
Context-driven Movement Primitive Adaptation

Master-Thesis von Daniel Wilbers aus Rheine

Tag der Einreichung:

1. Gutachten: Prof. Jan Peters, Ph.D.
2. Gutachten: M.Sc. Rudolf Lioutikov



TECHNISCHE
UNIVERSITÄT
DARMSTADT



Context-driven Movement Primitive Adaptation

Vorgelegte Master-Thesis von Daniel Wilbers aus Rheine

1. Gutachten: Prof. Jan Peters, Ph.D.
2. Gutachten: M.Sc. Rudolf Lioutikov

Tag der Einreichung:

Für meine Eltern,
die immer für mich da sind.

Erklärung zur Master-Thesis

Hiermit versichere ich, die vorliegende Master-Thesis ohne Hilfe Dritter nur mit den angegebenen Quellen und Hilfsmitteln angefertigt zu haben. Alle Stellen, die aus Quellen entnommen wurden, sind als solche kenntlich gemacht. Diese Arbeit hat in gleicher oder ähnlicher Form noch keiner Prüfungsbehörde vorgelegen.

Darmstadt, den 18. Juli 2016

(Daniel Wilbers)

Abstract

Humanlike robot skills, like cleaning a table or handing over a plate, can often be generalized to different task variations. Usually, these are start-/goal position, and trained environment changes. We investigate how to modify motion primitives to context changes, which are not included in the training data. Specifically, we focus on maintaining humanlike motion characteristics and generalizability. Therefore, we present an optimization technique, which allows to maximize the expected return while minimizing the Kullback-Leibler Divergence to the demonstrations. Simultaneously, our algorithm learns how to linearly combine the adapted primitive with the demonstrations, such that only relevant parts of the primitive are adapted. We evaluate our approach in obstacle avoidance and broken joint scenarios in simulation, as well as on a real robot.

Zusammenfassung

Menschenähnliche Fertigkeiten von Robotern, wie beispielsweise einen Tisch abzuwischen oder einen Teller zu überreichen sind meistens generalisierbar zu leichten Veränderungen einer Aufgabe. Typischerweise sind diese Veränderungen unterschiedliche Anfangs- und Endpositionen, sowie zuvor gelernte Umgebungsänderungen.

Wir untersuchen, wie sich Bewegungsprimitive an Veränderungen des Kontexts, welche nicht durch Trainingsdaten abgedeckt werden, anpassen lassen. Dabei konzentrieren wir uns auf die Beibehaltung der menschlichen Charakteristiken und Anpassbarkeit der Primitive. Wir eine Optimierungsmethode vor, welche sowohl die erwartete Belohnung maximiert, als auch die Kullback-Leibler Divergenz zu den Demonstrationen minimiert. Gleichzeitig ist unser Algorithmus in der Lage eine lineare Kombination des angepassten Primitives mit den Demonstrationen zu lernen, sodass nur relevante Abschnitte verändert werden.

Wir evaluieren unseren Ansatz in Hindernisumgehungs- und defekten-Gelenk Szenarien, sowohl in der Simulation, als auch auf einem echten Roboter.

Contents

1. Introduction and Related Work	2
2. Primitive Optimization for Context-Adaptation	4
2.1. Notation	4
2.2. Problem Statement	4
2.3. Approximation with Samples	6
2.4. Combination of Primitives	7
3. Analysis: Different Aspects of our Approach	9
3.1. Primitive Combination	9
3.2. Staying close to Demonstrations	10
3.3. Tuning γ	10
4. Evaluation in Simulation and on a Real Robot	13
4.1. Hole-Reaching Task	13
4.1.1. Obstacle Avoidance	13
4.1.2. Broken Joint Scenario	13
4.2. Table-Cleaning Task	14
5. Conclusion	18
Bibliography	19
A. Appendix	21
A.1. Gaussian Distribution	21
A.2. Kullback-Leibler Divergence	21
A.3. Derivation of Blending	22
A.4. Derivation of our Algorithm without Sub-Policies	23
A.5. Detailed Derivation of our full Algorithm	26
A.5.1. Derivation of the Dual	27
A.5.2. Expectation Approximation	28

Figures

List of Figures

1.1. Illustration of the Table-Cleaning task. The robot should avoid the obstacle with the sponge staying on the table. The dashed black arrow denotes a demonstrated trajectory. The solid red arrow illustrates the context-adapted primitive, which the robot will follow to avoid the flowers on the table.	3
2.1. The demonstrated policy $\pi_d(w)$ is represented as the blue shaded area with mean and two times standard deviation. In green, the optimized policy $\pi^*(w)$ avoids the obstacle, while maintaining the shape of the distribution. In red, the combination of both policies $\pi^+(w)$ avoids the obstacle, but also exactly matches $\pi_d(w)$ in the beginning and end. The corresponding activation function, shown in the second plot, is parametrized with a difference of sigmoid functions and learned accordingly as the sub-policy $\lambda(a)$	8
3.1. Linear combination of two primitives, which move up and right respectively, to achieve new behavior. a) - b): Simultaneous constant activation a over time results in a diagonal movement. Depending on the magnitude of the activation a , the resulting primitives reach further into the upper right. c) - e): Examples for the effect of different activation functions, which are changing over time. Depending on the characteristics of the activation function the combined primitive can represent completely different behavior.	9
3.2. Comparison between REPS ($\gamma = 0$ case) and our algorithm. Each image contains 50 red ellipses which denote gaussian distributions after optimization. N_i denotes the number of samples per iteration. I is the number of iterations. The black ellipse denotes the target distribution. a) - c): REPS randomly is biased to one option or maintains a wide distribution over all options. d) - f): Our algorithm stays close to the target distribution. With a small sample set $N_i = 10$ the solutions are biased towards one option. With enough samples a single option is found near the target distribution.	10
3.3. a) - f): Solutions of the sub-policy $\pi(w)$ in parameter space for six γ variations. All red ellipses contain 90% of the probability mass. The black ellipses denote the target distribution $\pi_d(w)$. The red ellipses are the found solutions $\pi^*(w)$ after 20 iterations. All experiments are initialized with a much larger variance than the solutions. g): KL-Divergence $\mathcal{D}(\pi(w) \pi_d(w))$ during the iterative optimization for different γ	11
3.4. Comparison of the effects of γ on the optimal sub-policies $\pi^*(w)$. The blue distribution is the target policy $\pi_d(w)$. Sub-Policy $\lambda(a)$ is independent of the reward function and neglected here. The gray rectangle is a restricted area with undesirable parameters.	12
4.1. Obstacle Avoidance in the Hole-Reaching task. The robot must reach the bottom of the hole without touching the obstacle or ground. b) and c): Mean trajectories of $\pi^+(w)$	14
4.2. Broken joint scenario in the Hole-Reaching task. b): Optimized Mean Trajectory $\pi^+(w)$ with learned activations, so that the solution is close to the demonstrations in the beginning. c): Standard REPS solution $\pi^*(w)$ with a different start position since no activations are learned. The corresponding solutions in joint space are given in Figure 4.3. . .	15
4.3. The five joint dimensions in the broken joint case of the Hole-Reaching task. The working joints are forced to change in order to compensate for the broken 5th joint. The corresponding solutions in task space are given in Figure 4.2. . .	15
4.4. KL-Divergence $\mathcal{D}(\pi^*(w) \pi_d(w))$ during optimization in the broken joint case of the Hole-Reaching task. The KL-Divergence of the bounded policy case $\gamma = 50$ is lower than the one of the unbounded policy.	16
4.5. Obstacle Avoidance in the Table-Cleaning task. a): The adapted policy $\pi^+(w)$ locally avoids the obstacle and exactly matches the demonstrated policy $\pi_d(w)$ when possible. b): Policy $\pi^*(w)$ is only activated in regions where obstacle avoidance is necessary.	17

Abbreviations, Symbols and Operators

List of Symbols

Notation	Description
a	vector of activation parameters
a_i	activations for each basis function
A	diagonal matrix with activation factors for each basis function
w_i	weight for basis function i
w	vector of basis function weights
κ	context to which a robot skill needs to be adapted
μ	mean of a gaussian distribution
Σ	covariance matrix of a gaussian distribution
ϵ_1, ϵ_2	constraint parameters, which limit the exploration during optimization
$\alpha_1, \alpha_2, \eta_1, \eta_2$	lagrangian multipliers
γ	tuning parameter, which determines how much the target distribution is incorporated
I	number of iterations
N_i	number of samples per iteration
$\pi(w)$	policy, which hierarchically represents a distribution over trajectories
$\pi_d(w)$	target policy learned from the demonstrations
$\pi^*(w)$	optimal policy after optimization, which limits $\mathcal{D}(\pi^*(w) \pi_d(w))$
$\pi^+(w)$	combined policy between $\pi^*(w)$ and $\pi_d(w)$
$\lambda(a)$	policy of activations, a can also be parameterized
$q_a(a)$	old estimate of policy $\lambda(a)$, used as a sampling distribution
$q_w(w)$	old estimate of policy $\pi(w)$, used as a sampling distribution
y_t	state of trajectory at time t
τ	trajectory

List of Operators

Notation	Description	Operator
\exp	exponential function	$\exp(\bullet)$
e	exponential function	e^\bullet
\mathbb{E}	expectation	$\mathbb{E}[\bullet]$
Var	variance	$\text{Var}[\bullet]$
\mathcal{N}	gaussian distribution	$\mathcal{N}(\bullet)$
\mathcal{D}	KL-Divergence	$\mathcal{D}(\bullet \bullet)$
\log	natural logarithm	$\log(\bullet)$
R	reward function	$R(\bullet)$

List of Abbreviations

Notation	Description
KL	Kullback-Leibler divergence
s.t.	subject to

1 Introduction and Related Work

Humanlike motions of humanoid robots are desirable in many different areas. In human-robot collaboration settings it is important that the robots intentions are clear to the human. The legibility of robot motions is a requirement for a seamless integration of the robot as a helpful assistant [8]. In settings, where a robot arm is attached to a human as an exoskeleton [11] it is mandatory that the robots movements are humanlike. In this paper, we investigate how to maintain humanlike characteristics in motion primitives, which are learned from demonstrations of a human teacher.

Learning from demonstrations (also known as imitation learning) is a well-established approach for programming robots. While optical motion capturing techniques are very popular, they also suffer from the correspondence problem [10]. The humans and robot body may significantly vary, which makes it hard to match the movements of a teacher to the joints of a robot. To leverage the problematic, *kinesthetic teaching* has been successfully applied in various scenarios to learn an initial policy [5, 9, 16]. Instead of passively observing movements, a human teacher guides the robot arm by moving it by hand. The recorded demonstrations can then be used to extract movement characteristics and learn a policy, which represents the robot skill. We assume, that the reenacted robot trajectories produced by the learned skill from demonstrations naturally inherit humanlike characteristics.

In general, using expert demonstrations limits the search space and therefore can significantly bootstrap the policy search. However, directly using demonstrations of a specific task to successfully teach a robot might not be sufficient. Further policy improvement can be necessary for successfully solving a task, e.g. to compensate for the teachers actuation [16].

In comparison to common Motion Planning techniques like STOMP [14], CHOMP [26], and RRTs [21] we do not focus on generating completely new motion plans, but rather on adjusting existing ones. In our case we also do not try to generate single trajectories for a specific task instance, but optimize a distribution over trajectories. We use these trajectory distributions to represent a motion primitive, which is a solution to variations of a demonstrated task. Our primitive representation relates to the Probabilistic Movement Primitives [23] framework, but our approach should be applicable to other primitive frameworks as well, e.g. Dynamic Movement Primitives [13].

A key issue of robot skill teaching is the ability of primitives to generalize to different situations. For example a primitive which moves a chess figure forward should be applicable to all tiles on a chess board [5]. Often, the generalization ability corresponds to the adaptation to new start and/or goal positions, while preserving the trajectories shape [27, 15, 24]. Particularly, this means that the primitives can only be adapted to specific pre-trained changes.

Here, we distinguish between two forms of generalizability. The ones described so far are specifically engineered or trained, whereas in this paper we focus on modifying skills to completely new contexts, e.g. new obstacles or broken joints. Various reinforcement learning approaches have been successfully applied to handle different contexts. Work by Pastor et.al. [24] encodes goal-parameters in the primitive representation itself [24], whereas different approaches learn a distribution over meta-parameters [15, 20]. In the latter, variations of the context must be present in the demonstration, so that they can be learned. An advantage of these approaches is, that the adaptation itself is fast to compute, once learned. No further task-relearning is necessary as long as the context changes are covered in the demonstrations.

At the same time, a primitive must be relearned if an adaptation is not possible. In this case, any policy improvement method could be applied, depending on the requirements.

In this paper, we aim to adapt a primitive to an unseen situation in such a way, that it is still applicable to the context which it was already generalized to. An example is adapting a primitive to avoid an obstacle while maintaining possible start and goal positions. In relation to our work, Lim et al. [22] propose a framework, in which motion primitives are learned from demonstrations with PCA and combined to produce humanlike motions, which are assessed qualitatively. More recently, Huang et. al. [12] introduced a graph- and sample-based motion planner, which can combine two different motions. In their case, multiple upper and lower body motion primitives are coordinated and sequenced to produce a motion graph. The humanlike characteristics arise only from blending between demonstrated motions. Instead of sequencing multiple primitives in order to adapt to a new situation, we investigate how to partially change a single primitive. Closely related to our work Ye and Alterovitz combine motion planning with imitation learning to find task solutions [28]. These can be outside the demonstrations if the demonstrations are blocked. This approach focus on automatically extracting time-dependent task constraints, which are satisfied even after optimization. In our approach we assume the task-constraints are already captured in the given primitives. We punish violations of the task-constraints inside a corresponding reward function.

