

Latent Space Reinforcement Learning

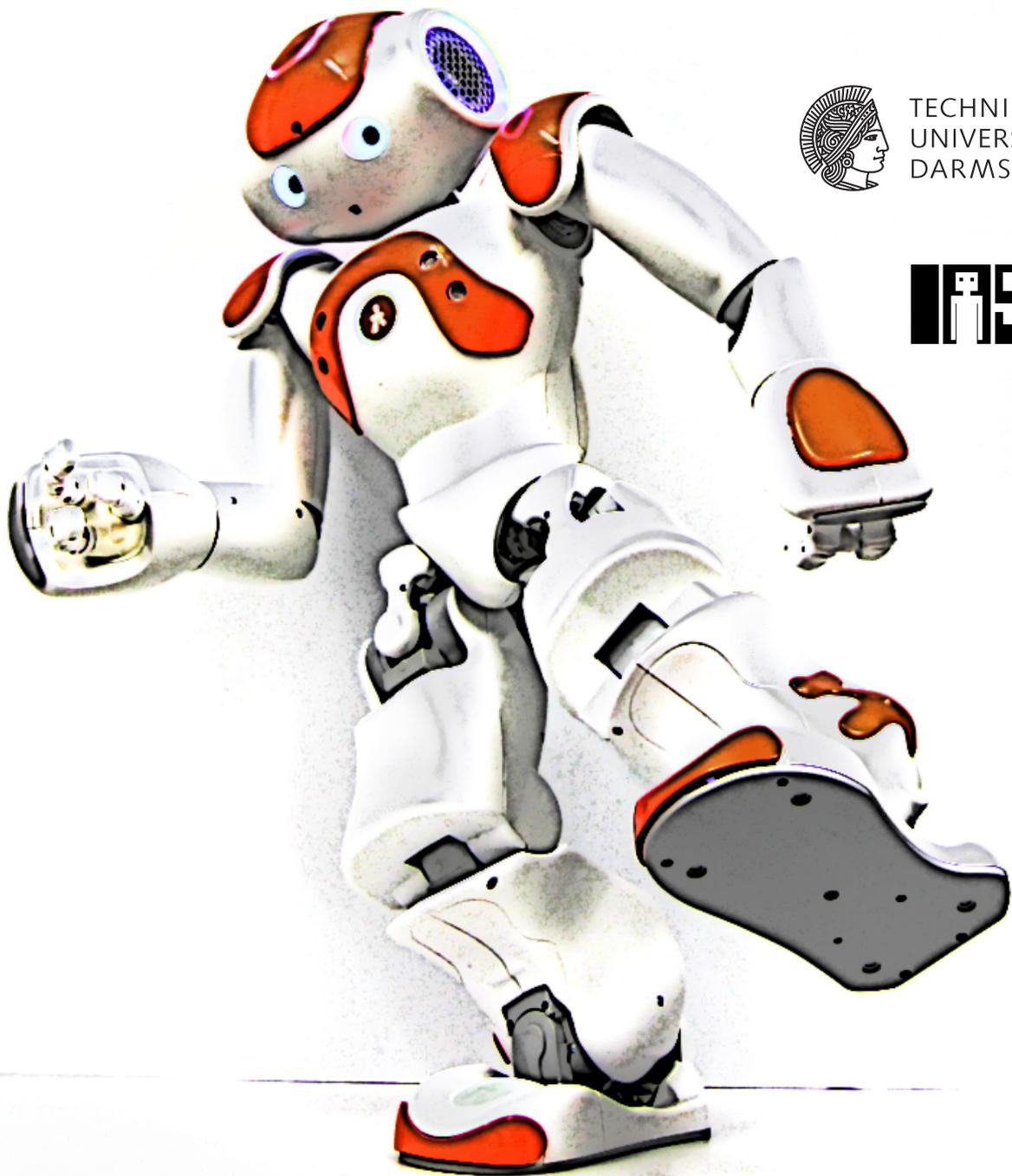
Selbstverstärkendes Lernen in latenten Vektorräumen

Bachelor-Thesis von Kevin Sebastian Luck aus Hannoversch Münden

Mai 2014



TECHNISCHE
UNIVERSITÄT
DARMSTADT



Latent Space Reinforcement Learning
Selbstverstärkendes Lernen in latenten Vektorräumen

Vorgelegte Bachelor-Thesis von Kevin Sebastian Luck aus Hannoversch Münden

1. Gutachten: Prof. Dr. Jan Peters
2. Gutachten: Dr. Heni Ben Amor
3. Gutachten: Dr. Gerhard Neumann

Tag der Einreichung:

