is negligible.

Requirements and Limitations of Model-Free and Model-Based Policy Search. Model-free policy search methods typically require that the policy can be represented with less than 100 parameters. In addition, an initialization for the policy parameters needs to be determined, e.g., by imitation learning. In addition, model-free policy search methods are inherently local search methods and might get stuck in a local optimum.

The major advantage of model-based policy search is at the same time its major limitation: the availability of the model. Before we can exploit the model as a simulator, a sufficiently good model needs to be learned from data. By explicitly describing posterior uncertainty about the learned model itself, the effect of model errors can be reduced substantially [72, 20]. Despite the fact that non-parametric models are a very rich class of models, in some way we always need to impose smoothness assumptions to the underlying system dynamics.

4.3 Future Challenges and Research Topics.

Model-free and model-based policy search methods have, so far, been developed mostly in isolation. However, the combination of model-free policy search with learned models seems to be a promising approach, a recent example is given in [41]. For example, most model-based approaches greedily exploit the learned model by using gradient-based approaches. Model-free policy update strategies could be used to avoid a greedy optimization, and, hence, the additional exploration might in the end improve the quality of the learned models. Another promising approach is to smoothly switch from model-based policy search in the initial learning phase to model-free policy search when sufficiently much data is available for model-free policy updates. In such case, the model could guide the initial exploration into relevant areas of the parameter space. For fine-tuning the policy, real trajectory samples are

used, and, hence, the policy update is not affected from model errors.

We believe that one of the most important challenges in robot learning is to incorporate structured and hierarchical learning methods into policy search. Many motor tasks are structured. For example, many tasks can be decomposed into elemental movements, also called movement primitives or options. Such a decomposition suggests a modular control approach wherein options can be adapted to the current context, activated simultaneously and sequentially in time. While there have been first approaches in model-free policy search to adapt options to the current context as well as to select options with a gating network, hierarchical approaches have so far not been explored for model-based reinforcement learning. Extending model-free methods with learned models seems promising in the field of hierarchical policy search. In general, we believe that the use of hierarchical robot control policies used in robot learning is largely unexplored and will become a new important subfield in robot learning.

In model-free policy search, we think it is important to generate a unified framework for imitation learning and reinforcement learning, such that data from both approaches can be used within one algorithm. Further goals in model-free policy search include developing a principled way to combine the advantages of step-based policy search algorithms and episode-based algorithms. This means to combine the effective use of sampled trajectory data from step-based methods with more sophisticated exploration strategies and the extensibility to hierarchical policies. Furthermore, the advantages of step-based and episode-based algorithms need to be explored in more detail, e.g., for which policy representation which type of algorithm is more useful. In addition, we believe that a principled treatment of exploration as individual objective for policy search algorithms is a promising approach to achieve a more stable learning process and avoid heuristics such as adding additional noise terms to sustain exploration.

For model-based approaches, the main challenges are choosing an appropriate model class and performing long-term predictions with this model. Non-parametric methods, such as Gaussian processes or LWBR already provide the necessary flexibility but might be difficult to scale to high-dimensional systems. In order to model discontinuities,

the use of hierarchical forward models is promising approach. The hierarchy can either be defined by the user or directly inferred from the data [27]. Once an appropriate model class has been chosen, the model can be used for long-term predictions and policy evaluation. For nonlinear model classes, approximate inference is required. Depending on the chosen model class and the number of policy parameters, we can choose between stochastic approximate inference, i.e., sampling, or deterministic approximate inference, e.g., moment matching. Since the policy representation, the learned model, and the policy update strategy are inherently connected, all these components need to fit well together in order to make the overall learning process successful.

A key challenge in learning for robots is to deal with sensory information: First, sensor data is typically noisy. Especially for model learning purposes, noisy data is challenging since not only the measurements but also the training inputs are noisy. This fact is often tacitly ignored in practice. Second, sensor data, such as images, can be high dimensional. Dimensionality reduction and feature learning can be used for a lower-dimensional compact representation of the data. Therefore, robot learning with noisy and high-dimensional data can also be phrased in the context of learning and solving partially observable Markov decision processes (MDPs) with continuous state and control spaces.

Finally, the field of robot learning needs to move to more complex applications. So far, many applications included single-stroke tasks such as hitting a baseball or catching a ball with a cup. The next step is to integrate robot learning into large-scale tasks which require the execution of a multitude of single movements. Challenging examples are given by dexterous manipulation, legged locomotion on uneven terrain, playing a full game of table tennis against a human champion, or learning to play soccer with humanoid robots. For these complex tasks, it will be necessary to exploit their modularity to simplify the learning problem. Ideally, we would automatically create re-usable submodules for policies and models that can be combined to complex policies and models.

Acknowledgments

The research leading to these results has received funding from the European Community's Seventh Framework Programme (FP7/2007–2013) under grant agreement #270327.

References

- [1] P. Abbeel, M. Quigley, and A. Y. Ng. Using Inaccurate Models in Reinforcement Learning. In *Proceedings of the 23rd International Conference on Machine Learning*, pages 1–8, Pittsburgh, PA, USA, June 2006.
- [2] E. W. Aboaf, S. M. Drucker, and C. G. Atkeson. Task-Level Robot Learning: Juggling a Tennis Ball More Accurately. In *Proceedings of the International Conference on Robotics and Automation*, 1989.
- [3] S. Amari. Natural Gradient Works Efficiently in Learning. *Neural Computation*, 10:251–276, February 1998.
- [4] B. D. O. Anderson and J. B. Moore. *Optimal Filtering*. Dover Publications, Mineola, NY, USA, 2005.
- [5] K. J. Aström and B. Wittenmark. *Adaptive Control*. Dover Publications, 2008.
- [6] C. G. Atkeson and J. C. Santamaría. A Comparison of Direct and Model-Based Reinforcement Learning. In *Proceedings of the International Conference on Robotics and Automation*, 1997.
- [7] J. A. Bagnell and J. G. Schneider. Autonomous Helicopter Control using Reinforcement Learning Policy Search Methods. In *Proceed-*