Being able to re-learn skills through primitive-combination and/or primitive-optimization, such that they can adapt to completely

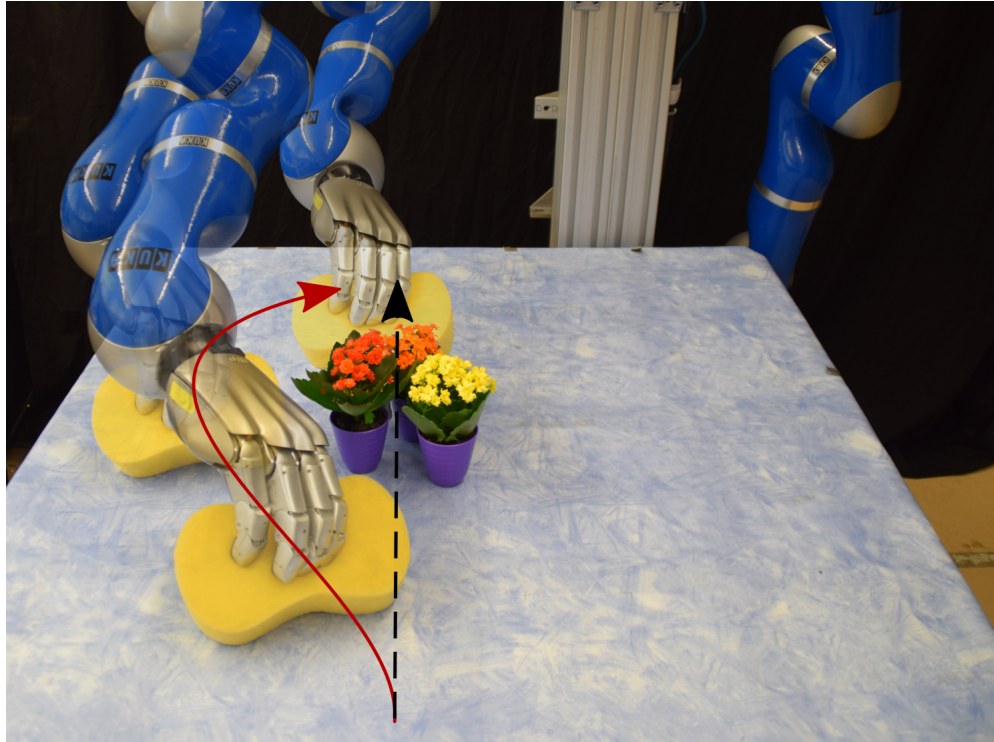


Figure 1.1.: Illustration of the Table-Cleaning task. The robot should avoid the obstacle with the sponge staying on the table. The dashed black arrow denotes a demonstrated trajectory. The solid red arrow illustrates the context-adapted primitive, which the robot will follow to avoid the flowers on the table.

new situations is essential for building and extending skill databases. Our contribution is that we optimize primitives to unseen situations, while binding the solution to stay close to the demonstrations. We also learn how to linearly combine two primitives in order to exactly match the demonstrated primitive whenever possible in the new context.

2 Primitive Optimization for Context-Adaptation

We develop a reinforcement learning strategy, based on the Relative Entropy Policy Search [25] algorithm. The two main differences are the simultaneous optimization of two sub-policies and the explicit minimization of the KL-Divergence [19] to a target distribution of one sub-policy.

2.1 Notation

Throughout this paper a primitive is defined as a policy $\pi(w)$. The policy represents a gaussian distribution over parameters $w \sim \mathcal{N}(\mu_w, \Sigma_w)$. Realizations of w can be used to generate a trajectory $\tau(w)$, so that by placing a distribution over w we hierarchically define a distribution over trajectories. In our case we use S radial basis functions to approximate trajectories $\tau(w) = \{y_1, \dots, y_T\}$ as

$$y_t(w) = \sum_i^S w_i \exp(-k(t - c_i)^2), \quad (2.1)$$

where c_i defines the basis center and k modulates the basis width. In order to adapt a trajectory, we learn how to adjust the weights w . Hence, any optimization of a trajectory distribution is equal to adjusting the distribution over w .

In relation to the REPS formalism, we assume that our policy is a joint distribution, which can be split into $\pi(w)$ and $\lambda(a)$. We refer to these as sub-policies of a gaussian joint-policy

$$\mathcal{N} \left(\begin{bmatrix} \mu_w \\ \mu_a \end{bmatrix}, \begin{bmatrix} \Sigma_w & 0 \\ 0 & \Sigma_a \end{bmatrix} \right) = \pi(w)\lambda(a). \quad (2.2)$$

By doing so, we can improve two different parameter sets, while both jointly determine a reward function $R(w, a)$. Specifically, we use $\pi(w)$ to represent a trajectory distribution and $\lambda(a) = \mathcal{N}(\mu_a, \Sigma_a)$ as distribution over activation parameters (see Section 2.4). The distribution $\pi_d(w)$ is a target distribution, to which we want to stay close. In general, the target distribution could be arbitrary. In our case it is the policy learned from demonstrations. We denote the context to which we want to adapt the primitive with κ . In the obstacle avoidance scenario κ represents the obstacles position and shape parameters.

2.2 Problem Statement

Given a target policy $\pi_d(w)$ and an unseen context κ find an optimized joint policy $(\pi(w)\lambda(a))^*$, which maximizes the expected reward of $R(w, a, \kappa)$ while minimizing the KL-divergence $\mathcal{D}(\pi(w)||\pi_d(w))$. The full optimization problem is given as

$$\begin{aligned} \max_{\pi, \lambda} \quad & J = \int_w \int_a R(w, a, \kappa) \pi(w) \lambda(a) da dw \\ & - \gamma \mathcal{D}(\pi(w)||\pi_d(w)) \\ \text{s.t.} \quad & \int_w \pi(w) \log \left(\frac{\pi(w)}{q_w(w)} \right) dw \leq \epsilon_1, \\ & \int_a \lambda(a) \log \left(\frac{\lambda(a)}{q_a(a)} \right) da \leq \epsilon_2, \\ & \int_w \pi(w) dw = 1, \\ & \int_a \lambda(a) da = 1. \end{aligned} \quad (2.3)$$

The distributions q_w and q_a represent the current estimates of the sub-policies, from which we can sample. The first two constraints ϵ_1 and ϵ_2 limit the exploration of the policy update by limiting the KL-Divergence between the current estimate and the new policy. These constraints prevent that the policy is destroyed by an too excessive update step [25]. The last two constraints make sure that each sub-policy is a probability distribution and sum up to one. Solving the optimization problem (see Appendix A.5) with the method of lagrangian multipliers yields the sub-policy update rules

$$\begin{aligned}\pi^*(w) &= \frac{q_w(w)^{\frac{\eta_1}{\gamma+\eta_1}} \pi_d(w)^{\frac{\gamma}{\gamma+\eta_1}} \exp\left(\frac{R(w)}{\gamma+\eta_1}\right)}{\int_w q_w(w)^{\frac{\eta_1}{\gamma+\eta_1}} \pi_d(w)^{\frac{\gamma}{\gamma+\eta_1}} \exp\left(\frac{R(w)}{\gamma+\eta_1}\right) dw}, \\ \lambda^*(a) &= \frac{q_a(a) \exp\left(\frac{R(a)}{\eta_2}\right)}{\int_a q_a(a) \exp\left(\frac{R(a)}{\eta_2}\right) da},\end{aligned}\tag{2.4}$$

with the terms $R(a)$ and $R(w)$ representing the expected reward for specific parameters w and a given the other sub-policy

$$\begin{aligned}R(w) &= \int_a R(w, a, \kappa) \lambda(a) da, \\ R(a) &= \int_w R(w, a, \kappa) \pi(w) dw.\end{aligned}\tag{2.5}$$

Due to the recursive dependencies respectively on the other sub-policy we can only approximate $R(a)$ and $R(w)$. Both terms measure how good the parameters w and a perform locally. Given that we iteratively update our policies and restrict the KL-Divergence between update steps we can locally test our parameters against the policies from the previous iteration.

The update of the unrestricted sub-policy $\lambda^*(a)$ is equal to the standard REPS formulation [25]. It is an exponential re-weighting of the old sampling distribution. The restricted sub-policy $\pi^*(w)$ is a geometric average of the sampling distribution, the target distribution, and the exponential returns. By setting $\gamma = 0$ we obtain the standard REPS formulation and both updates have the same form. The parameters η_1 and η_2 are lagrangian multipliers and can be viewed as temperature. The dual formulation leads to the optima η_1, η_2 respectively and is given as

$$\begin{aligned}g(\eta_1, \eta_2) &= \\ &- \mathbb{E}[R(w, a, \kappa)]_{\pi^* \lambda^*} + \eta_1 \epsilon_1 + \eta_2 \epsilon_2 \\ &+ (\gamma + \eta_1) \log \left(\int_w q_w(w) [e^{R(w)} \pi_d(w)^\gamma q_w(w)^{-\gamma}]^{\frac{1}{\gamma+\eta_1}} dw \right) \\ &+ \eta_2 \log \left(\int_a q_a(a) e^{\left(\frac{R(a)}{\eta_2}\right)} da \right).\end{aligned}\tag{2.6}$$

Solving Equation 2.3 reduces to minimizing the dual function

$$\begin{aligned}\underset{\eta_1, \eta_2}{\text{minimize}} \quad & g(\eta_1, \eta_2) \\ \text{s.t.} \quad & \eta_i \geq 0, i = 1, 2,\end{aligned}\tag{2.7}$$

which is much easier to optimize.

2.3 Approximation with Samples

In the following we explain how to approximate the integrals and the expectation term in the dual with samples. We estimate the expectation $\mathbb{E}[R(w, a, \kappa)]_{\pi^* \lambda^*}$, which is part of the original formulation Equation 2.3, with importance sampling, using $q_w(a)$ and $q_a(a)$ as the sampling distribution. The resulting approximation with importance weights $\psi(w_i, a_i)$ and N samples is

$$\mathbb{E}[R(w, a, \kappa)]_{\pi^* \lambda^*} \approx \frac{\frac{1}{N} \sum_i R(w_i, s_i, \kappa) \psi(w_i, a_i)}{\frac{1}{N} \sum_i \psi(w_i, s_i)}, \quad (2.8)$$

$$\psi(w_i, a_i) = \left[e^{R(w_i)} \pi_d(w_i)^\gamma q_1(w_i)^{-\gamma} \right]^{\frac{1}{\gamma + \eta_1}} e^{\left(\frac{R(a_i)}{\eta_2} \right)}.$$

We can further replace the integrals inside both logarithms of Equation 2.6 using samples based on $\int_y p(y) f(y) dy \approx \frac{1}{N} \sum_{i=1}^N f(y_i)$. Given these approximations the dual is purely sample-based and can be solved with any constraint optimizer. For numerical stability we suggest to rewrite the sample-based dual according to the log-sum-exp¹ and exp-normalize² identities. In our case we exclusively use gaussians for all policies, so that the new means and variances can be computed with closed-form reward-weighted maximum likelihood updates of the samples s_i [7] as

$$\mu_{\text{new}} = \frac{\sum_{i=1}^N \phi_i s_i}{\sum_{i=1}^N \phi_i},$$

$$\Sigma_{\text{new}} = \frac{\sum_{i=1}^N \phi_i (s_i - \mu)(s_i - \mu)^T}{Z}, \quad (2.9)$$

$$\text{with } Z = \frac{\left(\sum_{i=1}^N \phi_i \right)^2 - \sum_{i=1}^N (\phi_i)^2}{\sum_{i=1}^N \phi_i}.$$

The term Z is used to calculate an unbiased covariance-estimate. Algorithm-Box 1 gives an overview of the iterative procedure.

Algorithm 1: Primitive Optimization

Input: Context Situation κ , Target Policy $\pi_d(w)$

Output: Optimal sub-policies $\pi^*(w)$ and $\lambda^*(a)$

while not converged **do**

begin Policy Evaluation

Sampling:

 Generate N sample pairs

$w_i \sim \pi(w)$, $a_i \sim \lambda(a)$

Evaluation:

 Compute rewards $R(w_i, a_i, \kappa)$ for each sampled pair i

 Approximate $R(w_i, \kappa)$ and $R(a_i, \kappa)$

Optimize: Minimize dual:

$(\eta_1^*, \eta_2^*) = \text{argmin } g(\eta_1, \eta_2)$

end

begin Policy Improvement

Sub-policy updates $\pi^*(w)$ and $\lambda^*(a)$: Compute weights ϕ for weighted ML

$\phi_w(w_i) = \left[\exp(R(w_i)) \pi_d(w_i)^\gamma q_w(w_i)^{-\gamma} \right]^{\frac{1}{\gamma + \eta_1}}$

$\phi_a(a_i) = \exp\left(\frac{R(a_i)}{\eta_2}\right)$

 For gaussian distributions:

 Compute improved μ_w, Σ_w and μ_a, Σ_a (see Equation 2.9)

end

end

¹ $\log\left(\sum_i \exp(x_i)\right) = k + \log \sum_i \exp(x_i - k)$ with $k = \max x_i$

² $\frac{\sum_{i=1}^N a_i \exp(z_i)}{\sum_{i=1}^N \exp(z_i)} = \frac{\sum_{i=1}^N a_i \exp(z_i + k)}{\sum_{i=1}^N \exp(z_i + k)}$ with $k = -\max(z_i)$

In our case we use $\pi(w)$ and $\pi_d(w)$ for representing trajectory distributions and $\lambda(a)$ for activation parameters, these distributions could be arbitrary gaussian distributions for a completely different setting.