Please cite this document with:

URN: [urn:nbn:de:tuda-tuprints-38321](https://nbn-resolving.org/urn:nbn:de:tuda-tuprints-38321)

URL: <http://tuprints.ulb.tu-darmstadt.de/id/eprint/3832>

Dieses Dokument wird bereitgestellt von tuprints,

E-Publishing-Service der TU Darmstadt

<http://tuprints.ulb.tu-darmstadt.de>

tuprints@ulb.tu-darmstadt.de



This publication is licensed under the following Creative Commons License:

Attribution 4.0 International

<http://creativecommons.org/licenses/by/4.0/>

In memory of Cornelia Lutz

Erklärung zur Bachelor-Thesis

Hiermit versichere ich, die vorliegende Bachelor-Thesis ohne Hilfe Dritter und 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.

In der abgegebenen Thesis stimmen die schriftliche und elektronische Fassung überein.

Darmstadt, den 6. Mai 2014

(Kevin S. Luck)

Thesis Statement

I herewith formally declare that I have written the submitted thesis independently. I did not use any outside support except for the quoted literature and other sources mentioned in the paper. I clearly marked and separately listed all of the literature and all of the other sources which I employed when producing this academic work, either literally or in content. This thesis has not been handed in or published before in the same or similar form.

In the submitted thesis the written copies and the electronic version are identical in content.

Darmstadt, May 6, 2014

(Kevin S. Luck)

Abstract

Often we have to handle high dimensional spaces if we want to learn motor skills for robots. In policy search tasks we have to find several parameters to learn a desired movement. This high dimensionality in parameters can be challenging for reinforcement algorithms, since more samples for finding an optimal solution are needed with every additional dimension. On the other hand, if the robot has a high number of actuators, an inherent correlation between these can be found for a specific motor task, which we can exploit for a faster convergence. One possibility is to use techniques to reduce the dimensionality of the space, which is used as a pre-processing step or as an independent process in most applications. In this thesis we present a novel algorithm which combines the theory of policy search and probabilistic dimensionality reduction to uncover the hidden structure of high dimensional action spaces. Evaluations on an inverse kinematics task indicate that the presented algorithm is able to outperform the reference algorithms PoWER and CMA-ES, especially in high dimensional spaces. Furthermore we evaluate our algorithm on a real-world task. In this task, a NAO robot learns to lift his leg while keeping balance. The issue of collecting samples for learning on a real robot in such a task, which is often very time and cost consuming, is considered in here by using a small number of samples in each iteration.

Zusammenfassung

Versucht man für Roboter Bewegungsabläufe mithilfe von Verstärkendem Lernen, dem sogenannten Reinforcement Learning, zu lernen, muss man sich häufig mit hochdimensionalen Räumen auseinandersetzen. So hat beispielsweise bereits der NAO Roboter 26 Freiheitsgrade, hier Gelenkwinkel, mithilfe derer eine Bewegung umgesetzt werden kann. Zusätzlich kommen bei Policy Search Problemen noch Parameter für sog. Features hinzu, welche beispielsweise Gauss-Kurven sein können, die abhängig von der Zeit sind. Versuchen wir nun einen Bewegungsablauf zu erlernen, erreichen wir mit steigender Zahl der einzelnen Bewegungen relativ schnell eine sehr hohe Anzahl an Parametern aus der die Bewegungen abgeleitet werden. Diese Parameter müssen für einen optimalen Bewegungsablauf geschätzt werden und spannen einen hochdimensionalen Raum auf.

In hochdimensionalen Räumen, in denen die einzelnen (Gelenk-) Konfigurationen als Punkte liegen, kann man jedoch oft Unterräume finden, die eine stark reduzierte Anzahl an Dimensionen aufweisen. Die Dimensionsachsen dieser niedrig dimensionalen Räume kodieren hier Korrelationen zwischen den verschiedenen Parametern, wobei wir unser Augenmerk auf Korrelationen zwischen den Aktoren, z.B. Gelenke, legen werden. Ein intuitives Beispiel für solche Korrelationen sind z.B. Bewegungen einer (menschlichen) Hand, bei der die Gelenke der Finger oftmals in Abhängigkeit zueinander stehen.