- ings of the International Conference on Robotics and Automation, pages 1615–1620. IEEE Press, 2001.
- [8] J. A. Bagnell and J. G. Schneider. Covariant Policy Search. In *Proceedings of the International Joint Conference on Artificial Intelligence*, August 2003.
 - [9] J. Baxter and P. Bartlett. Direct Gradient-Based Reinforcement Learning: I. Gradient Estimation Algorithms. Technical report, 1999.
 - [10] J. Baxter and P. L. Bartlett. Infinite-Horizon Policy-Gradient Estimation. *Journal of Artificial Intelligence Research*, 2001.
 - [11] D. P. Bertsekas. *Dynamic Programming and Optimal Control*, volume 1 of *Optimization and Computation Series*. Athena Scientific, Belmont, MA, USA, 3rd edition, 2005.
 - [12] C. M. Bishop. *Pattern Recognition and Machine Learning*. Information Science and Statistics. Springer-Verlag, 2006.
 - [13] J. A. Boyan. Least-Squares Temporal Difference Learning. In *Proceedings of the Sixteenth International Conference on Machine Learning*, pages 49–56, 1999.
 - [14] S. Boyd and L. Vandenberghe. *Convex Optimization*. Cambridge University Press, 2004.
 - [15] W. S. Cleveland and S. J. Devlin. Locally-Weighted Regression: An Approach to Regression Analysis by Local Fitting. *Journal of the American Statistical Association*, 83(403):596–610, 1988.
 - [16] R. Coulom. *Reinforcement Learning Using Neural Networks, with Applications to Motor Control*. PhD thesis, Institut National Polytechnique de Grenoble, 2002.
 - [17] C. Daniel, G. Neumann, and J. Peters. Hierarchical Relative Entropy Policy Search. In N. Lawrence and M. Girolami, editors, *Proceedings of the International Conference of Artificial Intelligence and Statistics*, pages 273–281, 2012.
 - [18] C. Daniel, G. Neumann, and J. Peters. Learning Concurrent Motor Skills in Versatile Solution Spaces. In *IEEE/RSJ International Conference on Intelligent Robots and Systems*, 2012.
 - [19] P. Dayan and G. E. Hinton. Using Expectation-Maximization for Reinforcement Learning. *Neural Computation*, 9(2):271–278, 1997.
 - [20] M. P. Deisenroth and C. E. Rasmussen. PILCO: A Model-Based

- and Data-Efficient Approach to Policy Search. In *Proceedings of the International Conference on Machine Learning*, pages 465–472, New York, NY, USA, June 2011. ACM.
- [21] M. P. Deisenroth, C. E. Rasmussen, and D. Fox. Learning to Control a Low-Cost Manipulator using Data-Efficient Reinforcement Learning. In *Proceedings of the International Conference on Robotics: Science and Systems*, Los Angeles, CA, USA, June 2011.
- [22] M. P. Deisenroth, C. E. Rasmussen, and J. Peters. Gaussian Process Dynamic Programming. *Neurocomputing*, 72(7–9):1508–1524, March 2009.
- [23] K. Doya. Reinforcement Learning in Continuous Time and Space. *Neural Computation*, 12(1):219–245, January 2000.
- [24] G. Endo, J. Morimoto, T. Matsubara, J. Nakanishi, and G. Cheng. Learning CPG-based Biped Locomotion with a Policy Gradient Method: Application to a Humanoid Robot. *International Journal of Robotics Research*, 2008.
- [25] S. Fabri and V. Kadiramanathan. Dual Adaptive Control of Non-linear Stochastic Systems using Neural Networks. *Automatica*, 34(2):245–253, 1998.
- [26] A. A. Fel'dbaum. Dual Control Theory, Parts I and II. *Automation and Remote Control*, 21(11):874–880, 1961.
- [27] E. B. Fox and D. B. Dunson. Multiresolution Gaussian Processes. In *Advances in Neural Information Processing Systems*. The MIT Press, 2012.
- [28] N. Hansen, S. Muller, and P. Koumoutsakos. Reducing the Time Complexity of the Derandomized Evolution Strategy with Covariance Matrix Adaptation (CMA-ES). *Evolutionary Computation*, 11(1):1–18, 2003.
- [29] V. Heidrich-Meisner and C. Igel. Hoeffding and Bernstein Races for Selecting Policies in Evolutionary Direct Policy Search. In *Proceedings of the 26th Annual International Conference on Machine Learning*, pages 401–408. ACM, 2009.
- [30] V. Heidrich-Meisner and C. Igel. Neuroevolution Strategies for Episodic Reinforcement Learning. *Journal of Algorithms*, 64(4):152–168, oct 2009.
- [31] A. J. Ijspeert and S. Schaal. Learning Attractor Landscapes for