2.4 Combination of Primitives

Combining primitives to generate more complex behavior, which can solve new tasks is a requirement for building motion libraries. In the following, we describe how we combine primitives in a way, such that we maintain humanlike characteristics. Staying close to demonstrations and maintaining humanlike characteristics can be seen from two different perspectives. We can either maintain the shape of a trajectory distribution or exactly match the trajectory distribution from demonstrations whenever possible. To achieve the latter we suggest to linearly combine different policies. An illustration of both approaches is given in Figure 2.1. We make use of some gaussian properties to combine skills. First, the sum of two gaussian random variables $x \sim \mathcal{N}(\mu_x, \Sigma_x)$ and $y \sim \mathcal{N}(\mu_y, \Sigma_y)$ is a gaussian distribution $z \sim \mathcal{N}(\mu_x + \mu_y, \Sigma_x + \Sigma_y)$. Second, the affine transformation $y = c + Bx$ of a gaussian x is again a gaussian distribution $y \sim \mathcal{N}(c + B\mu, B\Sigma B^T)$. For the combined policy $\pi^+(w)$ we get

$$\begin{aligned} \mathcal{N}(\mu_{\pi^+}, \Sigma_{\pi^+}) &= \mathbf{A}\mathcal{N}(\mu_{\pi^*}, \Sigma_{\pi^*}) + (1 - \mathbf{A})\mathcal{N}(\mu_{\pi_d}, \Sigma_{\pi_d}) \\ \text{with } \mu_{\pi^+} &= \mathbf{A}\mu_{\pi_d} + (1 - \mathbf{A})\mu_{\pi^*} \\ \Sigma_{\pi^+} &= \mathbf{A}\Sigma_{\pi_d}\mathbf{A}' + (1 - \mathbf{A})\Sigma_{\pi^*}(1 - \mathbf{A})'. \end{aligned} \tag{2.10}$$

The combination itself is performed in the weight space of the trajectories. In comparison to the blending approach from [23] we are still able to generate smooth trajectory samples after the combination since we maintain a complete distribution in the weight space. The elements of the diagonal matrix $\mathbf{a} = \text{diag}(\mathbf{A})$ represent the activation factors for each basis weight of the policy. We obtain a combined trajectory distribution $\pi^+(w)$. Typically, we place a probability distribution over the activations directly, so that $\dim(a) = \dim(w)$ or parametrize the activations. For example, if we temporarily want to switch the active robot skill, a useful parametrization would be the difference of two sigmoid functions (e.g. in Figure 2.1). By doing so we can significantly reduce the number of dimensions. On the other side we must take prior knowledge into account which may also limit the performance. In relation to our algorithm we can optimize the activations represented as the sub-policy $\lambda(a)$.

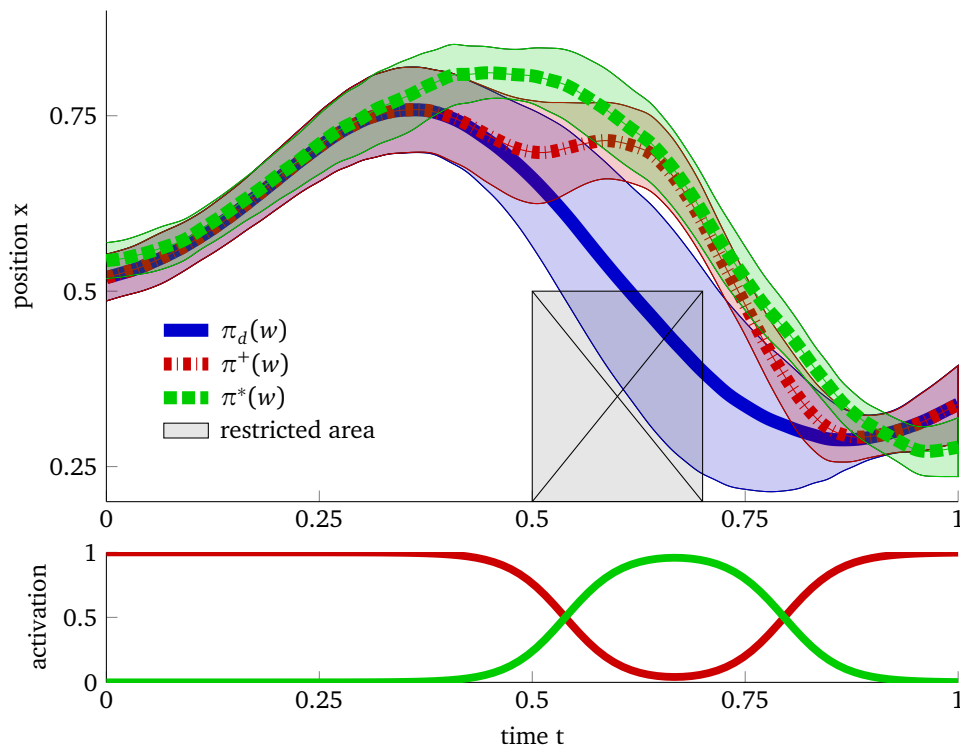


Figure 2.1.: The demonstrated policy $\pi_d(w)$ is represented as the blue shaded area with mean and two times standard deviation. In green, the optimized policy $\pi^*(w)$ avoids the obstacle, while maintaining the shape of the distribution. In red, the combination of both policies $\pi^+(w)$ avoids the obstacle, but also exactly matches $\pi_d(w)$ in the beginning and end. The corresponding activation function, shown in the second plot, is parametrized with a difference of sigmoid functions and learned accordingly as the sub-policy $\lambda(a)$.

3 Analysis: Different Aspects of our Approach

In the following section we individually discuss and demonstrate various aspects of our algorithm.

3.1 Primitive Combination

With our linear combination approach, we can achieve a different behavior, based on how we choose the activations. We neglect optimizing a trajectory distribution here and only focus on the combination itself. Assume we are given two primitives which we want to combine to achieve new behavior. In Figure 3.1 we show such two planar primitives, which can move up and right respectively. None of them can reach the upper right corner by itself, whereas it is possible with the combination. Simply taking the average as in the $a = 0.5$ case is not enough to fully exploit the combination. The blending approach from [23] fails in this case because it follows the regions with smallest variance. Blending can be useful for fulfilling task-constraints, but not for exploring completely new behavior. The examples in Figure 3.1c to Figure 3.1e show further combinations, yielding completely different results.

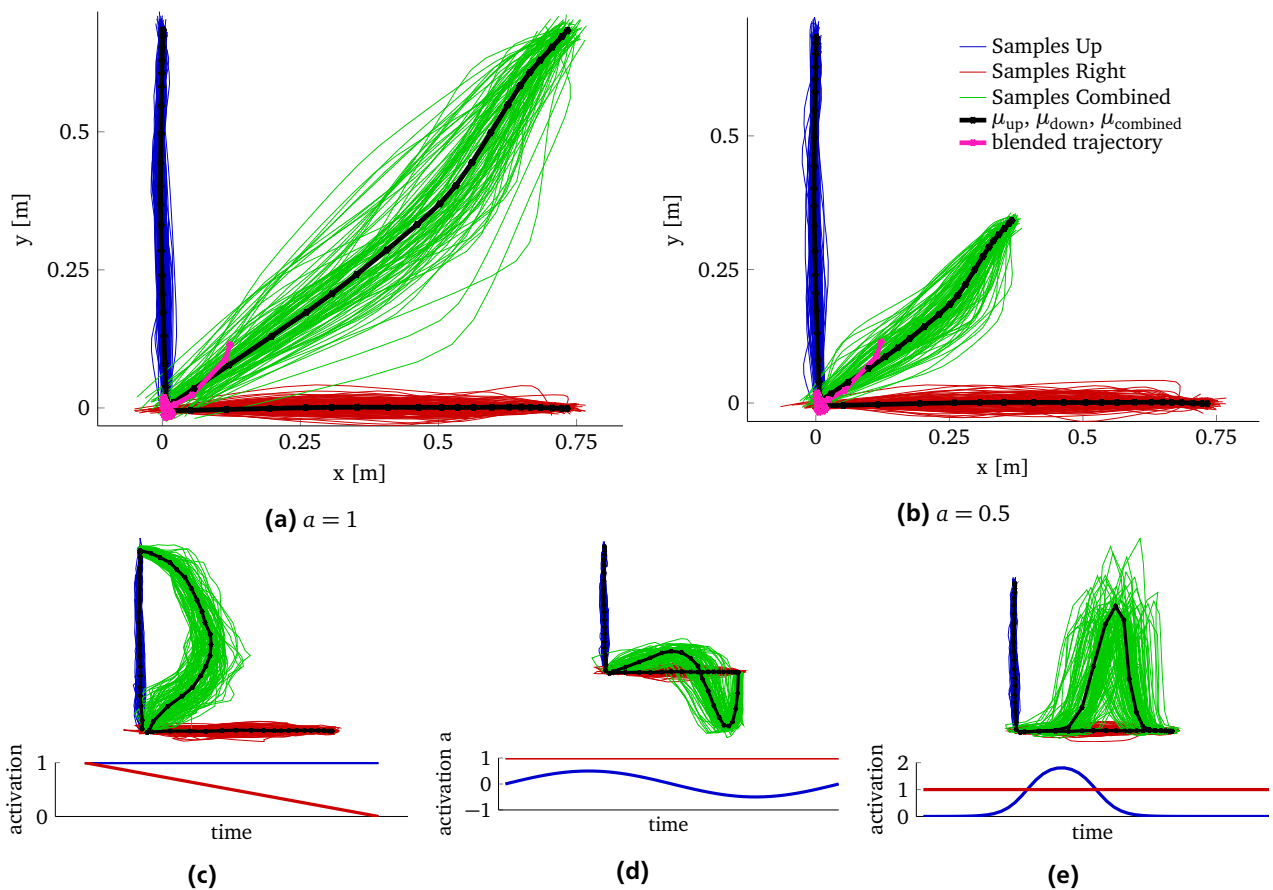


Figure 3.1.: Linear combination of two primitives, which move up and right respectively, to achieve new behavior. a) - b): Simultaneous constant activation a over time results in a diagonal movement. Depending on the magnitude of the activation a , the resulting primitives reach further into the upper right. c) - e): Examples for the effect of different activation functions, which are changing over time. Depending on the characteristics of the activation function the combined primitive can represent completely different behavior.

3.2 Staying close to Demonstrations

To emphasize the effects of minimizing the KL-Divergence to a target distribution we give a simplified example with a highly multimodal reward function, where all optima are equally good. In Figure 3.2 the bright yellow areas represent a high reward, whereas blue areas are much worse. We now assume that a demonstrated distribution in the upper left area is not available anymore. In order to show the effects of our approach sub-policy $\lambda(a)$ does not influence the reward. As we can see, limiting the KL-Divergence yields solutions which are closer to the black target distribution in the upper left. If we compare the solutions on the right side (Figure 3.2c and Figure 3.2f) we see that with enough iterations we can also match the targets variance and hence find a single option.