In dieser Abschlussarbeit wird ein neuartiger Policy Search Algorithmus vorgestellt, der die versteckte latente Struktur in einem solchen hochdimensionalen Parameterraum ausnutzt und damit Reinforcement Learning und Dimensionsreduktion in einer Theorie vereint. Im Gegensatz zu früheren Ansätzen wird hierbei die Dimensionsreduktion nicht als ein Vorverarbeitungsschritt oder als ein unabhängiger Prozess eingesetzt, sondern direkt im Lernalgorithmus durchgeführt. Wie die durchgeführten Evaluationen mit einer Aufgabe der inversen Kinematik zeigen, kann insbesondere bei einer hohen Anzahl an Gelenkwinkel der vorgestellte Algorithmus eine deutlich bessere Konvergenz in Richtung eines optimalen Ergebnisses vorweisen als die zur Referenz verwendeten Algorithmen PoWER und CMA-ES. Weiterhin wird die Möglichkeit der Anwendung in realen Lern-Szenarien anhand eines Experiments mit einem NAO Roboter aufgezeigt, bei dem der Roboter die Fähigkeit erwerben soll, auf einem Bein zu stehen. In beiden Evaluationen wurde der normalerweise begrenzten Möglichkeit zur Erzeugung von Testläufen auf Robotern Rechnung getragen, indem eine möglichst kleine Anzahl an Testläufen pro Iteration gewählt wurde.

Acknowledgements

First, I want to thank my supervisor Heni Ben Amor, who gave me the opportunity and the freedom to work on the topic of this thesis and to develop my ideas. Especially for the transatlantic supervising, which was sometimes difficult but beneficial in the end.

I would also like to thank Gerhard Neumann, who was open for mathematical discussions all the time and has given me plenty of ideas for the future.

I owe the idea for the Lift-One-Leg Task on the NAO Robot to Erik Berger, who helped me with the simulation tool and did the physically test runs on the NAO robot in Freiberg.

Jan Peters for establishing my first contact to Heni when I was new to the group.

My two colleagues Hongh Linh Thai and especially Thomas Hesse, for all the discussions late in the night, when we had to work in the office to solve our problems for the deadlines. Actually, the name of the PePPeEr Algorithm can be traced back to one of this conversations with Thomas.

Furthermore I gratefully acknowledge the many helpful suggestions of Annemarie Arnold, Judith Bönninghausen, Inga Schmidt and Wolfgang Heenes during the preparation of this thesis.

And all my friends, my family and especially my sister Susan, that they have let me work in solitude and privacy in the past months with my formulas.

Contents

1. Introduction	1
1.1. Motivation	2
1.2. Outline	2
2. Foundation	3
2.1. Matrix Variate Normal Distributions	3
2.2. Expectation Maximization	4
2.3. Dimensionality Reduction with Probabilistic Principal Component Analysis	5
2.4. Policy Search	5
2.5. The Reinforcement Learning Framework of PoWER	6
3. Reinforcement Learning in Spaces with Low Intrinsic Dimensionality	9
3.1. Using the PoWER Framework for Estimating Policies with Latent Variables	9
3.2. Low Dimensional and Stochastic Policy Models	9
3.3. Linear Model Properties	10
3.4. Expectation and Maximization of the Parameters	11
3.4.1. Expectation	11
3.4.2. Maximization	11
3.4.3. Complete Algorithm	14
3.5. Diagonal Matrix Extension	14
4. Experiments	17
4.1. Learning Inverse Kinematics	17
4.2. Learning to Stand on One Leg	18
4.3. Experimental Intuition	20
5. Conclusion and Future Work	23
Bibliography	24
A. Appendix	27
A.1. Derivatives of Scalar Forms	27
A.2. Marginalization Rule	27

Figures

List of Figures

1.1. Two different and common used approaches to use dimensionality reduction techniques in policy search.	1
1.2. The presented approach combines dimensionality reduction and policy search in one method.	2
4.2. The rows represent the Gaussians and the columns the specific time step. The values of the Gaussians are normalized for each time step.	17
4.1. A robot arm with 8 DOF in a 2D Tracking Task. These movements were learned with PePP _c Er within 1000 iterations and have an error near zero. The red circle is the desired trajectory of the end-effector.	17
4.3. Comparison between PePP _c Er , PoWER and CMA-ES on the inverse kinematic task with a 26-linked robot arm. In each iteration we generated 50 parameterized joint configurations and used the 25 best samples. Every Algorithm was performed with 30 trials and plotted with his mean.	18
4.4. Different comparisons between PePP _c Er and one algorithm from Fig. 4.3, each algorithm result plotted with its mean and variance.	19
4.5. The final poses of two extreme solutions learned by the PePP _c Er Algorithm, which are still stable.	19
4.6. Two different policies for standing on one leg learned using latent space policy search. The PePP _c Er algorithm needed only 100 executions for the first policy.	20
4.7. The results of the first five iterations of Policy 5 learned by the PePP _c Er Algorithm	20
4.8. The resulting poses after 1, 5, 10 and 15 iterations with the PePP _c Er Algorithm. Each row is one independent trial.	21
4.9. The sixth Bukin function.	21
4.10. The PePP _c Er Algorithm performed on the sixth Bukin function. The red star is the current mean in the given iteration and the black line the principal component. Both are calculated from five samples, given by the white crosses. The global minima of the Bukin function is given by the red dot.	22