- Learning Motor Primitives. In *Advances in Neural Information Processing Systems*, pages 1523–1530. MIT Press, Cambridge, MA, 2003.
- [32] S. J. Julier and J. K. Uhlmann. Unscented Filtering and Nonlinear Estimation. *Proceedings of the IEEE*, 92(3):401–422, March 2004.
- [33] H. Kimura and S. Kobayashi. Efficient Non-Linear Control by Combining Q-learning with Local Linear Controllers. In *Proceedings of the 16th International Conference on Machine Learning*, pages 210–219, 1999.
- [34] J. Ko, D. J. Klein, D. Fox, and D. Haehnel. Gaussian Processes and Reinforcement Learning for Identification and Control of an Autonomous Blimp. In *Proceedings of the International Conference on Robotics and Automation*, pages 742–747, 2007.
- [35] J. Kober, B. J. Mohler, and J. Peters. Learning Perceptual Coupling for Motor Primitives. In *Intelligent Robots and Systems*, pages 834–839, 2008.
- [36] J. Kober, K. Mülling, O. Kroemer, C. H. Lampert, B. Schölkopf, and J. Peters. Movement Templates for Learning of Hitting and Batting. In *International Conference on Robotics and Automation*, pages 853–858, 2010.
- [37] J. Kober, E. Oztop, and J. Peters. Reinforcement Learning to adjust Robot Movements to New Situations. In *Proceedings of the 2010 Robotics: Science and Systems Conference*, 2010.
- [38] J. Kober and J. Peters. Policy Search for Motor Primitives in Robotics. *Machine Learning*, pages 1–33, 2010.
- [39] N. Kohl and P. Stone. Policy Gradient Reinforcement Learning for Fast Quadrupedal Locomotion. In *Proceedings of the International Conference on Robotics and Automation*, 2003.
- [40] P. Kormushev, S. Calinon, and D. G. Caldwell. Robot Motor Skill Coordination with EM-based Reinforcement Learning. In *Proceedings of the IEEE/RSJ International Conference on Intelligent Robots and Systems*, 2010.
- [41] A. Kupcsik, M. P. Deisenroth, J. Peters, and G. Neumann. Data-Efficient Generalization of Robot Skills with Contextual Policy Search. In *Proceedings of the AAAI Conference on Artificial Intelligence*, 2013.

- [42] M. G. Lagoudakis and R. Parr. Least-Squares Policy Iteration. *Journal of Machine Learning Research*, 4:1107–1149, December 2003.
- [43] D. J. C. MacKay. *Information Theory, Inference, and Learning Algorithms*. Cambridge University Press, The Edinburgh Building, Cambridge CB2 2RU, UK, 2003.
- [44] D. C. McFarlane and K. Glover. *Lecture Notes in Control and Information Sciences*, volume 138, chapter Robust Controller Design using Normalised Coprime Factor Plant Descriptions. Springer-Verlag, 1989.
- [45] J. Morimoto and C. G. Atkeson. Minimax Differential Dynamic Programming: An Application to Robust Biped Walking. In S. Becker, S. Thrun, and K. Obermayer, editors, *Advances in Neural Information Processing Systems*. The MIT Press, 2003.
- [46] R. Neal and G. E. Hinton. A View Of The EM Algorithm That Justifies Incremental, Sparse, And Other Variants. In *Learning in Graphical Models*, pages 355–368. Kluwer Academic Publishers, 1998.
- [47] J. A. Nelder and R. Mead. A Simplex Method for Function Minimization. *Computer Journal*, 7:308–313, 1965.
- [48] G. Neumann. Variational Inference for Policy Search in Changing Situations. In *Proceedings of the 28th International Conference on Machine Learning*, pages 817–824, New York, NY, USA, June 2011. ACM.
- [49] G. Neumann and J. Peters. Fitted Q-Iteration by Advantage Weighted Regression. In *Neural Information Processing Systems*. MA: MIT Press, 2009.
- [50] A. Y. Ng. Stanford Engineering Everywhere CS229—Machine Learning, Lecture 20, 2008. <http://see.stanford.edu/materials/aimlcs229/transcripts/MachineLearning-Lecture20.html>.
- [51] A. Y. Ng, A. Coates, M. Diel, V. Ganapathi, J. Schulte, B. Tse, E. Berger, and E. Liang. Autonomous Inverted Helicopter Flight via Reinforcement Learning. In M. H. Ang Jr. and O. Khatib, editors, *International Symposium on Experimental Robotics*, vol-

ume 21 of *Springer Tracts in Advanced Robotics*, pages 363–372. Springer, 2004.

- [52] A. Y. Ng and M. Jordan. PEGASUS: A Policy Search Method for Large MDPs and POMDPs. In *Proceedings of the 16th Conference on Uncertainty in Artificial Intelligence*, pages 406–415, 2000.
- [53] A. Y. Ng, H. J. Kim, M. I. Jordan, and S. Sastry. Autonomous Helicopter Flight via Reinforcement Learning. In S. Thrun, L. K. Saul, and B. Schölkopf, editors, *Advances in Neural Information Processing Systems*, Cambridge, MA, USA, 2004. The MIT Press.
- [54] D. Nguyen-Tuong, M. Seeger, and J. Peters. Model Learning with Local Gaussian Process Regression. *Advanced Robotics*, 23(15):2015–2034, 2009.
- [55] B. Øksendal. *Stochastic Differential Equations: An Introduction with Applications (Universitext)*. Springer, 6th edition, Sept. 2010.
- [56] J. Peters, M. Mistry, F. E. Udawadia, J. Nakanishi, and S. Schaal. A Unifying Methodology for Robot Control with Redundant DOFs. *Autonomous Robots*, (1):1–12, 2008.
- [57] J. Peters, K. Mülling, and Y. Altun. Relative Entropy Policy Search. In *Proceedings of the 24th National Conference on Artificial Intelligence*. AAAI Press, 2010.
- [58] J. Peters and S. Schaal. Policy Gradient Methods for Robotics. In *Proceedings of the 2006 IEEE/RSJ International Conference on Intelligent Robotics Systems*, pages 2219–2225, Beijing, China, 2006.
- [59] J. Peters and S. Schaal. Applying the Episodic Natural Actor-Critic Architecture to Motor Primitive Learning. In *Proceedings of the European Symposium on Artificial Neural Networks*, 2007.
- [60] J. Peters and S. Schaal. Natural Actor-Critic. *Neurocomputation*, 71(7-9):1180–1190, 2008.
- [61] J. Peters and S. Schaal. Reinforcement Learning of Motor Skills with Policy Gradients. *Neural Networks*, (4):682–97, 2008.
- [62] J. Peters, S. Vijayakumar, and S. Schaal. Reinforcement Learning for Humanoid Robotics. In *IEEE-RAS International Conference on Humanoid Robots*. IEEE, September 2003.
- [63] J. Quiñonero-Candela and C. E. Rasmussen. A Unifying View of Sparse Approximate Gaussian Process Regression. *Journal of*