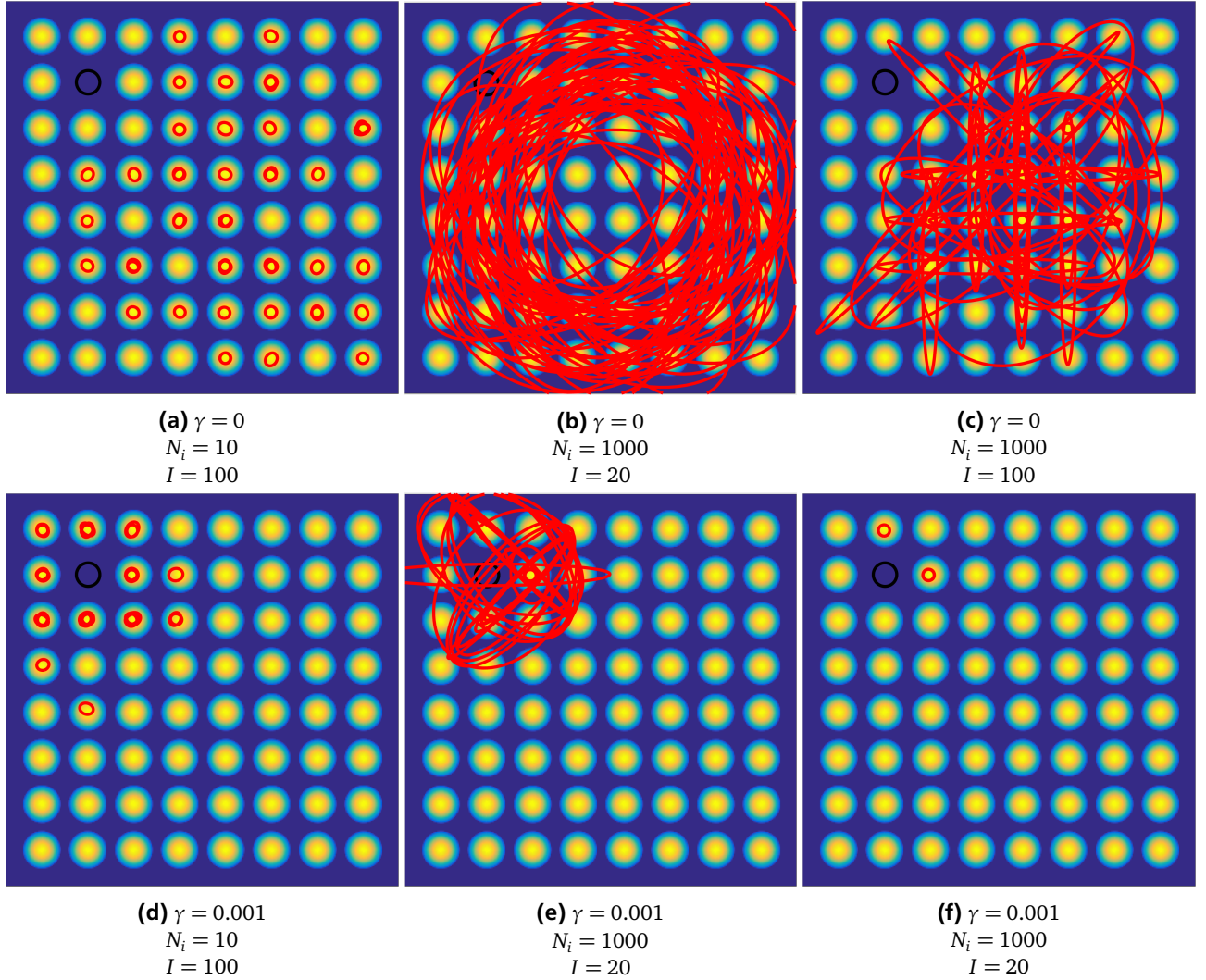


Figure 3.2.: Comparison between REPS ($\gamma = 0$ case) and our algorithm. Each image contains 50 red ellipses which denote gaussian distributions after optimization. N_i denotes the number of samples per iteration. I is the number of iterations. The black ellipse denotes the target distribution.

- a) - c): REPS randomly is biased to one option or maintains a wide distribution over all options.
- d) - f): Our algorithm stays close to the target distribution. With a small sample set $N_i = 10$ the solutions are biased towards one option. With enough samples a single option is found near the target distribution.

3.3 Tuning γ

To emphasize the effects of our algorithm we demonstrate how the solutions with our algorithm are changing when we modify the γ -parameter (Equation 2.3). Basically, γ influences how close to the target distribution we want to be. Setting γ to zero cancels the

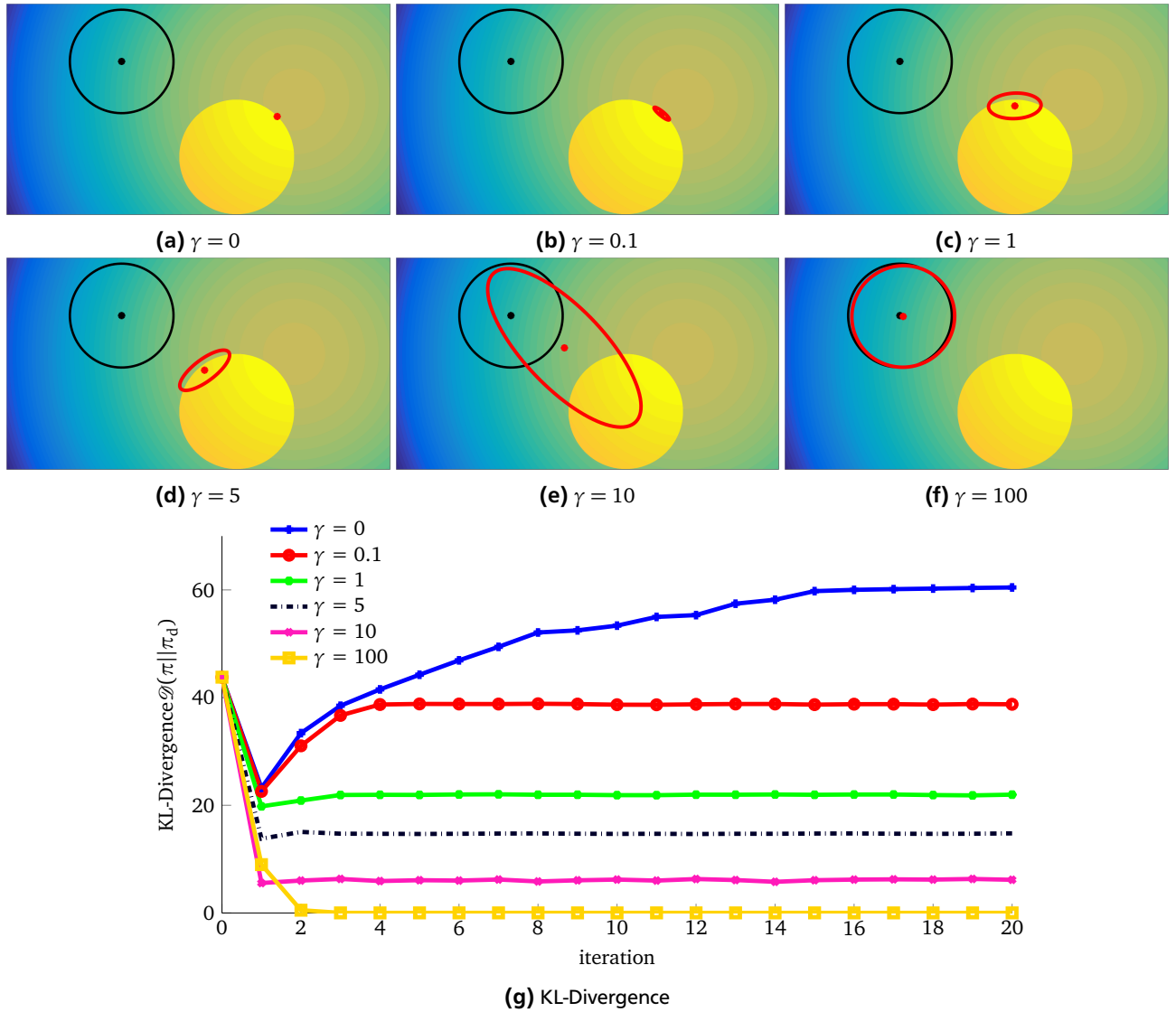


Figure 3.3.: a) - f): Solutions of the sub-policy $\pi(w)$ in parameter space for six γ variations. All red ellipses contain 90% of the probability mass. The black ellipses denote the target distribution $\pi_d(w)$. The red ellipses are the found solutions $\pi^*(w)$ after 20 iterations. All experiments are initialized with a much larger variance than the solutions. g): KL-Divergence $\mathcal{D}(\pi(w)||\pi_d(w))$ during the iterative optimization for different γ .

effect of the KL-Divergence in the optimization. A characteristic of our approach is that we are comparing reward and KL-divergence against each other, which must not necessarily be of the same magnitude. Depending on the reward function tuning γ in an adjusted parameter range is necessary. In Figure 3.3 we give a simplified example based on a two-dimensional reward function. The toy function is quadratic with an additional circular discontinuity. In order to show the effects of limiting the KL-Divergence sub-policy $\lambda(a)$ has no effect on the reward. The higher the γ -value the more the solutions get pulled towards the target distribution. In the $\gamma = 0$ case the result reduces to REPS and reaches the global optimum of the reward-function. We see the effect in Figure 3.3g: The higher γ the lower the KL-Divergence. Notably, the reward is lower the higher γ is. When optimizing we need to cope with this trade-off. As shown in Section 3.2 if many local optima are available the reward can still be equally good.

In relation to trajectory distributions Figure 3.4 shows a simplified example based on a one-dimensional trajectory distribution $\pi(w)$ to directly show the effects on the solution.

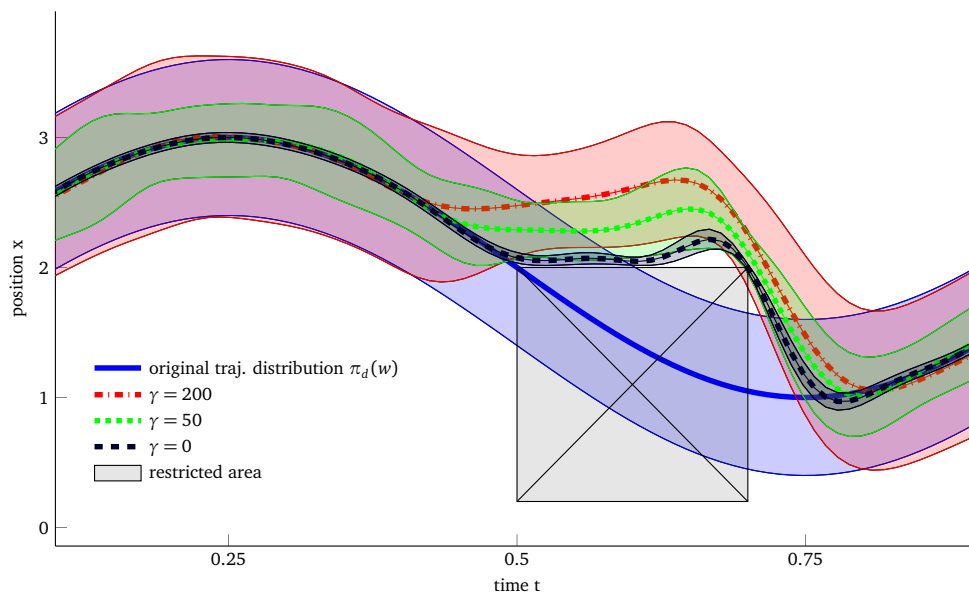


Figure 3.4.: Comparison of the effects of γ on the optimal sub-policies $\pi^*(w)$. The blue distribution is the target policy $\pi_d(w)$. Sub-Policy $\lambda(a)$ is independent of the reward function and neglected here. The gray rectangle is a restricted area with undesirable parameters.

4 Evaluation in Simulation and on a Real Robot

We demonstrate different aspects of our algorithm and apply it to various problems. Therefore we evaluate settings in task space, as well as in joint space. In addition to simulations, we also illustrate the execution on a 7-DOF real robot arm.

4.1 Hole-Reaching Task

We apply our algorithm on a planar 5-DOF robot. Each link is one meter long and has no joint limits. The robot's end-effector must reach the bottom of a hole, which is two meters away, one meter deep, and 30-60cm wide. With ten basis functions for each degree of freedom we have a 50-dimensional parameter vector for sub-policy $\pi(w)$. Additionally, we learn the activations a without parameterizing them, but assume all dimensions have the same activation. In total, we optimize 60 parameters. The reward depends on a collision cost, acceleration punishment, and a reward for reaching the goal position. We initialize our optimization with the policy learned from the demonstrations. In the demonstrations the robot always starts from a roughly upright position and reaches the bottom of the hole.

Obtaining a suitable covariance matrix is especially challenging in this task. Due to the high dimensional parameter space and limited number of available samples per iteration, the weighted-ML updates during the optimization are most likely biased. Because this task is operating in joint space it is crucial that we obtain a proper estimate of the covariances between joints. The joints at the beginning of the joint chain highly influence the suitable angles of the following joints. Due to this high correlation a precise estimate of the covariance matrix is necessary for finding successful solutions. Inspired by the CECER approach [1], which in practice is hard to tune, we modify the covariance matrix estimate Σ_w after the weighted-ML update. We use a convex combination of the current estimate Σ_i with the covariance matrix Σ_{i-1} from the last iteration to limit the covariance shrinkage $\Sigma_w^{\text{new}} = \delta \Sigma_w^{i-1} + (1 - \delta) \Sigma_w^i$. No further adjustments were needed for Σ_a . We apply this setting in two different context scenarios.

4.1.1 Obstacle Avoidance

In the first scenario an obstacle is added to the scene. The robot is supposed to reach the 60 cm wide hole without colliding with the obstacle. As it can be seen in Figure 4.1 we can successfully learn such a behavior. The optimization is performed in joint space. Thus, we use forward kinematics to calculate collisions with the obstacle, which is given in task space. The obstacle is one-by-one meter square block. We compare the $\gamma = 0$ case with the best working $\gamma \geq 0$ case ($\gamma = 10^{-7}$). Furthermore, the activation function was only allowed to be either one or zero and was clipped at the first and last value to stay close to the demonstrations. By doing so, we ensure that we specifically learn where to switch between the adapted skill and the one from demonstrations. We don't need to encode a start or goal position inside the reward function, as it's maintained through the combination with demonstrations. Both cases were executed with identical parameters except γ . We repeated the experiment five times with each 40 iterations and 250 samples per iteration. Following the ideas of [14], the covariance matrix Σ_w was initialized with a scaled version of the one learned from demonstrations. Therefore only the variances belonging to the middle of the trajectory were scaled up. Additionally, for this task $\delta = 0.9$ was found practical. After we adapted the skill, we tested our results by sampling from the learned trajectory distributions. Sampling 500 trajectories from the original distribution resulted in only 0.008% collision free ones. In the $\gamma = 0$ case 0.45% were collision free, whereas our approach succeeded in average 68% of the time, which is close to one standard deviation.