Abbreviations, Symbols and Operators

List of Abbreviations

CMA-ES	Covariance Matrix Adaption Evolution Strategy
DOF	degrees-of-freedom
EM	Expectation Maximization
KL	Kullback-Leibler
ML	Machine Learning
MLE	Maximum Likelihood Estimation
PCA	Principle Component Analysis
p.d.f.	probability density function
PePPER	'Policy Search with Probabilistic Principle Component Exploration in the Action Space'
PL	Policy Learning
PoWER	'Policy learning by Weighting Exploration with the Returns'
PPCA	Probabilistic Principal Component Analysis
RL	Reinforcement Learning

List of Symbols

Notation	Description
a	a single action
\mathbf{a}	the action vector
ϵ	a normal distributed (isotropic) exploration
\mathbf{E}	a matrix variate normal distributed exploration
ϕ	the feature vector
\mathbf{I}	the identity matrix

\mathbf{Z}	the latent matrix which contains different latent vectors
\mathbf{z}	the latent vector or variable
\mathbf{W}	the basis transformation matrix from the latent subspace into original
$\boldsymbol{\mu}$	the mean vector of a probability distribution
\mathbf{M}	the mean matrix of a probability distribution
\mathbf{D}	diagonal covariance matrix of a probability distribution
σ^2	variance of a probability distribution
$\boldsymbol{\theta}$	a set of parameters
r	a single reward for one step
\mathbf{s}	the current state
$\boldsymbol{\tau}$	a trajectory
$\mathbf{0}$	a matrix or vector which only contains zeros

List of Operators

Notation	Description	Operator
\det	the determinant of a matrix	$\det(\cdot)$
diag	delivers the diagonal matrix of a matrix	$\text{diag}(\cdot)$
\mathbb{E}	the expectation operator, it will be used if the context is clear	$\mathbb{E}[\cdot]$
$\mathbb{E}_{\mathbf{x}}$	the expectation in respect of a random variable \mathbf{x}	$\mathbb{E}_{\mathbf{x}}[\cdot]$
$\mathbb{E}_{\mathbf{x} \mathbf{y}}$	the expectation in respect of a random variable \mathbf{x} conditioned on \mathbf{y}	$\mathbb{E}_{\mathbf{x} \mathbf{y}}[\cdot]$
\otimes	the kronecker product of two matrices \mathbf{A} and \mathbf{B}	$\mathbf{A} \otimes \mathbf{B}$
KL	Kullback-Leibler divergence	$\text{KL}(\cdot \parallel \cdot)$
L	the likelihood function	$L(\cdot)$
\mathcal{L}	the lower bound of a likelihood	$\mathcal{L}_{\boldsymbol{\theta}}(\cdot)$
\mathcal{N}	(multivariate) normal distribution with mean $\boldsymbol{\mu}$ and covariance Σ	$\mathcal{N}(\boldsymbol{\mu}, \Sigma)$
$\mathcal{N}_{p,n}$	matrix variate normal distribution with mean matrix \mathbf{M} and covariance matrix Σ	$\mathcal{N}_{p,n}(\mathbf{M}, \Sigma)$
\sim	the random variable \mathbf{x} is distributed with distribution \mathcal{N}	$\mathbf{x} \sim \mathcal{N}$
π	stochastic policy	$\pi(\cdot \cdot)$
p	the probability of \mathbf{x} (conditioned on \mathbf{y})	$p(\mathbf{x} \mathbf{y})$
Q_t^π	the weighting based on the reward	$Q_t^\pi(\mathbf{s}_t, \mathbf{a})$
trace	delivers the trace of a matrix	$\text{trace}(\cdot)$
vec	the vectorization of a matrix	$\text{vec}(\cdot)$
var	the variance operator	$\text{var}(\cdot)$
cov	the covariance operator	$\text{cov}(\cdot, \cdot)$

1 Introduction

Closely related to artificial intelligence is the vision of fully autonomous robots, that can perform tasks without the help or supervision of humans. A core component of this autonomous behavior is the ability of the robot to improve itself during the execution of tasks based on a reward. Developing algorithmic approaches that realize such learning abilities is the goal of the reinforcement learning (RL) community. Reinforcement learning has been applied to a broad field of applications in recent years including helicopter flight [1], robot table tennis [2], or quadruped locomotion [3].