- Machine Learning Research*, 6(2):1939–1960, 2005.
- [64] T. Raiko and M. Tornio. Variational Bayesian Learning of Non-linear Hidden State-Space Models for Model Predictive Control. *Neurocomputing*, 72(16–18):3702–3712, 2009.
- [65] C. E. Rasmussen and C. K. I. Williams. *Gaussian Processes for Machine Learning*. Adaptive Computation and Machine Learning. The MIT Press, Cambridge, MA, USA, 2006.
- [66] L. Rozo, S. Calinon, D. G. Caldwell, P. Jimenez, and C. Torras. Learning collaborative impedance-based robot behaviors. In *AAAI Conference on Artificial Intelligence*, Bellevue, Washington, USA, 2013.
- [67] T. Rückstieß, M. Felder, and J. Schmidhuber. State-Dependent Exploration for Policy Gradient Methods. In *European Conference on Machine Learning*, pages 234–249, 2008.
- [68] T. Rückstieß, F. Sehnke, T. Schaul, D. Wierstra, Y. Sun, and J. Schmidhuber. Exploring Parameter Space in Reinforcement Learning. *Paladyn*, 1(1):14–24, March 2010.
- [69] S. J. Russel and P. Norvig. *Artificial Intelligence: A Modern Approach*. Pearson Education, 2003.
- [70] S. Schaal and C. G. Atkeson. Constructive Incremental Learning from only Local Information. *Neural Computation*, 10(8):2047–2084, 1998.
- [71] S. Schaal, J. Peters, J. Nakanishi, and A. Ijspeert. Learning Movement Primitives. In *International Symposium on Robotics Research*, pages 561–572, 2003.
- [72] J. G. Schneider. Exploiting Model Uncertainty Estimates for Safe Dynamic Control Learning. In *Advances in Neural Information Processing Systems*. Morgan Kaufman Publishers, 1997.
- [73] F. Sehnke, C. Osendorfer, T. Rückstieß, A. Graves, J. Peters, and J. Schmidhuber. Policy Gradients with Parameter-based Exploration for Control. In *Proceedings of the International Conference on Artificial Neural Networks*, 2008.
- [74] F. Sehnke, C. Osendorfer, T. Rückstieß, A. Graves, J. Peters, and J. Schmidhuber. Parameter-Exploring Policy Gradients. *Neural Networks*, 23(4):551–559, 2010.

- [75] C. Shu, H. Ding, and N. Zhao. Numerical Comparison of Least Square-Based Finite-Difference (LSFD) and Radial Basis Function-Based Finite-Difference (RBFDD) Methods. *Computers & Mathematics with Applications*, 51(8):1297–1310, April 2006.
- [76] E. Snelson and Z. Ghahramani. Sparse Gaussian Processes using Pseudo-inputs. In Y. Weiss, B. Schölkopf, and J. C. Platt, editors, *Advances in Neural Information Processing Systems 18*, pages 1257–1264. The MIT Press, Cambridge, MA, USA, 2006.
- [77] F. Stulp and O. Sigaud. Path Integral Policy Improvement with Covariance Matrix Adaptation. In *Proceedings of the 29th International Conference on Machine Learning*, 2012.
- [78] Y. Sun, D. Wierstra, T. Schaul, and J. Schmidhuber. Efficient Natural Evolution Strategies. In *Proceedings of the 11th Annual conference on Genetic and Evolutionary Computation*, pages 539–546, New York, NY, USA, 2009. ACM.
- [79] R. Sutton, D. McAllester, S. Singh, and Y. Mansour. Policy Gradient Methods for Reinforcement Learning with Function Approximation. In *Neural Information Processing Systems*, 1999.
- [80] R. S. Sutton and A. G. Barto. *Reinforcement Learning*. The MIT Press, Boston, MA, 1998.
- [81] E. Theodorou, J. Buchli, and S. Schaal. A Generalized Path Integral Control Approach to Reinforcement Learning. *Journal of Machine Learning Research*, (11):3137–3181, 2010.
- [82] S. Thrun, W. Burgard, and D. Fox. *Probabilistic Robotics*. The MIT Press, Cambridge, MA, USA, 2005.
- [83] E. Todorov. Optimal Control Theory. *Bayesian Brain*, 2006.
- [84] M. Toussaint. Robot Trajectory Optimization using Approximate Inference. In *Proceedings of the 26th International Conference on Machine Learning*, 2009.
- [85] N. Vlassis and M. Toussaint. Model-Free Reinforcement Learning as Mixture Learning. In *International Conference on Machine Learning*, page 136, 2009.
- [86] N. Vlassis, M. Toussaint, G. Kontes, and S. Piperidis. Learning Model-Free Robot Control by a Monte Carlo EM Algorithm. *Autonomous Robots*, 27(2):123–130, 2009.

- [87] P. Wawrzynski and A. Pacut. Model-Free Off-Policy Reinforcement Learning in Continuous Environment. In *Proceedings of the INNS-IEEE International Joint Conference on Neural Networks*, pages 1091–1096, 2004.
- [88] D. Wierstra, T. Schaul, J. Peters, and J. Schmidhuber. Natural Evolution Strategies. In *IEEE Congress on Evolutionary Computation*, pages 3381–3387, 2008.
- [89] R. J. Williams. Simple Statistical Gradient-Following Algorithms for Connectionist Reinforcement Learning. *Machine Learning*, 8:229–256, 1992.
- [90] B. Wittenmark. Adaptive Dual Control Methods: An Overview. In *In Proceedings of the 5th IFAC Symposium on Adaptive Systems in Control and Signal Processing*, pages 67–72, 1995.
- [91] K. Xiong, H.-Y. Zhang, and C. W. Chan. Performance Evaluation of UKF-based Nonlinear Filtering. *Automatica*, 42:261–270, 2006.