4.1.2 Broken Joint Scenario

In the second scenario we assume that the last link of the robot is broken and cannot move anymore. Therefore the skill must be adapted, so that the robot can still reach the hole. We enforce parameters of the broken joint dimension to be zero. The reward function is equal to the setting above, but without any additional obstacle. The target distribution is the same distribution, which was learned from demonstrations with all joints working. Figure 4.2b shows our results after 30 iterations with each 2500 samples. In comparison to the REPS case $\pi^*(w)$ with $\gamma = 0$ and without activation, our solution still maintains the start position (see Figure 4.3) and still is able to move the arm into the hole. Figure 4.4 shows the KL-Divergences of both cases.

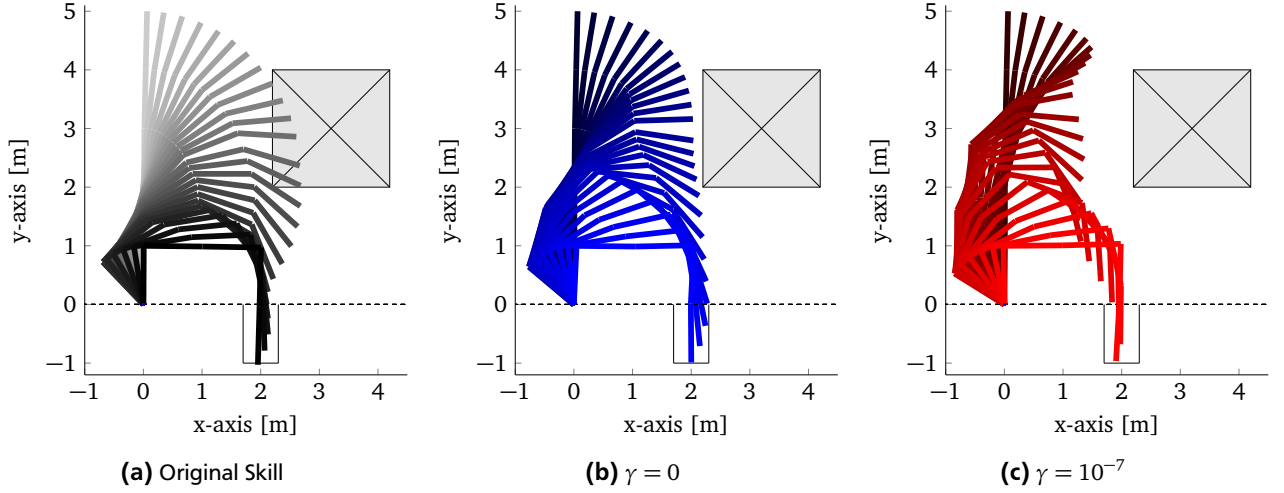


Figure 4.1.: Obstacle Avoidance in the Hole-Reaching task. The robot must reach the bottom of the hole without touching the obstacle or ground. b) and c): Mean trajectories of $\pi^+(w)$.

4.2 Table-Cleaning Task

We also tested our approach on a 7-DOF real robot arm. The robot is supposed to pick up a sponge and wipe multiple times over a table. We demonstrate the task three times via kinesthetic teaching with an empty table. We extract the contact points with the table from the learned trajectory distribution to specify the task constraints. During optimization these points should still be in contact with the table. Our goal is to reenact the skill even if items are still present on the table as illustrated in Figure 1.1. We optimize in joint space, but currently only check for collisions of the end-effector in task space. With 50 basis functions per degree of freedom and additional 50-activation parameters we optimize a 400 dimensional parameter vector in total. The activation parameters are clipped to either zero or one. As shown in Figure 4.5 the activations are learned and the obstacles are avoided accordingly. The solution was found after 30 iterations with each 500 samples. Unlike the necessary adjustments of the covariance matrix Σ_w in the Hole-Reaching Task, it was not required for the Table-Cleaning even if the number of parameters is much higher. In the obstacle avoidance case of the Hole-Reaching task a wide range of the trajectory distribution needs to be changed to successfully avoid the obstacle. In comparison to that, in this specific instance of the Table-Cleaning task only a limited region of the trajectory distribution needs to be changed in order to avoid the obstacle. During the optimization this is learned by the sub-policy $\lambda(a)$, such that the optimization can locally concentrate on adapting the important regions. Hence, for the solution $\pi^+(w)$ most of the parts of difficult to estimate covariance matrix Σ_w are not considered. Due to that valid solution can still be found even if Σ_w is biased.

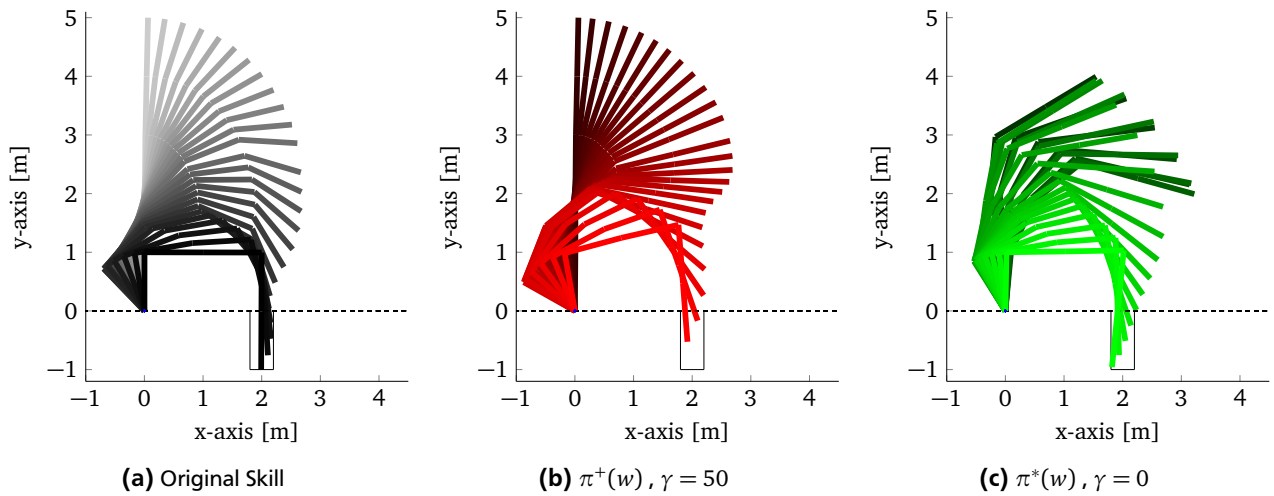


Figure 4.2.: Broken joint scenario in the Hole-Reaching task.

b): Optimized Mean Trajectory $\pi^+(w)$ with learned activations, so that the solution is close to the demonstrations in the beginning. c): Standard REPS solution $\pi^*(w)$ with a different start position since no activations are learned.

The corresponding solutions in joint space are given in Figure 4.3.

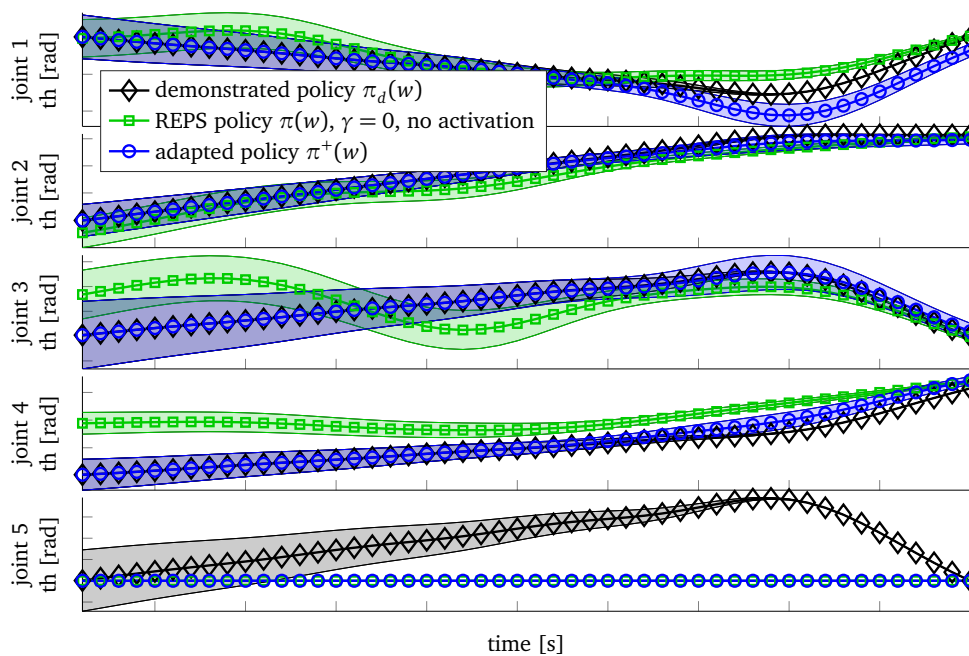


Figure 4.3.: The five joint dimensions in the broken joint case of the Hole-Reaching task. The working joints are forced to change in order to compensate for the broken 5th joint. The corresponding solutions in task space are given in Figure 4.2.

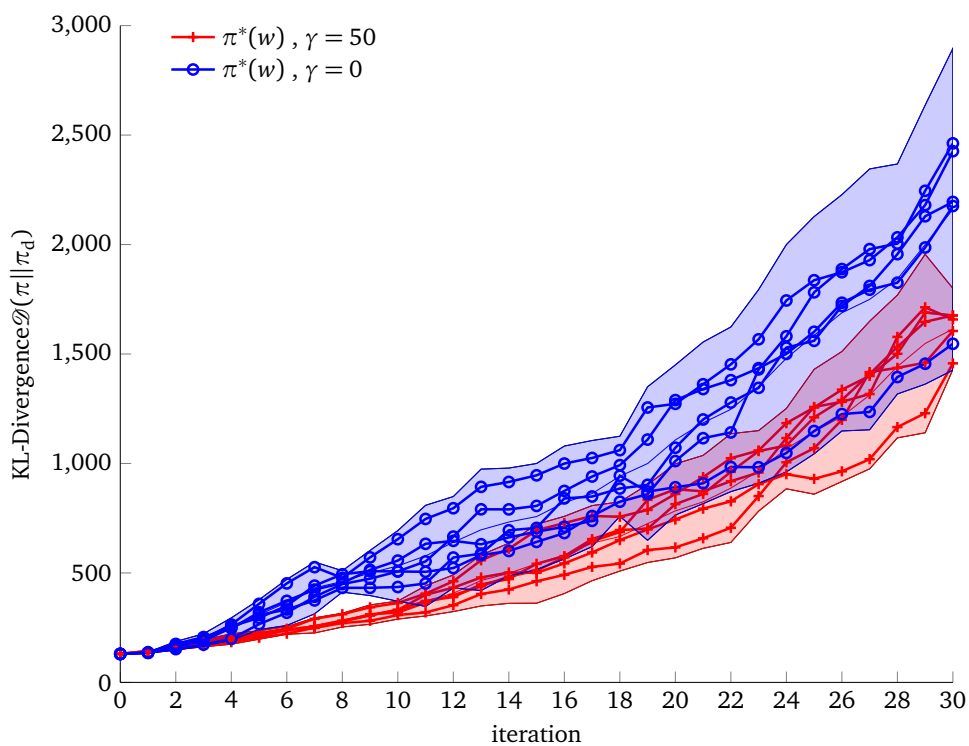
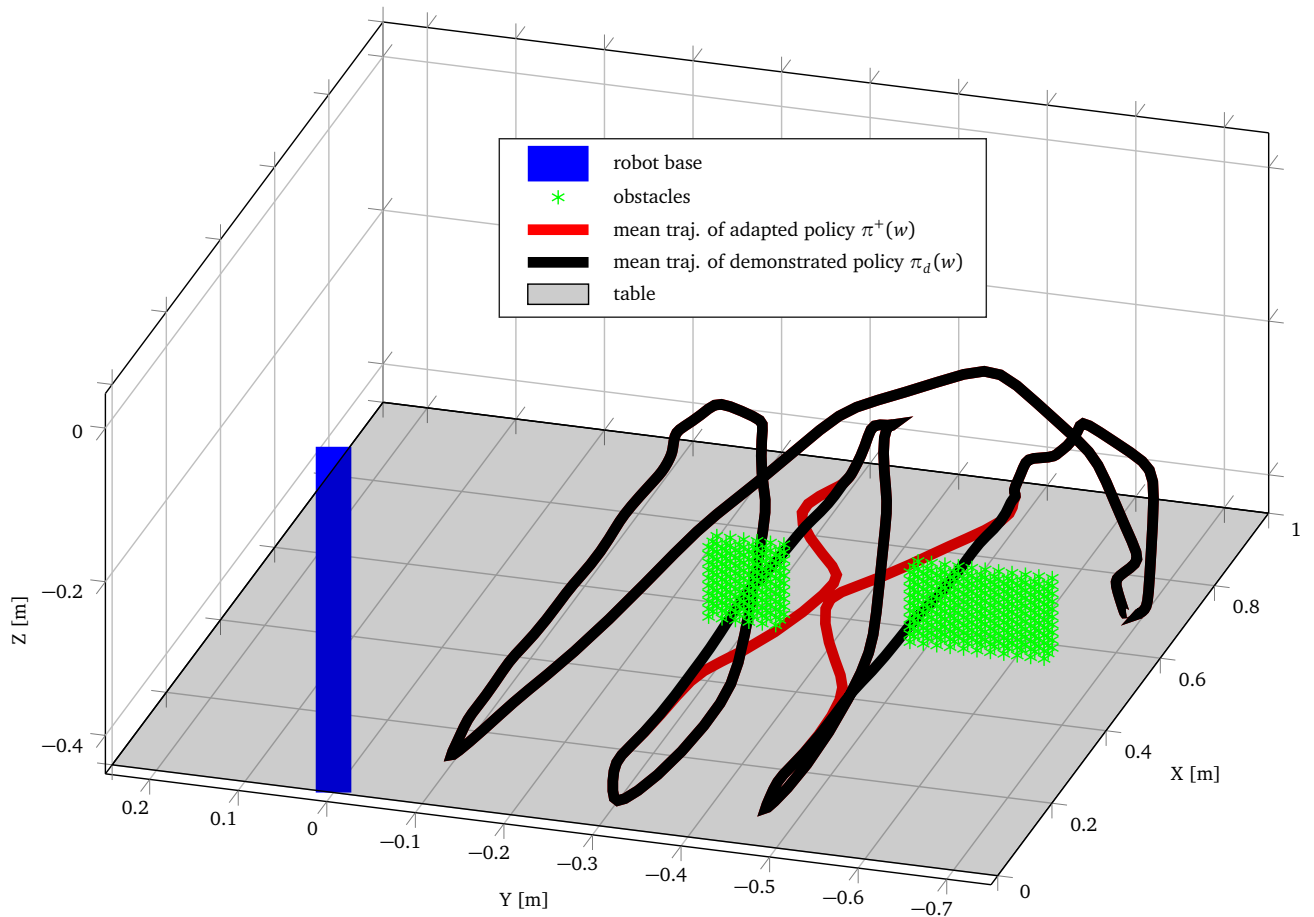
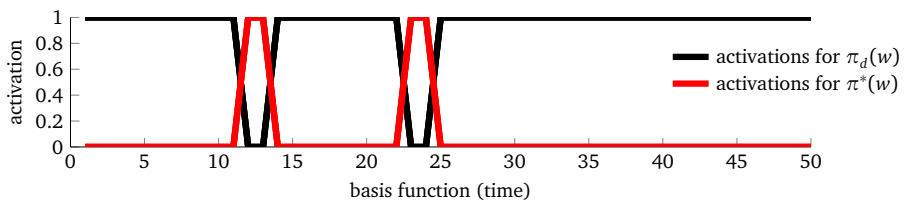