An instance of RL-based techniques is *policy search*. Policy search methods try to find a policy π , which maximizes the expected long-term reward. In the case of a parameterized policy, i.e. linear Gaussian Models or Dynamic Motor Primitives [4], we have to determine several parameters for each degree-of-freedom. Thus, there is a strong connection between the number of parameters and degrees-of-freedom of a robot, which results in high dimensional parameter spaces in most real world tasks.

Admittedly, high dimensional parameter spaces often introduce problems, namely the *curse of dimensionality* [5]. With every new dimension, more trials need to be made to determine an optimal policy with policy search. This is due to the reinforcement learning nature of policy search, where we have to explore the parameter space to find improvements for the policy. In case of a naïve approach we need exponentially more samples. Furthermore, the simulation or the real trial on a robot is very time consuming and costly. Hence, we want to minimize the number of trials.

This gap between efficiency in terms of the number of trials and the effort for exploration in high dimensional spaces can be closed by the usage of dimensionality reduction techniques. The idea behind such techniques is the assumption that there are correlations between the axes of the high dimensional space such that the valid parameter vectors lie on a low dimensional manifold. A well-known method for finding low dimensional manifolds this is the Principle Component Analysis (PCA) [6] and its probabilistic version the Probabilistic Principal Component Analysis (PPCA) [7]. For example, by using PCA we can find such a manifold with 8 dimensions for the human hand, which has around 23 degrees-of-freedom (DOF). Solving optimization problems on this manifold is significantly easier than on the original high dimensional space.

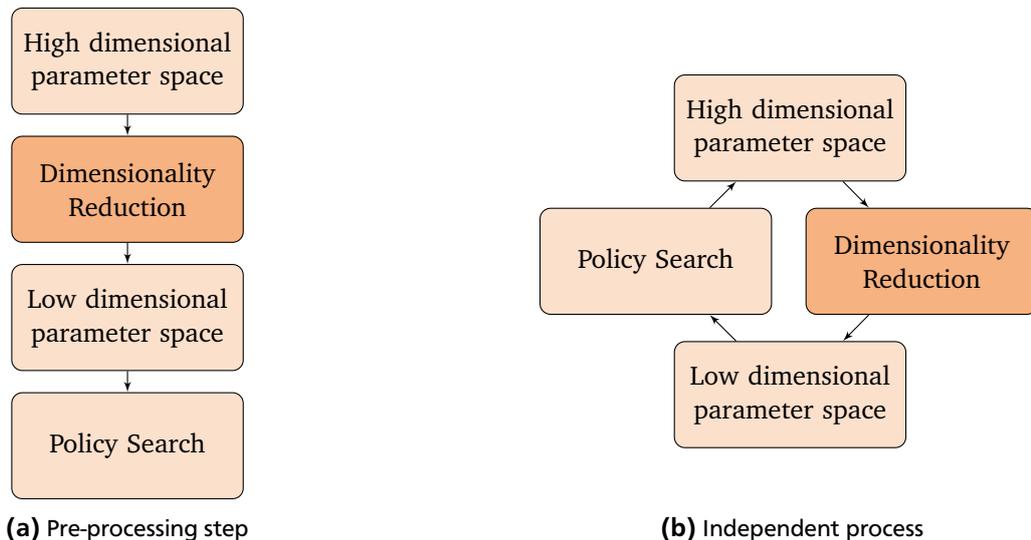


Figure 1.1.: Two different and common used approaches to use dimensionality reduction techniques in policy search.

However, in current applications these techniques are usually used as a pre-processing step (see Fig. 1.1(a)), before reinforcement learning or as an independent process (see Fig. 1.1(b)): In [8], Deisenroth et al. projected the original 25-dimensional training inputs onto an 8-dimensional latent space for predicting a walking gait for a biped robot. For the approach presented in [9], Reduced Rank Regression was used to identify a low dimensional

subspace, which was used for learning on the real robot. In a similar manner, Bitzer et al. [10] used user-provided training data to learn a latent space with linear as well as nonlinear dimensionality reduction for robot learning.