A Gradients of Frequently Used Policies

All used policies use Gaussian distributions to generate exploration. Here we state the most frequently used gradients for Gaussian policies w.r.t the mean and the covariance matrix of the Gaussian. The gradients are always stated for policies in action space. However, the policies which are defined in parameter space have of course the same gradient.

The log-likelihood of a Gaussian policy $\pi_{\theta}(\mathbf{u}|\mathbf{x}) = \mathcal{N}(\mathbf{u}|\boldsymbol{\mu}, \boldsymbol{\Sigma})$ is given by

$$\log \pi_{\theta}(\mathbf{u}|\mathbf{x}) = -\frac{d}{2} \log 2\pi - \frac{1}{2} \log |\boldsymbol{\Sigma}| - \frac{1}{2} (\mathbf{u} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{u} - \boldsymbol{\mu}).$$

Constant Mean. If the policy is given as $\pi_{\theta}(\mathbf{u}|\mathbf{x}) = \mathcal{N}(\mathbf{u}|\boldsymbol{\mu}, \boldsymbol{\Sigma})$ and $\boldsymbol{\mu}$ is part of $\boldsymbol{\theta}$, then

$$\nabla_{\boldsymbol{\mu}} \log \pi_{\theta}(\mathbf{u}|\mathbf{x}) = (\mathbf{u} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}.$$

The gradient simplifies if $\boldsymbol{\Sigma} = \text{diag}(\boldsymbol{\sigma}^2)$,

$$\nabla_{\mu_d} \log \pi_{\theta}(\mathbf{u}|\mathbf{x}) = (u_d - \mu_d)/\sigma_d^2,$$

for the d -th dimension of $\boldsymbol{\mu}$.

Linear Mean. If the policy is given as $\pi_{\theta}(\mathbf{u}_t|\mathbf{x}_t) = \mathcal{N}(\mathbf{u}|\boldsymbol{\phi}_t(\mathbf{x})^T \mathbf{M}, \boldsymbol{\Sigma})$ and \mathbf{M} is part of $\boldsymbol{\theta}$, then

$$\nabla_{\mathbf{M}} \log \pi_{\theta}(\mathbf{u}|\mathbf{x}) = (\mathbf{u} - \boldsymbol{\phi}_t(\mathbf{x})^T \mathbf{M})^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\phi}_t(\mathbf{x}),$$

or for a diagonal covariance matrix

$$\nabla_{\mathbf{m}_d} \log \pi_{\theta}(\mathbf{u}|\mathbf{x}) = (u_d - \boldsymbol{\phi}_t(\mathbf{x})^T \mathbf{m}_d) \boldsymbol{\phi}_t(\mathbf{x}) / \sigma_d^2,$$

where \mathbf{m}_d corresponds to the d -th column of \mathbf{M} .

Diagonal Covariance Matrix. For a diagonal covariance matrix $\boldsymbol{\Sigma} = \text{diag}(\boldsymbol{\sigma}^2)$ the derivative is given by

$$\nabla_{\sigma_d} \log \pi_{\theta}(\mathbf{u}|\mathbf{x}) = -\frac{1}{\sigma_d^2} + \frac{(u_d - \mu_d)^2}{\sigma_d^3}$$

Full Covariance Matrix. For representing the full covariance matrix, typically the Cholesky decomposition of the covariance matrix $\Sigma = \mathbf{A}^T \mathbf{A}$, where \mathbf{A} is an upper-triangular matrix, is used as parametrization [78]. The parametrization with the Cholesky decomposition exploits the symmetry of the covariance matrix and enforces that Σ is positive definite. The gradient of $\log \pi_{\theta}(\mathbf{u}|\mathbf{x})$ w.r.t \mathbf{A} is given by

$$\begin{aligned} \partial_{a_{i,j}} \log \pi_{\theta}(\mathbf{u}|\mathbf{x}) &= -\frac{1}{2} \partial_{a_{i,j}} \log |\mathbf{A}^T \mathbf{A}| - \\ &\quad \frac{1}{2} \partial_{a_{i,j}} (\mathbf{A}^{-T} (\mathbf{u} - \boldsymbol{\mu}))^T (\mathbf{A}^{-T} (\mathbf{u} - \boldsymbol{\mu})) \\ &= a_{i,j}^{-1} \delta_{i,j} + s_{i,j}, \end{aligned}$$

where $\delta_{i,j}$ is the dirac-delta function which is one if $i = j$ and zero elsewhere and $s_{i,j}$ is the (i, j) th element of the matrix \mathbf{S} ,

$$\mathbf{S} = (\mathbf{u} - \boldsymbol{\mu}) (\mathbf{u} - \boldsymbol{\mu})^T \mathbf{A}^{-1} \mathbf{A}^{-T} \mathbf{A}^{-1}.$$

B Weighted ML Estimates of Frequently Used Policies

We assume a data-set \mathcal{D} given in the following form

$$\mathcal{D} = \left\{ \mathbf{x}^{[i]}, \mathbf{u}^{[i]}, d^{[i]} \right\}_{i=1 \dots N},$$

where $d^{[i]}$ denotes the weighting of the i th sample. In this section we will state the solution of weighted ML estimation, i.e.,

$$\boldsymbol{\theta}^* = \operatorname{argmax}_{\boldsymbol{\theta}} \sum_{i=1}^N \log \pi_{\boldsymbol{\theta}} \left(\mathbf{u}^{[i]} | \mathbf{x}^{[i]} \right),$$

for the most frequently used policies.