Figure 4.4.: KL-Divergence $\mathcal{D}(\pi^*(w)||\pi_d(w))$ during optimization in the broken joint case of the Hole-Reaching task. The KL-Divergence of the bounded policy case $\gamma = 50$ is lower than the one of the unbounded policy.



(a) Task-space trajectories of the Table-Cleaning task



(b) Learned Activations

Figure 4.5.: Obstacle Avoidance in the Table-Cleaning task.

a): The adapted policy $\pi^+(w)$ locally avoids the obstacle and exactly matches the demonstrated policy $\pi_d(w)$ when possible.

b): Policy $\pi^*(w)$ is only activated in regions where obstacle avoidance is necessary.

5 Conclusion

In this paper, we presented an approach for adapting robot skills to new situations. Our algorithm makes use of three different aspects. The first one is, that we simultaneously optimize the expected reward over two sub-policies, instead of one single policy. Second, we explicitly bind the KL-Divergence of one of them to maintain humanlike motion characteristics. And third, we learn a linear combination with the demonstrations to exactly match the demonstrated policy where possible.

As our results show, we are able to produce trajectories, which are adapted to a context and still stay close to the demonstrated trajectories.

In the future, we aim to extend our approach to work as well with mixture models, which in our case is related to the HiREPS approach [6]. Furthermore, it seems promising to integrate entropy constraints similar to [2], which can alleviate premature convergence. Additionally, to improve the approximations made in this paper we will investigate learning a reward model during optimization. We will further examine the possibility of sampling the covariance matrix Σ_w with a Wishart distribution. Related to which it might be beneficial to apply an Inverse-Wishart distribution as a conjugate prior to Σ_w .

Bibliography

- [1] A. Abdolmaleki, N. Lau, L. P. Reis, and G. Neumann. Regularized covariance estimation for weighted maximum likelihood policy search methods. In *Humanoid Robots (Humanoids)*, 2015.
- [2] A. Abdolmaleki, R. Lioutikov, J. Peters, N. Lau, L. P. Reis, and G. Neumann. Model-based relative entropy stochastic search. In *Advances in Neural Information Processing Systems (NIPS)*. Curran Associates, Inc., 2015.
- [3] D. Barber. *Bayesian Reasoning and Machine Learning*. Cambridge University Press, 2012.
- [4] C. M. Bishop. *Pattern Recognition and Machine Learning*. Springer-Verlag New York, Secaucus, NJ, USA, 2006.
- [5] S. Calinon, F. Guenter, and A. Billard. On learning, representing, and generalizing a task in a humanoid robot. *IEEE Transactions on Systems, Man, and Cybernetics, Part B (Cybernetics)*, 37:286–298, 2007.
- [6] C. Daniel, G. Neumann, O. Kroemer, and J. Peters. Hierarchical relative entropy policy search. *Journal of Machine Learning (JMLR)*, accepted.
- [7] M. P. Deisenroth, G. Neumann, and J. Peters. A survey on policy search for robotics. *Foundations and Trends in Robotics*, 2:1–142, 2013.
- [8] A. D. Dragan, K. C. Lee, and S. S. Srinivasa. Legibility and predictability of robot motion. In *Proceedings of the 8th ACM/IEEE International Conference on Human-robot Interaction (HRI)*, 2013.
- [9] A. D. Dragan, K. Muelling, J. A. Bagnell, and S. S. Srinivasa. Movement primitives via optimization. In *International Conference on Robotics and Automation (ICRA)*, 2015.
- [10] P. Englert, A. Paraschos, M. P. Deisenroth, and J. Peters. Probabilistic model-based imitation learning. *Adaptive Behavior*, 21:388–403, 2013.
- [11] A. Frisoli, C. Loconsole, R. Bartalucci, and M. Bergamasco. A new bounded jerk on-line trajectory planning for mimicking human movements in robot-aided neurorehabilitation. *Robotics and Autonomous Systems*, 61:404–415, 2013.
- [12] Y. Huang, M. Mahmudi, and M. Kallmann. Planning humanlike actions in blending spaces. In *International Conference on Intelligent Robots and Systems (IROS)*, 2011.
- [13] A. J. Ijspeert, J. Nakanishi, H. Hoffmann, P. Pastor, and S. Schaal. Dynamical movement primitives: Learning attractor models for motor behaviors. *Neural Computation*, 25:328–373, 2013.
- [14] M. Kalakrishnan, S. Chitta, E. Theodorou, P. Pastor, and S. Schaal. Stomp: Stochastic trajectory optimization for motion planning. In *International Conference on Robotics and Automation (ICRA)*, 2011.
- [15] J. Kober, E. Oztop, and J. Peters. Reinforcement learning to adjust robot movements to new situations. In *Proceedings of Robotics: Science and Systems (R:SS)*, 2010.
- [16] J. Kober and J. Peters. Policy search for motor primitives in robotics. In *Advances in Neural Information Processing Systems (NIPS)*, 2009.
- [17] D. Koller and N. Friedman. *Probabilistic Graphical Models: Principles and Techniques*. MIT Press, 2009.
- [18] S. Kullback. *Information Theory and Statistics*. Courier Corporation, 1968.
- [19] S. Kullback and R. A. Leibler. On information and sufficiency. *The Annals of Mathematical Statistics*, 22:79–86, 1951.

-
- [20] A. G. Kupcsik, M. P. Deisenroth, J. Peters, and G. Neumann. Data-efficient generalization of robot skills with contextual policy search. In *Proceedings of the Twenty-Seventh AAAI Conference on Artificial Intelligence*, 2013.
- [21] S. M. LaValle. Rapidly-exploring random trees: A new tool for path planning. Technical report, Computer Science Dept., Iowa State University, 1998.
- [22] B. Lim, S. Ra, and F. C. Park. Movement primitives, principal component analysis, and the efficient generation of natural motions. In *International Conference on Robotics and Automation (ICRA)*, 2005.
- [23] A. Paraschos, C. Daniel, J. Peters, and G. Neumann. Probabilistic movement primitives. In *Advances in Neural Information Processing Systems (NIPS)*, 2013.
- [24] P. Pastor, H. Hoffmann, T. Asfour, and S. Schaal. Learning and generalization of motor skills by learning from demonstration. In *International Conference on Robotics and Automation (ICRA)*, 2009.
- [25] J. Peters, K. Muelling, and Y. Altun. Relative entropy policy search. In *Proceedings of the Twenty-Fourth National Conference on Artificial Intelligence (AAAI), Physically Grounded AI Track*, 2010.
- [26] N. Ratliff, M. Zucker, J. A. D. Bagnell, and S. Srinivasa. Chomp: Gradient optimization techniques for efficient motion planning. In *International Conference on Robotics and Automation (ICRA)*, 2009.
- [27] S. Schaal. Is imitation learning the route to humanoid robots? *Trends in Cognitive Sciences*, 3:233–242, 1999.
- [28] G. Ye and R. Alterovitz. Demonstration-guided motion planning. In *Proceedings of International Symposium on Robotics Research (ISRR)*, 2011.

A Appendix

A.1 Gaussian Distribution

We make extensive use of the gaussian distribution. For some of our derivations it can be advisable to make use of the form (e.g. shown in [17, 7.1.1])

$$\begin{aligned}\mathcal{N}(x|\mu, \Sigma) &= \frac{1}{(2\pi)^{n/2}|\Sigma|^{1/2}} \exp\left(-\frac{1}{2}\left[x^T \Sigma^{-1}x - 2x^T \Sigma^{-1}\mu + \mu^T \Sigma^{-1}\mu\right]\right) \\ &\propto \exp\left(-\frac{1}{2}x^T \Sigma^{-1}x + x^T \Sigma^{-1}\mu\right).\end{aligned}\tag{A.1}$$

We can use the *linearity of expectation* property $\mathbb{E}[X + Y] = \mathbb{E}[X] + \mathbb{E}[Y]$ together with $\text{Var}[X + Y] = \text{Var}[X] + \text{Var}[Y]$ [17, p.32] to show that the sum of two independent gaussians is

$$\begin{aligned}X &\sim \mathcal{N}(\mu_X, \Sigma_X), \\ Y &\sim \mathcal{N}(\mu_Y, \Sigma_Y), \\ Z = X + Y &\implies Z \sim \mathcal{N}(\mu_X + \mu_Y, \Sigma_X + \Sigma_Y).\end{aligned}\tag{A.2}$$

Another useful property is that a linear transformation $Y = MX + c$ of a gaussian $X \sim \mathcal{N}(\mu_X, \Sigma_X)$ with independent gaussian noise $\mathcal{N}(\mu_c, \Sigma_c)$ is again distributed as a gaussian [3, p.178] resulting in

$$p(Y) = \mathcal{N}(M\mu_X + \mu_c, M\Sigma_X M^T + \Sigma_c).\tag{A.3}$$

We made use of these properties for the linear combination of two primitives (Section 2.4).

A.2 Kullback-Leibler Divergence

The Kullback-Leibler Divergence [19] is a measure for the difference between two probability distributions $p(x)$ and $q(x)$. For continuous distributions it is defined as

$$\mathcal{D}(p(x)||q(x)) = \int p(x) \log \frac{p(x)}{q(x)} dx.\tag{A.4}$$

Although it is non-symmetric $\mathcal{D}(p||q) \neq \mathcal{D}(q||p)$ it has very useful properties. It equals zero if, and only if, $p(x) = q(x)$ and always satisfies $\mathcal{D}(p||q) \geq 0$ [4].