However, this use of dimensionality reduction has serious limitations. First, using dimensionality reduction as a pre-processing step requires a significantly large training set of solutions for an accurate approximation of the assumed manifold. Often, such a large training set is not available or requires human demonstrations. Furthermore, we have to deal with the problem that if the latent space is chosen wrongly for the given task, the reinforcement algorithm may not be able to find a satisfying solution since it is not possible to change the parameters of the latent space anymore. Moreover, if we use human demonstrations, e.g., recorded joint configurations, we may have to deal with over-fitting if we assume a one-to-one mapping between the human and the robot body and hence neglect the *correspondence problem* [11]. On the other hand, using the dimensionality reduction as an independent process can decrease the possible learning efficiency if we are not using the reward information for the identification of the subspace. Furthermore it is debatable how to consider the dimensionality reduction in the policy model such that overfitting does not occur while learning in the low dimensional subspace. This can be a problem if we use only the low dimensional manifold for sampling such that the following adaptations of the manifold would depend on the initialization.

1.1 Motivation

The issues described naturally lead to the question, whether it is possible to merge the idea of policy search and dimensionality reduction to one technique. This is the main question which shall be answered by this thesis. Our aim is to find a theory, based on Expectation Maximization (EM) policy search, which allows us to develop a reinforcement learning algorithm able to find a latent space to perform a directed exploration. Thus the samples, which are used for policy search, lie on the manifold created by the previous samples with the highest reward.

This means that the dimensionality reduction is no longer a pre-processing step or done independently (see Fig. 1.2). Thus our policy search algorithm would be able to find a satisfying policy and a low dimensional manifold together in one process. The ability to find such a latent space online is satisfying the constraint that in the most cases we do not have a training set of possible solutions for initializations and that we are uncertain about the integrity of our current samples with respect to the representation of the possible manifold. Additionally, we will be able to use the rewards in a direct way for suggesting the latent space.

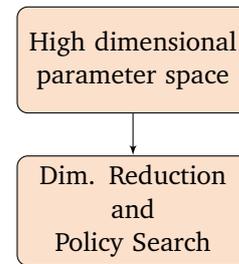


Figure 1.2.: The presented approach combines dimensionality reduction and policy search in one method.

1.2 Outline

The first chapter of this thesis will describe the required mathematical foundations and starts with a short introduction into the matrix variate normal distributions and EM algorithm. Subsequently, the theory behind PPCA will be highlighted, which serves as the main idea for the dimensionality reduction technique used later in this thesis and is based on EM. The discussion of the reinforcement learning framework of PoWER follows an overview of policy search from a probabilistic point of view and will be the basis of our theory for latent space reinforcement learning. The derivation of this theory is given in the second chapter, where the PePPeR Algorithm is introduced. The chapter is closed with a diagonal matrix extension for the PePPeR Algorithm. In the fourth chapter, we evaluate the introduced algorithm on a toy task where we need to learn inverse kinematics of a high-dimensional robot. Subsequently, we present the results for a real-world task on a NAO robot with the application of the PePPeR Algorithm. The goal was to let the robot lift one leg without losing balance. This thesis is closed with a conclusion and a short description of potential future work.

Parts of this thesis were submitted on the International Conference on Intelligent Robots and Systems (IROS) 2014 under the title *Latent Space Policy Search for Robotics* [12].

2 Foundation

2.1 Matrix Variate Normal Distributions

In this section we introduce matrix variate normal distributions such that we are able to use them for the formulation of the PePPeEr Algorithm later on. Furthermore, the use of them leads to an elegant way for the solve of several expectations that occur while using the EM algorithm.

The standard normal distribution $\mathcal{N}(\mu, \sigma^2)$, given by the probability density function (p.d.f.)

$$f(x) = (2\pi\sigma^2)^{-\frac{1}{2}} \exp\left\{-\frac{1}{2\sigma^2}(x-\mu)^2\right\}, \quad (2.1)$$

can be extended to vectors and is then called the *multivariate* normal distribution. The p.d.f. for a vector $\mathbf{x} \in \mathbb{R}^n$ is given by

$$f(\mathbf{x}) = (2\pi)^{-\frac{n}{2}} \det(\Sigma)^{-\frac{1}{2}} \exp\left\{-\frac{1}{2}(\mathbf{x}-\boldsymbol{\mu})\Sigma^{-1}(\mathbf{x}-\boldsymbol{\mu})^T\right\} \quad (2.2)$$

with a mean vector $\boldsymbol{\mu} \in \mathbb{R}^n$ and a covariance matrix $\Sigma \in \mathbb{R}^{n \times n}$, which is symmetric and positive definite. We denote this distribution $\mathcal{N}(\boldsymbol{\mu}, \Sigma)$. Now let us assume, that we have a set of p samples $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_p$, each vector sampled from such a multivariate normal distribution $\mathcal{N}(\boldsymbol{\mu}, \Sigma)$. If we organize this sample set as a matrix like