For the episode-based formulation of Policy-Search, the states $\mathbf{x}^{[i]}$ are exchanged with the contexts $\mathbf{s}^{[i]}$ and the actions $\mathbf{u}^{[i]}$ are exchanged by the parameters $\boldsymbol{\theta}^{[i]}$ of the lower level controller.

Gaussian Policy, Constant Mean. Consider a policy which is given by $\pi(\mathbf{u}) = \mathcal{N}(\mathbf{u}|\boldsymbol{\mu}, \Sigma)$, i.e., we do not have a state or a context. Such a policy is for example useful to model the upper-level policy without contexts. The weighted ML-solution for $\boldsymbol{\mu}$ and Σ is given by

$$\boldsymbol{\mu} = \frac{\sum_{i=1}^N d^{[i]} \mathbf{u}^{[i]}}{\sum_{i=1}^N d^{[i]}}, \quad \boldsymbol{\Sigma} = \frac{\sum_{i=1}^N d^{[i]} (\mathbf{u}^{[i]} - \boldsymbol{\mu}) (\mathbf{u}^{[i]} - \boldsymbol{\mu})^T}{Z}, \quad (4.1)$$

where

$$Z = \frac{\left(\sum_{i=1}^N d^{[i]}\right)^2 - \sum_{i=1}^N (d^{[i]})^2}{\sum_{i=1}^N d^{[i]}}$$

is used to obtain an unbiased estimate of the covariance. The elements σ of a diagonal covariance matrix $\boldsymbol{\Sigma} = \text{diag}(\boldsymbol{\sigma})$ can be obtained by

$$\sigma_h = \frac{\sum_{i=1}^N d^{[i]} \left(\mathbf{u}_h^{[i]} - \mu_h\right)^2}{Z}. \quad (4.2)$$

Gaussian Policy, Linear Mean. The policy is given by $\pi(\mathbf{u}|x) = \mathcal{N}(\mathbf{u}|\mathbf{W}^T \boldsymbol{\phi}(x), \boldsymbol{\Sigma})$. The weighted ML-solution for \mathbf{W} is determined by the weighted pseudo-inverse

$$\mathbf{W}_{\text{new}} = (\boldsymbol{\Phi}^T \mathbf{D} \boldsymbol{\Phi})^{-1} \boldsymbol{\Phi}^T \mathbf{D} \mathbf{U}, \quad (4.3)$$

where $\boldsymbol{\Phi} = [\boldsymbol{\phi}^{[1]}, \dots, \boldsymbol{\phi}^{[N]}]$ contains the feature vectors for all samples and \mathbf{D} is the diagonal weighting matrix containing the weightings $d_i^{[i]}$. The covariance matrix $\boldsymbol{\Sigma}$ is obtained by

$$\boldsymbol{\Sigma} = \frac{\sum_{i=1}^N d^{[i]} (\mathbf{u}^{[i]} - \mathbf{W}^T \boldsymbol{\phi}(x^{[i]})) (\mathbf{u}^{[i]} - \mathbf{W}^T \boldsymbol{\phi}(x^{[i]}))^T}{Z}, \quad (4.4)$$

where Z is defined as in Eq.(4.1).

Gaussian Policy, Linear Mean, State-Dependent Variance

The policy is given by $\pi(u|x) = \mathcal{N}(u|\mathbf{w}^T \boldsymbol{\phi}(x), \boldsymbol{\phi}(x)^T \boldsymbol{\Sigma}_{\mathbf{w}} \boldsymbol{\phi}(x))$. Here, we consider only the scalar case as the multi-dimensional distribution case is more complicated and only rarely used. The weighted ML-solution for \mathbf{w} is determined by the weighted pseudo-inverse where the weights $\tilde{d}^{[i]}$ are given as the product of the state-dependent precision $\boldsymbol{\phi}(x^{[i]})^T \boldsymbol{\Sigma}_{\mathbf{w}} \boldsymbol{\phi}(x^{[i]})$ and the actual weighting $d^{[i]}$. The solution for \mathbf{w} is equivalent to Equation (4.3) where \mathbf{D} is set to the new weightings $\tilde{d}^{[i]}$.

C Derivations of the Dual Functions for REPS

In this section we will briefly discuss constraint optimization, Lagrangian multipliers and dual-functions in general. Subsequently, we will derive the dual-functions for the different REPS formulations.

Lagrangian Function. Consider the following general constraint optimization problem with equality and inequality constraints

$$\begin{aligned} \max_{\mathbf{y}} \quad & f(\mathbf{y}) \\ \text{s.t:} \quad & \mathbf{a}(\mathbf{y}) = \mathbf{0} \\ & \mathbf{b}(\mathbf{y}) \leq \mathbf{0} \end{aligned} \quad (4.5)$$

Such optimization problem can be solved by finding the saddle-points of the Lagrangian

$$L = f(\mathbf{y}) + \boldsymbol{\lambda}_1^T \mathbf{a}(\mathbf{y}) + \boldsymbol{\lambda}_2^T \mathbf{b}(\mathbf{y}). \quad (4.6)$$

The optimization problem has a local maximum if the direction of the gradient $\partial_{\mathbf{y}} f(\mathbf{y})$ is aligned with the normal of the constraints $\boldsymbol{\lambda}_1^T \partial_{\mathbf{y}} \mathbf{a}(\mathbf{y})$ and $\boldsymbol{\lambda}_2^T \partial_{\mathbf{y}} \mathbf{b}(\mathbf{y})$. Such a point can be found by differentiating the Lagrangian w.r.t \mathbf{y} and setting it to zero.

$$\partial_{\mathbf{y}} f(\mathbf{y}) = \boldsymbol{\lambda}_1^T \partial_{\mathbf{y}} \mathbf{a}(\mathbf{y}) + \boldsymbol{\lambda}_2^T \partial_{\mathbf{y}} \mathbf{b}(\mathbf{y}). \quad (4.7)$$

Dual Function. Optimizing the dual function of an optimization problem is, under certain conditions equivalent to solving the original optimization problem [14]. However, the dual function is often easier to optimize. The dual function is obtained by finding $\mathbf{y} = \mathbf{c}(\boldsymbol{\lambda}_1, \boldsymbol{\lambda}_2)$ which satisfies the saddle-point condition given in Equation (4.7). This solution is in turn set back into the Lagrangian, which results in the dual-function $g(\boldsymbol{\lambda}_1, \boldsymbol{\lambda}_2)$.