Typically $p(x)$ is referred to as the true distribution and is approximated by $q(x)$, which is called the model distribution. If $p(x) = \mathcal{N}_1(\mu_1, \Sigma_1)$ and $q(x) = \mathcal{N}_2(\mu_2, \Sigma_2)$ are gaussian distributions the KL-Divergence can be expressed as

$$\mathcal{D}(\mathcal{N}_1||\mathcal{N}_2) = \frac{1}{2} \left(\log \frac{|\Sigma_2|}{|\Sigma_1|} + \text{tr}(\Sigma_2^{-1}\Sigma_1) + (\mu_2 - \mu_1)^T (\Sigma_2^{-1})(\mu_2 - \mu_1) - n \right),\tag{A.5}$$

with ($n = \text{dimension}(x)$) [18, 9.1]. In cases, where its infeasible to compute the determinant of Σ_1 and/or Σ_2 it is useful to compute the symmetric KL-Divergence $\mathcal{D}_{sym}(p, q) = \mathcal{D}(p||q) + \mathcal{D}(q||p)$ for gaussians as

$$\mathcal{D}_{sym}(\mathcal{N}_1||\mathcal{N}_2) = \frac{1}{4} \left(\text{tr}(\Sigma_2^{-1}\Sigma_1) + \text{tr}(\Sigma_1^{-1}\Sigma_2) + (\mu_2 - \mu_1)^T (\Sigma_1^{-1} + \Sigma_2^{-1})(\mu_2 - \mu_1) - 2n \right).\tag{A.6}$$

A.3 Derivation of Blending

For reasons of completeness we show how the blending of primitives used in [23] can be derived. The distribution over trajectories is represented as

$$p(\tau) \propto \prod_t \prod_i p_i(y_t)^{\alpha^i}, \quad (\text{A.7})$$

with (we omit subscript t in the following and use $x = y_t$)

$$\begin{aligned} p_i(x) &= \mathcal{N}(x|\mu_i, \Sigma_i) = \frac{1}{Z} \exp\left(\frac{1}{2}x^T \Sigma_i^{-1} x + x^T \Sigma_i^{-1} \mu_i + \frac{1}{2}\mu_i^T \Sigma_i^{-1} \mu_i\right) \\ &= \frac{1}{Z} \exp\left(\frac{1}{2}x^T \Sigma_i^{-1} x + x^T \Sigma_i^{-1} \mu_i + \text{const}\right) \\ &= Z' \exp\left(\frac{1}{2}x^T \Sigma_i^{-1} x + x^T \Sigma_i^{-1} \mu_i\right). \end{aligned} \quad (\text{A.8})$$

Inserting Equation A.8 in Equation A.7 gives us

$$\begin{aligned} p(\tau) &\propto \prod_t \prod_i [Z' \exp\left(\frac{1}{2}x^T \Sigma_i^{-1} x + x^T \Sigma_i^{-1} \mu_i\right)]^{\alpha^i} \\ &= \prod_t \prod_i [Z']^{\alpha^i} [\exp\left(\frac{1}{2}x^T \Sigma_i^{-1} x + x^T \Sigma_i^{-1} \mu_i\right)]^{\alpha^i} \\ &= \prod_t \prod_i \tilde{Z} \exp\left(\frac{1}{2}x^T \Sigma_i^{-1} x + x^T \Sigma_i^{-1} \mu_i\right)^{\alpha^i} \\ &= \prod_t \prod_i \tilde{Z} \exp\left[\left(\frac{1}{2}x^T \Sigma_i^{-1} x + x^T \Sigma_i^{-1} \mu_i\right) \alpha^i\right] \\ &= \prod_t \left(\prod_i \tilde{Z}\right) \left(\prod_i \exp\left[\left(\frac{1}{2}x^T \Sigma_i^{-1} x + x^T \Sigma_i^{-1} \mu_i\right) \alpha^i\right]\right) \\ &= \prod_t \left(\prod_i \tilde{Z}\right) \exp\left(\sum_i \left[\frac{1}{2}x^T \Sigma_i^{-1} x + x^T \Sigma_i^{-1} \mu_i\right] \alpha^i\right) \\ &= \prod_t \hat{Z} \exp\left(\sum_i \frac{1}{2}x^T \Sigma_i^{-1} x \alpha^i + \sum_i x^T \Sigma_i^{-1} \mu_i \alpha^i\right) \\ &= \prod_t \hat{Z} \exp\left(\frac{1}{2}x^T \sum_i \left(\frac{\Sigma_i}{\alpha^i}\right)^{-1} x + x^T \sum_i \left(\frac{\Sigma_i}{\alpha^i}\right)^{-1} \mu_i\right), \end{aligned} \quad (\text{A.9})$$

which equals the form of the resulting distribution $\mathcal{N}(y_t^*|\mu_t^*, \Sigma_t^*) = \hat{Z} \exp\left(\frac{1}{2}x^T (\Sigma_t^*)^{-1} x + x^T (\Sigma_t^*)^{-1} \mu_t^*\right)$. The covariance is now given as

$$\Sigma_t^* = \left(\sum_i (\Sigma_t^i / \alpha^i)^{-1}\right)^{-1}, \quad (\text{A.10})$$

and for the mean μ_t^* it must hold

$$\begin{aligned} (\Sigma_t^*)^{-1} \mu_t^* &= \sum_i (\Sigma_t^i / \alpha^i)^{-1} \mu_t^i \\ \mu_t^* &= \Sigma_t^* \left(\sum_i (\Sigma_t^i / \alpha^i)^{-1} \mu_t^i\right). \end{aligned} \quad (\text{A.11})$$

A.4 Derivation of our Algorithm without Sub-Policies

In relation to our full optimization algorithm we neglect the separate joint policies here and only optimize a single one. In the following we show how to derive a REPS version, which stays close to the demonstrations $\pi_d(w)$. In general, this target distribution could be arbitrary. How strongly this distribution is incorporated into the final solution can be modified by tuning γ . The optimization problem is formulated as

$$\begin{aligned} \max_{\pi} \quad & J_{\lambda} = \int_w R(w)\pi(w) - \gamma\pi(w) \log\left(\frac{\pi(w)}{\pi_d(w)}\right) dw \\ \text{s.t.} \quad & \int_w \pi(w) \log\left(\frac{\pi(w)}{q(w)}\right) \leq \epsilon, \\ & \int_w \pi(w) dw = 1. \end{aligned} \tag{A.12}$$

Formulating this with the method of lagrangian multipliers η and α and rearranging terms gives us

$$\begin{aligned} L &= \int_w R(w)\pi(w) - \gamma\pi(w) \log\left(\frac{\pi(w)}{\pi_d(w)}\right) dw - \eta \left[\int_w \pi(w) \log\left(\frac{\pi(w)}{q(w)}\right) - \epsilon \right] - \alpha \left[\int_w \pi(w) dw - 1 \right] \\ &= \int_w \left[R(w)\pi(w) - \gamma\pi(w) \log\left(\frac{\pi(w)}{\pi_d(w)}\right) - \eta\pi(w) \log\left(\frac{\pi(w)}{q(w)}\right) - \alpha\pi(w) \right] dw + \eta\epsilon + \alpha \\ &= \int_w \pi(w) \left[R(w) - \gamma \log\left(\frac{\pi(w)}{\pi_d(w)}\right) - \eta \log\left(\frac{\pi(w)}{q(w)}\right) - \alpha \right] dw + \eta\epsilon + \alpha. \end{aligned} \tag{A.13}$$

For the optimal policy we set the derivative of L to zero

$$\frac{dL}{d\pi(w)} = R(w) - \gamma \log\left(\frac{\pi(w)}{\pi_d(w)}\right) - \eta \log\left(\frac{\pi(w)}{q(w)}\right) - \eta - \alpha - \gamma \stackrel{!}{=} 0, \tag{A.14}$$

and solve for $\pi^*(w)$

$$\begin{aligned} \pi^*(w) &= \exp\left[\frac{R(w) + \gamma \log \pi_d(w) + \eta \log q(w) - \eta - \alpha - \gamma}{\gamma + \eta}\right] \\ &= [\exp[R(w) + \log(\pi_d(w)^\gamma) + \eta \log(q(w)^\eta) - \eta - \alpha - \gamma]]^{\frac{1}{\gamma + \eta}} \\ &= [\pi_d(w)^\gamma q(w)^\eta \exp(R(w) - \eta - \alpha - \gamma)]^{\frac{1}{\gamma + \eta}} \\ &= [\pi_d(w)^\gamma q(w)^\eta \exp(R(w))]^{\frac{1}{\gamma + \eta}} \exp(-\eta - \alpha - \gamma)^{\frac{1}{\gamma + \eta}}. \end{aligned} \tag{A.15}$$

To get rid of α we use the fact that π is a probability distribution and sums up to one

$$\begin{aligned} 1 &= \int_w \pi^*(w) dw = \int_w [\pi_d(w)^\gamma q(w)^\eta \exp(R(w))]^{\frac{1}{\gamma + \eta}} \exp(-\eta - \alpha - \gamma)^{\frac{1}{\gamma + \eta}} dw \\ &= \int_w [\pi_d(w)^\gamma q(w)^\eta \exp(R(w))]^{\frac{1}{\gamma + \eta}} dw \exp(-\eta - \alpha - \gamma)^{\frac{1}{\gamma + \eta}} \\ \exp(-\eta - \alpha - \gamma)^{\frac{1}{\gamma + \eta}} &= \frac{1}{\int_w [\pi_d(w)^\gamma q(w)^\eta \exp(R(w))]^{\frac{1}{\gamma + \eta}} dw}. \end{aligned} \tag{A.16}$$

The exponential term can be plugged into Equation A.15 and gives us the policy update rule

$$\pi^*(w) = \frac{[\pi_d(w)^\gamma q(w)^\eta \exp(R(w))]^{\frac{1}{\gamma + \eta}}}{\int_w [\pi_d(w)^\gamma q(w)^\eta \exp(R(w))]^{\frac{1}{\gamma + \eta}} dw}. \tag{A.17}$$

To compute η we plug the optimal policy (Equation A.15) into the lagrangian (Equation A.13)¹ and derive the dual function $g(\eta)$

$$\begin{aligned}
g(\eta) &= \int_w \pi(w) \left[R(w) - \gamma \log\left(\frac{\pi^*(w)}{\pi_d(w)}\right) - \eta \log\left(\frac{\pi^*(w)}{q(w)}\right) - \alpha \right] dw + \eta\epsilon + \alpha \\
&= \int_w \pi(w) \left[R(w) - \gamma \log\left(\frac{[\pi_d(w)^\gamma q(w)^\eta \exp(R(w) - \eta - \alpha - \gamma)]^{\frac{1}{\gamma+\eta}}}{\pi_d(w)}\right) \right. \\
&\quad \left. - \eta \log\left(\frac{[\pi_d(w)^\gamma q(w)^\eta \exp(R(w) - \eta - \alpha - \gamma)]^{\frac{1}{\gamma+\eta}}}{q(w)}\right) - \alpha \right] dw + \eta\epsilon + \alpha \\
&= \int_w \pi(w) \left[R(w) + (-\gamma - \eta) \log([\pi_d(w)^\gamma q(w)^\eta \exp(R(w) - \eta - \alpha - \gamma)]^{\frac{1}{\gamma+\eta}}) + \log(\pi_d(w)^\gamma) + \eta \log(q(w)) - \alpha \right] dw + \eta\epsilon + \alpha \\
&= \int_w \pi(w) \left[R(w) - \log(\pi_d(w)^\gamma q(w)^\eta \exp(R(w) - \eta - \alpha - \gamma)) + \log(\pi_d(w)^\gamma) + \eta \log(q(w)) - \alpha \right] dw + \eta\epsilon + \alpha \\
&= \int_w \pi(w) (\eta + \gamma) dw + \eta\epsilon + \alpha \\
&= (\eta + \gamma) \int_w \pi(w) dw + \eta\epsilon + \alpha.
\end{aligned} \tag{A.18}$$

With the fact, that $\int_w \pi(w) dw = 1$ we get

$$g(\eta) = \eta\epsilon + \eta + \alpha + \gamma. \tag{A.19}$$

We rewrite equation Equation A.16 to get rid of α as

$$\eta + \alpha + \gamma = (\gamma + \eta) \log\left(\int_w [\pi_d(w)^\gamma q(w)^\eta \exp(R(w))]^{\frac{1}{\gamma+\eta}} dw\right). \tag{A.20}$$

Plugging the result into Equation A.19 we get a representation for the dual, which does not depend on α anymore

$$g(\eta) = \eta\epsilon + (\gamma + \eta) \log\left(\int_w [\pi_d(w)^\gamma q(w)^\eta \exp(R(w))]^{\frac{1}{\gamma+\eta}} dw\right). \tag{A.21}$$

We can not analytically solve the integral in close form. Hence, we approximate it with importance sampling. We reuse the samples from the old distribution $q(w)$ and therefore first rewrite the integral as

$$\begin{aligned}
\int_w [\pi_d(w)^\gamma q(w)^\eta \exp(R(w))]^{\frac{1}{\gamma+\eta}} dw &= \int_w q(w) \frac{[\pi_d(w)^\gamma q(w)^\eta \exp(R(w))]^{\frac{1}{\gamma+\eta}}}{q(w)} dw \\
&= \int_w q(w) [\pi_d(w)^\gamma q(w)^{-\gamma} \exp(R(w))]^{\frac{1}{\gamma+\eta}} dw.
\end{aligned} \tag{A.22}$$

Using importance sampling with N samples we get

$$\tilde{g}(\eta) = \eta\epsilon + (\gamma + \eta) \log\left(\frac{1}{N} \sum_{i=1}^N [\pi_d(w_i)^\gamma q(w_i)^{-\gamma} \exp(R(w_i))]^{\frac{1}{\gamma+\eta}}\right). \tag{A.23}$$

¹ Only plug in inside the log since the rest cancels out later.