$$(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_p)^T = \begin{pmatrix} x_{11} & x_{12} & \cdots & x_{1n} \\ x_{21} & x_{22} & \cdots & x_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ x_{p1} & x_{p2} & \cdots & x_{pn} \end{pmatrix} \quad (2.3)$$

it is possible to regard this as one random matrix instead of several multivariate samples. Furthermore, it is possible to describe a distribution based on random matrices like a multivariate distribution. In fact the theory of matrix variate distributions can be traced back to multivariate distributions.

The following definitions and theorems are taken from [13], where the interested reader can find the proofs and a deeper discussion of the theory of matrix variate distributions.

To introduce the matrix variate normal distribution, we need first the definitions of the Kronecker product between two matrices and the vectorization of a matrix.

Definition 1 Given two matrices $\mathbf{A} \in \mathbb{R}^{m \times n}$ and $\mathbf{B} \in \mathbb{R}^{p \times q}$ the Kronecker product $\mathbf{A} \otimes \mathbf{B}$ is defined by

$$\mathbf{A} \otimes \mathbf{B} = \begin{pmatrix} a_{11}\mathbf{B} & a_{12}\mathbf{B} & \cdots & a_{1n}\mathbf{B} \\ a_{21}\mathbf{B} & a_{22}\mathbf{B} & \cdots & a_{2n}\mathbf{B} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1}\mathbf{B} & a_{m2}\mathbf{B} & \cdots & a_{mn}\mathbf{B} \end{pmatrix} \in \mathbb{R}^{mp \times nq}.$$

Definition 2 Given a matrix $(\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_n) = \mathbf{A} \in \mathbb{R}^{m \times n}$, the vectorization of this matrix $\text{vec}(\mathbf{A})$ is defined by

$$\text{vec}(\mathbf{A}) = \begin{pmatrix} \mathbf{a}_1 \\ \vdots \\ \mathbf{a}_n \end{pmatrix} \in \mathbb{R}^{nm}.$$

Now we can use the definition of the multivariate normal distribution to extend it to the matrix variate normal distributions.

Definition 3 We call a random matrix $\mathbf{X} \in \mathbb{R}^{p \times n}$ distributed with a matrix variate distribution $\mathcal{N}_{p,n}(\mathbf{M}, \Sigma \otimes \Psi)$, specified by the mean matrix $\mathbf{M} \in \mathbb{R}^{p \times n}$ and the covariance matrix $\Sigma \otimes \Psi$ with positive definite matrices $\Sigma \in \mathbb{R}^{p \times p}$ and $\Psi \in \mathbb{R}^{n \times n}$, if $\text{vec}(\mathbf{X}^T) \sim \mathcal{N}(\text{vec}(\mathbf{M}^T), \Sigma \otimes \Psi)$ holds.

In the case that Σ and Ψ are positive semidefinite matrices, an alternative definition has to be used which can be found in [13]. To distinguish between the multivariate normal distribution and the matrix variate distribution we will write $\mathbf{X} \sim \mathcal{N}_{p,n}(\mathbf{M}, \Sigma \otimes \Psi)$ in the case of matrix variate normal distributions for random matrices $\mathbf{X} \in \mathbb{R}^{p \times n}$. Furthermore it is worthwhile to see that in addition to the given definition above a random matrix \mathbf{X} is called matrix variate distributed with $\mathcal{N}_{p,n}(\mathbf{M}, \Sigma \otimes \Psi)$ if we have $\text{vec}(\mathbf{X}) \sim \mathcal{N}(\text{vec}(\mathbf{M}), \Psi \otimes \Sigma)$ under the described conditions.

Theorem 1 Let \mathbf{X} be a random matrix distributed with $\mathbf{X} \sim \mathcal{N}_{p,n}(\mathbf{M}, \Sigma \otimes \Psi)$, then the p.d.f. is defined by

$$f(\mathbf{X}) = (2\pi)^{-\frac{np}{2}} \det(\Sigma)^{-\frac{n}{2}} \det(\Psi)^{-\frac{p}{2}} \exp\left\{-\frac{1}{2} \text{trace}\left(\Sigma^{-1}(\mathbf{X} - \mathbf{M})\Psi^{-1}(\mathbf{X} - \mathbf{M})^T\right)\right\}$$

The proof for this theorem can be found in [13, Thm. 2.2.1]. In this thesis, we will also need the following theorem for the matrix variate second order expectation.