If the original optimization problem is maximized, the dual function needs to be minimized [14]. The dual function only depends on the Lagrangian multipliers and is therefore often easier to optimize. Each inequality constraint used in the original optimization problem introduces an inequality constraint for the Lagrangian multipliers. Hence,

the original optimization problem can also be solved by solving the following program

$$\begin{aligned} \min_{\boldsymbol{\lambda}_1, \boldsymbol{\lambda}_2} \quad & g(\boldsymbol{\lambda}_1, \boldsymbol{\lambda}_2) \\ \text{s.t:} \quad & \boldsymbol{\lambda}_2 \geq \mathbf{0} \end{aligned} \tag{4.8}$$

The solution for \mathbf{y} can subsequently be found by setting the Lagrangian parameters back into $\mathbf{c}(\boldsymbol{\lambda}_1, \boldsymbol{\lambda}_2)$.

Step-Based REPS. We denote $p(\mathbf{x}, \mathbf{u}) = \mu^\pi(\mathbf{x})\pi(\mathbf{u}|\mathbf{x})$ and $p(\mathbf{x}) = \sum_{\mathbf{u}} p(\mathbf{x}, \mathbf{u})$ for brevity of the derivations. To simplify the derivations, we will also write all integrals as sums. However, the derivations also hold for the formulation with the integrals. The Lagrangian for the program in Equation (2.79) with state features $\boldsymbol{\varphi}(\mathbf{x})$ is given by

$$\begin{aligned} L = & \left(\sum_{\mathbf{x}, \mathbf{u}} p(\mathbf{x}, \mathbf{u}) r(\mathbf{x}, \mathbf{u}) \right) + \eta \left(\epsilon - \sum_{\mathbf{x}, \mathbf{u}} p(\mathbf{x}, \mathbf{u}) \log \frac{p(\mathbf{x}, \mathbf{u})}{q(\mathbf{x}, \mathbf{u})} \right) \\ & + \mathbf{v}^T \sum_{\mathbf{x}'} \boldsymbol{\varphi}(\mathbf{x}') \left(\sum_{\mathbf{x}, \mathbf{u}} p(\mathbf{x}, \mathbf{u}) p(\mathbf{x}'|\mathbf{x}, \mathbf{u}) - \sum_{\mathbf{u}'} p(\mathbf{x}', \mathbf{u}') \right) \\ & + \lambda \left(1 - \sum_{\mathbf{x}, \mathbf{u}} p(\mathbf{x}, \mathbf{u}) \right), \end{aligned} \tag{4.9}$$

where η , \mathbf{v} and λ denote the Lagrangian multipliers. Rearranging terms results in

$$\begin{aligned} L = & \sum_{\mathbf{x}, \mathbf{u}} p(\mathbf{x}, \mathbf{u}) \left(r(\mathbf{x}, \mathbf{u}) - \eta \log \frac{p(\mathbf{x}, \mathbf{u})}{q(\mathbf{x}, \mathbf{u})} - \lambda - \mathbf{v}^T \boldsymbol{\varphi}(\mathbf{x}) \right. \\ & \left. + \mathbf{v}^T \sum_{\mathbf{x}'} p(\mathbf{x}'|\mathbf{x}, \mathbf{u}) \boldsymbol{\varphi}(\mathbf{x}') \right) + \eta \epsilon + \lambda. \end{aligned} \tag{4.10}$$

We substitute $V_{\mathbf{v}}(\mathbf{x}) = \mathbf{v}^T \boldsymbol{\varphi}(\mathbf{x})$. Differentiating the Lagrangian w.r.t $p(\mathbf{x}, \mathbf{u})$

$$\begin{aligned} \partial_{p(\mathbf{x}, \mathbf{u})} L = & r(\mathbf{x}, \mathbf{u}) - \eta \left(\log \frac{p(\mathbf{x}, \mathbf{u})}{q(\mathbf{x}, \mathbf{u})} + 1 \right) - \lambda \\ & - V_{\mathbf{v}}(\mathbf{x}) + \mathbb{E}_{p(\mathbf{x}'|\mathbf{x}, \mathbf{u})} [V_{\mathbf{v}}(\mathbf{x}')] = 0, \end{aligned} \tag{4.11}$$

and setting the derivative to zero yields the solution

$$p(\mathbf{x}, \mathbf{u}) = q(\mathbf{x}, \mathbf{u}) \exp\left(\frac{\delta_{\mathbf{v}}(\mathbf{x}, \mathbf{u})}{\eta}\right) \exp\left(\frac{-\eta - \lambda}{\eta}\right) \quad (4.12)$$

with $\delta_{\mathbf{v}}(\mathbf{x}, \mathbf{u}) = r(\mathbf{x}, \mathbf{u}) + \mathbb{E}_{p(\mathbf{x}'|\mathbf{x}, \mathbf{u})}[V_{\mathbf{v}}(\mathbf{x}')] - V_{\mathbf{v}}(\mathbf{x})$. Given that we require $\sum_{\mathbf{x}, \mathbf{u}} p(\mathbf{x}, \mathbf{u}) = 1$, it is necessary that

$$\exp\left(\frac{-\eta - \lambda}{\eta}\right) = \left(\sum_{\mathbf{x}, \mathbf{u}} q(\mathbf{x}, \mathbf{u}) \exp\left(\frac{\delta_{\mathbf{v}}(\mathbf{x}, \mathbf{u})}{\eta}\right)\right)^{-1}. \quad (4.13)$$