To improve numerical stability we first rearrange terms and then apply a log-sum-exp identity² as

$$\begin{aligned}
\tilde{g}(\eta) &= \eta\epsilon + (\gamma + \eta) \log \left(\sum_{i=1}^N \exp \left(\frac{1}{\gamma + \eta} (R(w_i) + \gamma [\log(\pi_d(w_i)) - \log(q(w_i))]) - \log(N) \right) \right) \\
&= \eta\epsilon + (\gamma + \eta) \left[\log \left(\sum_{i=1}^N \exp(s_i - k) \right) + k \right]
\end{aligned} \tag{A.24}$$

with

$$s_i = \frac{R(w_i) + \gamma [\log(\pi_d(w_i)) - \log(q(w_i))]}{\gamma + \eta} - \log(N)$$

$$k = \max(s_i).$$

The remaining optimization problem, which can be solved using any constrained optimizer, is given as

$$\operatorname{argmin}_{\eta} \tilde{g}(\eta) \quad \text{s.t.} \quad \eta > 0. \tag{A.25}$$

The derivative of $\tilde{g}(\eta)$, which is useful for optimization is

$$\begin{aligned}
\frac{d}{d\eta} \tilde{g}(\eta) &= \epsilon + \log \left(\frac{1}{N} \sum_{i=1}^N z_i \right) - \frac{1}{(\gamma + \eta)} \frac{\sum_{i=1}^N [R(w_i) + \gamma (\log(\pi_d(w_i)) - \log(q(w_i)))] z_i}{\sum_{i=1}^N z_i} \\
&\text{with} \\
z_i &= \exp \left(\frac{R(w_i) + \gamma [\log(\pi_d(w_i)) - \log(q(w_i))]}{\gamma + \eta} \right).
\end{aligned} \tag{A.26}$$

To make the gradient numerically stable, we can apply the exp-normalize trick³.

² The identity

$$\log \left(\sum_i \exp(x_i) \right) = k + \log \sum_i \exp(x_i - k)$$

holds. In our case $k = \max x_i$ is reasonable.

³

$$\frac{\sum_{i=1}^N a_i \exp(z_i)}{\sum_{i=1}^N \exp(z_i)} = \frac{\exp(k) \sum_{i=1}^N a_i \exp(z_i)}{\exp(k) \sum_{i=1}^N \exp(z_i)} = \frac{\sum_{i=1}^N a_i \exp(z_i + k)}{\sum_{i=1}^N \exp(z_i + k)}$$

with $k = -\max(z_i)$.

A.5 Detailed Derivation of our full Algorithm

For our full optimization problem we now assume that our policy is a joint distribution $\pi(w)\lambda(a)$. We maximize the expected reward, while the KL-divergence of $\pi(w)$ to the demonstrated policy $\pi_d(w)$ is minimized. The reward $R(w, a)$ depends on both sub-policies. The full optimization problem is given as

$$\begin{aligned}
 J &= \int_w \int_a R(w, a) \pi(w) \lambda(a) da dw - \gamma \int_w \pi(w) \log \left(\frac{\pi(w)}{\pi_d(w)} \right) dw \\
 \text{s.t. } &\int_w \pi(w) \log \left(\frac{\pi(w)}{q_1(w)} \right) dw \leq \epsilon_1 \\
 &\int_a \lambda(a) \log \left(\frac{\lambda(a)}{q_2(a)} \right) da \leq \epsilon_2 \\
 &\int_w \pi(w) dw = 1 \\
 &\int_a \lambda(a) da = 1.
 \end{aligned} \tag{A.27}$$

Using the method of lagrangian multipliers we can rewrite the problem as

$$L = J - \eta_1 \left[\int_w \pi(w) \log \left(\frac{\pi(w)}{q_1(w)} \right) dw - \epsilon_1 \right] - \eta_2 \left[\int_a \lambda(a) \log \left(\frac{\lambda(a)}{q_2(a)} \right) da - \epsilon_2 \right] - \alpha_1 \left[\int_w \pi(w) dw - 1 \right] - \alpha_2 \left[\int_a \lambda(a) da - 1 \right]. \tag{A.28}$$

We solve for the optimal policy by deriving $\frac{dL}{d\pi(w)}$ and $\frac{dL}{d\lambda(a)}$ and setting both to zero

$$\begin{aligned}
 \frac{dL}{d\pi(w)} &= \int_a R(w, a) \lambda(a) da - \gamma \log \left(\frac{\pi(w)}{\pi_d(w)} \right) - \gamma - \eta_1 \log \left(\frac{\pi(w)}{q_1(w)} \right) - \eta_1 - \alpha_1 \stackrel{!}{=} 0, \\
 \frac{dL}{d\lambda(a)} &= \int_w R(w, a) \pi(w) dw - \eta_2 \log \left(\frac{\lambda(a)}{q_2(a)} \right) - \eta_2 - \alpha_2 \stackrel{!}{=} 0.
 \end{aligned} \tag{A.29}$$

Solving for the optimal policies $\pi^*(w)$ and $\lambda^*(a)$ results in

$$\begin{aligned}
 \pi^*(w) &= [\exp(R(w)) \pi_d(w)^\gamma q_1(w)^{\eta_1}]^{\frac{1}{\gamma + \eta_1}} \exp(-\gamma - \eta_1 - \alpha_1)^{\frac{1}{\gamma + \eta_1}}, \\
 \lambda^*(a) &= q_2(a) \exp \left(\frac{R(a)}{\eta_2} \right) \exp(-\eta_2 - \alpha_2)^{\frac{1}{\eta_2}},
 \end{aligned} \tag{A.30}$$

where we have replaced the integral terms in the exponential functions with functions $R(a)$ and $R(w)$ respectively. Both represent the expected reward for a fixed parameter a or w and are given as

$$\begin{aligned}
 R(w) &= \int_a R(w, a) \lambda(a) da, \\
 R(a) &= \int_w R(w, a) \pi(w) dw.
 \end{aligned} \tag{A.31}$$

With the optimal policies we can solve for the policy update rule, which no longer depends on the lagrangian parameters (except temperatures η_1 and η_2).

We use the fact that $\pi(w)$ and $\lambda(a)$ are probability distributions and insert

$$\begin{aligned} 1 &= \int_w \pi^*(w) dw = \int_w [\exp(R(w)) \pi_d(w)^\gamma q_1(w)^{\eta_1}]^{\frac{1}{\gamma+\eta_1}} dw \exp(-\gamma - \eta_1 - \alpha_1)^{\frac{1}{\gamma+\eta_1}}, \\ 1 &= \int_a \lambda^*(a) da = \int_a q_2(a) \exp\left(\frac{R(a)}{\eta_2}\right) da \exp(-\eta_2 - \alpha_2)^{\frac{1}{\eta_2}}. \end{aligned} \quad (\text{A.32})$$

Rewriting the terms yields

$$\begin{aligned} \exp(-\gamma - \eta_1 - \alpha_1)^{\frac{1}{\gamma+\eta_1}} &= \frac{1}{\int_w [\exp(R(w)) \pi_d(w)^\gamma q_1(w)^{\eta_1}]^{\frac{1}{\gamma+\eta_1}} dw}, \\ \exp(-\eta_2 - \alpha_2)^{\frac{1}{\eta_2}} &= \frac{1}{\int_a q_2(a) \exp\left(\frac{R(a)}{\eta_2}\right) da}. \end{aligned} \quad (\text{A.33})$$

Plugging Equation A.33 into the optimal policy Equation A.30 gives us the policy update rules

$$\begin{aligned} \pi^*(w) &= \frac{[\exp(R(w)) \pi_d(w)^\gamma q_1(w)^{\eta_1}]^{\frac{1}{\gamma+\eta_1}}}{\int_w [\exp(R(w)) \pi_d(w)^\gamma q_1(w)^{\eta_1}]^{\frac{1}{\gamma+\eta_1}} dw} \propto q_1(w) [\exp(R(w)) \pi_d(w)^\gamma q_1(w)^{-\gamma}]^{\frac{1}{\gamma+\eta_1}}, \\ \lambda^*(a) &= \frac{q_2(a) \exp\left(\frac{R(a)}{\eta_2}\right)}{\int_a q_2(a) \exp\left(\frac{R(a)}{\eta_2}\right) da} \propto q_2(a) \exp\left(\frac{R(a)}{\eta_2}\right). \end{aligned} \quad (\text{A.34})$$

The optimal policies are weighted updates of their old policy, allowing us to develop an iterative optimization algorithm.

A.5.1 Derivation of the Dual

To solve for the unknown lagrangian parameters η_1 and η_2 we plug in the optimal policy Equation A.30 into the lagrangian approach Equation A.28 yielding

$$\begin{aligned} g(\eta_1, \eta_2) &= \int_w \int_a R(w, a) \pi^*(w) \lambda^*(a) da dw \\ &\quad - \eta_1 \left[\int_w \pi^*(w) \log\left(\frac{\pi^*(w)}{q_1(w)}\right) dw - \epsilon_1 \right] \\ &\quad - \eta_2 \left[\int_a \lambda^*(a) \log\left(\frac{\lambda^*(a)}{q_2(a)}\right) da - \epsilon_2 \right] \\ &\quad - \alpha_1 \left[\int_w \pi(w)^* dw - 1 \right] \\ &\quad - \alpha_2 \left[\int_a \lambda(a)^* da - 1 \right]. \end{aligned} \quad (\text{A.35})$$

Simplifying the expression results in

$$g = \int_w \int_a R(w, a) \pi^*(w) \lambda^*(a) da dw \quad (\text{A.36a})$$

$$- \int_w \pi^*(w) \left[\gamma \log\left(\frac{\pi^*(w)}{\pi_d(w)}\right) + \eta_1 \log\left(\frac{\pi^*(w)}{q_1(w)}\right) + \alpha_1 \right] dw \quad (\text{A.36b})$$

$$- \int_a \lambda^*(a) \left[\eta_2 \log\left(\frac{\lambda^*(a)}{q_2(a)}\right) + \alpha_2 \right] da \quad (\text{A.36c})$$

$$+ \eta_1 \epsilon_1 + \eta_2 \epsilon_2 + \alpha_1 + \alpha_2. \quad (\text{A.36d})$$

Further simplification give us

$$\begin{aligned}
g &= \int_w \int_a R(w, a) \pi^*(w) \lambda^*(a) da dw - \underbrace{\int_w \pi^*(w) R(w) dw}_{\text{Equation A.36b}} + \gamma + \eta_1 - \underbrace{\int_a \lambda^*(a) R(a) da}_{\text{Equation A.36c}} + \eta_2 + \eta_1 \epsilon_1 + \eta_2 \epsilon_2 + \alpha_1 + \alpha_2 \\
&= - \int_w \int_a R(w, a) \pi^*(w) \lambda^*(a) da dw + \eta_1 \epsilon_1 + \eta_2 \epsilon_2 + \alpha_1 + \alpha_2 + \eta_1 + \eta_2 + \gamma \\
&= -\mathbb{E}[R(w, a)]_{\pi^* \lambda^*} + \eta_1 \epsilon_1 + \eta_2 \epsilon_2 + \alpha_1 + \alpha_2 + \eta_1 + \eta_2 + \gamma.
\end{aligned} \tag{A.37}$$

We now solve Equation A.33 for $\gamma + \eta_1 + \alpha_1$ and $\eta_2 + \alpha_2$ and plug the results into the dual Equation A.37 resulting in

$$\begin{aligned}
g(\eta_1, \eta_2) &= -\mathbb{E}[R(w, a)]_{\pi^* \lambda^*} + \eta_1 \epsilon_1 + \eta_2 \epsilon_2 \\
&\quad + (\gamma + \eta_1) \log \left(\int_w q_1(w) [\exp(R(w)) \pi_d(w)^\gamma q_1(w)^{-\gamma}]^{\frac{1}{\gamma + \eta_1}} dw \right) \\
&\quad + \eta_2 \log \left(\int_a q_2(a) \exp\left(\frac{R(a)}{\eta_2}\right) da \right).
\end{aligned} \tag{A.38}$$

A.5.2 Expectation Approximation

To be able to solve the dual function Equation A.38 we approximate the expectation $\mathbb{E}[R(w, a)]_{\pi^* \lambda^*}$ with samples using importance sampling, which gives us

$$\mathbb{E}[R(w, a)]_{\pi^* \lambda^*} = \int_w \int_a \left[R(w, a) \frac{\pi^*(w) \lambda^*(a)}{q_1(w) q_2(a)} \right] q_1(w) q_2(a) da dw \approx \frac{\frac{1}{N} \sum_i R(w_i, s_i) \psi(w_i, s_i)}{\frac{1}{N} \sum_i \psi(w_i, s_i)}. \tag{A.39}$$

The importance weights $\psi(w_i, s_i)$ are calculated using the unnormalized distributions from the policy update rules Equation A.34 as

$$\begin{aligned}
\psi(w_i, s_i) &= \frac{\tilde{\pi}(w) \tilde{\lambda}(a)}{q_1(w) q_2(a)} \\
&= \frac{q_1(w) [\exp(R(w)) \pi_d(w)^\gamma q_1(w)^{-\gamma}]^{\frac{1}{\gamma + \eta_1}} q_2(a) \exp\left(\frac{R(a)}{\eta_2}\right)}{q_1(w) q_2(a)} \\
&= [\exp(R(w)) \pi_d(w)^\gamma q_1(w)^{-\gamma}]^{\frac{1}{\gamma + \eta_1}} \exp\left(\frac{R(a)}{\eta_2}\right).
\end{aligned} \tag{A.40}$$

For completeness, the log transformation of the importance weights $\log(\psi(w_j, s_j))$ is given as

$$\log(\psi(w_j, s_j)) = \frac{R(w) + \gamma (\log \pi_d(w) - \log q_1(w))}{\gamma + \eta_1} + \frac{R(a)}{\eta_2}. \tag{A.41}$$