Setting Equation (4.13) into Equation (4.12) yields the closed form solution for $p(\mathbf{x}, \mathbf{u})$. Reinserting Equation (4.12) into the Lagrangian (4.10)¹

$$g(\eta, \lambda) = \eta\epsilon + \eta + \eta\lambda = \eta\epsilon + \eta \log \exp\left(\frac{\eta + \lambda}{\eta}\right). \quad (4.14)$$

As we know from Equation (4.13) that λ depends on \mathbf{v} , we substitute (4.13) to get the formulation of the dual function, which depends on η and \mathbf{v} .

$$g(\eta, \mathbf{v}) = \eta\epsilon + \eta \log \sum_{\mathbf{x}, \mathbf{u}} q(\mathbf{x}, \mathbf{u}) \exp\left(\frac{\delta_{\mathbf{v}}(\mathbf{x}, \mathbf{u})}{\eta}\right). \quad (4.15)$$

Optimizing the Dual-Function. The dual function is in *log-sum-exp* form and therefore convex in \mathbf{v} . As we have one inequality constraint in the original optimization problem, we also get an inequality constraint for the dual problem which requires that $\eta > 0$. Hence, for a given set of samples $(\mathbf{x}^{[i]}, \mathbf{u}^{[i]})$, we have to solve the following problem²

$$\begin{aligned} \min_{\eta, \mathbf{v}} \quad & \eta\epsilon + \eta \log \sum_{\mathbf{x}^{[i]}, \mathbf{u}^{[i]}} \frac{1}{N} \exp\left(\frac{\delta_{\mathbf{v}}(\mathbf{x}^{[i]}, \mathbf{u}^{[i]})}{\eta}\right), \\ \text{s.t:} \quad & \eta > 0. \end{aligned} \quad (4.16)$$

¹It is easier to just insert Equation (4.12) into the $\log p(\mathbf{x}, \mathbf{u})$ term of the Lagrangian. All other terms connected to $p(\mathbf{x}, \mathbf{u})$ cancel out.

²For numerical accuracy, it is recommendable to subtract the maximum $\delta_{\mathbf{v}}$ inside the exp and add it again outside the log, i.e.,

$$g(\eta, \mathbf{v}) = \eta\epsilon + \max \delta_{\mathbf{v}} + \eta \log \sum_{\mathbf{x}^{[i]}, \mathbf{u}^{[i]}} \frac{1}{N} \exp\left(\frac{\delta_{\mathbf{v}}(\mathbf{x}^{[i]}, \mathbf{u}^{[i]}) - \max \delta_{\mathbf{v}}}{\eta}\right)$$

Any optimizer for constraint optimization problems can be used to solve this problem, for example `fmincon` in MATLAB. This optimization can typically be performed more efficiently by providing the optimization algorithm also the derivatives of g , which are given by

$$\partial_\eta g(\eta, \mathbf{v}) = \epsilon + \log \left(\sum_i \frac{1}{N} Z_i \right) - \frac{\sum_i Z_i \delta_{\mathbf{v}}(\mathbf{x}^{[i]}, \mathbf{u}^{[i]})}{\eta \sum_i Z_i}, \quad (4.17)$$

$$\partial_{\mathbf{v}} g(\eta, \mathbf{v}) = \frac{\sum_i Z_i \left(\mathbb{E}_{p(\mathbf{x}'|\mathbf{x}^{[i]}, \mathbf{u}^{[i]})} [\boldsymbol{\varphi}(\mathbf{x}')] - \boldsymbol{\varphi}(\mathbf{x}^{[i]}) \right)}{\sum_i Z_i}, \quad (4.18)$$

with $Z_i = \exp(\delta_{\mathbf{v}}(\mathbf{x}^{[i]}, \mathbf{u}^{[i]})/\eta)$.

Dual-Function of Episode-based REPS. The derivation of the dual-function for parameter-based REPS follows the derivation given for the infinite horizon REPS. For this reason, we will only state the resulting dual-function for the contextual policy search setup and skip the derivation. The dual-function is given by

$$g(\eta, \mathbf{v}) = \eta\epsilon + \mathbf{v}^T \hat{\boldsymbol{\varphi}} + \eta \log \sum_{\mathbf{s}, \boldsymbol{\theta}} q(\mathbf{s}, \boldsymbol{\theta}) \exp \left(\frac{\delta_{\mathbf{v}}(\mathbf{s}, \boldsymbol{\theta})}{\eta} \right), \quad (4.19)$$

where $\delta_{\mathbf{v}}(\mathbf{s}, \boldsymbol{\theta}) = R(\mathbf{s}, \boldsymbol{\theta}) - \mathbf{v}^T \boldsymbol{\varphi}(\mathbf{s})$.

Dual-Function of Hierarchical REPS. The dual function of HiREPS is given by

$$g(\eta, \xi, \mathbf{v}) = \eta\epsilon + \hat{H}_q \kappa \xi + \eta \log \sum_{\mathbf{s}, \boldsymbol{\theta}} \tilde{p}(o|\mathbf{s}, \boldsymbol{\theta})^{1+\frac{\xi}{\eta}} \exp \left(\frac{\delta_{\mathbf{v}}(\mathbf{s}, \boldsymbol{\theta})}{\eta} \right), \quad (4.20)$$

where $\delta_{\mathbf{v}}(\mathbf{x}, \mathbf{u})$ is defined as for the episode-based REPS algorithm and ξ is the Lagrangian multiplier connected with the constraint which prevents an overlapping of the options.