Theorem 2 Given a random matrix $\mathbf{X} \sim \mathcal{N}_{p,n}(\mathbf{M}, \Sigma \otimes \Psi)$ with $\mathbf{X} \in \mathbb{R}^{p \times n}$, where the distribution is specified by the mean matrix $\mathbf{M} \in \mathbb{R}^{p \times n}$ and the covariance with $\Sigma \in \mathbb{R}^{p \times p}$ and $\Psi \in \mathbb{R}^{n \times n}$. For the expectation of the term $\mathbf{X}^T \mathbf{A} \mathbf{X}$ with a matrix $\mathbf{A} \in \mathbb{R}^{p \times p}$ we can find that

$$E[\mathbf{X}^T \mathbf{A} \mathbf{X}] = \text{trace}(\Sigma \mathbf{A}^T) \Psi + \mathbf{M}^T \mathbf{A} \mathbf{M}.$$

This theorem simplifies various steps in the later proofs, as there is no need to consider each dimension of an action vector independently with the theory of multivariate normal distributions.

2.2 Expectation Maximization

In this section a short introduction is given into the Expectation Maximization (EM) algorithm based on the introduced notations which will follow closely [14, pp.450] and [1, pp.43]. The EM algorithm is used to determine the maximum likelihood solution of a stochastic model with latent variables, which are here denoted by \mathbf{z} . Let a data set be given by $\mathbf{A} = [\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_T]$ then the goal of EM is to maximize the probability of the likelihood

$$p_{\theta}(\mathbf{A}) = \prod_{t=1}^T p_{\theta}(\mathbf{a}_t) = \prod_{t=1}^T \int_{\mathcal{Z}} p_{\theta}(\mathbf{a}_t, \mathbf{z}) d\mathbf{z}, \quad (2.4)$$

where we assume independent and identically distributed data points \mathbf{a}_t . The distribution is determined by the parameters given by θ . Instead of the maximization of Eq. 2.4 we can maximize the log-likelihood given by

$$\ln p_{\theta}(\mathbf{A}) = \sum_{t=1}^T \ln \int_{\mathcal{Z}} p_{\theta}(\mathbf{a}_t, \mathbf{z}) d\mathbf{z}, \quad (2.5)$$

because the logarithm is a monotonically increasing function and so the maximum of $p_{\theta}(\mathbf{A})$ is the same as for $\ln p_{\theta}(\mathbf{A})$. Solving this equation for the parameters is difficult, while the logarithm is inside of the sum. Introducing the distribution $q(\mathbf{Z})$ and using the identity $p_{\theta}(\mathbf{A}) = \frac{p_{\theta}(\mathbf{A}, \mathbf{Z})}{p_{\theta}(\mathbf{Z}|\mathbf{A})}$ we get that

$$\begin{aligned} \ln p_{\theta}(\mathbf{A}) &= \int_{\mathcal{Z}} q(\mathbf{Z}) \ln p_{\theta}(\mathbf{A}) d\mathbf{Z} = \int_{\mathcal{Z}} q(\mathbf{Z}) \ln \frac{q(\mathbf{Z}) p_{\theta}(\mathbf{A}, \mathbf{Z})}{q(\mathbf{Z}) p_{\theta}(\mathbf{Z}|\mathbf{A})} d\mathbf{Z} \\ &= \int_{\mathcal{Z}} q(\mathbf{Z}) \ln \frac{p_{\theta}(\mathbf{A}, \mathbf{Z})}{q(\mathbf{Z})} d\mathbf{Z} - \int_{\mathcal{Z}} q(\mathbf{Z}) \ln \frac{p_{\theta}(\mathbf{Z}|\mathbf{A})}{q(\mathbf{Z})} d\mathbf{Z} \\ &= \mathcal{L}_{\theta}(\mathbf{A}) + \text{KL}(q(\mathbf{Z}) \| p_{\theta}(\mathbf{Z}|\mathbf{A})) \end{aligned} \quad (2.6)$$

holds (the verification can be found in [14, p.451]). The Kullback-Leibler (KL) divergence has the property that $\text{KL}(q \parallel p) \geq 0$ and $\text{KL}(q \parallel p) = 0$ if, and only if, the distributions q and p are equal. With this property we can now argue that the term $\mathcal{L}_\theta(q)$ is a lower bound of $\ln p_\theta(\mathbf{A})$. The two iterative steps of the EM algorithm is first to minimize the KL divergence by setting it to zero and then to maximize the lower bound $\mathcal{L}_\theta(q)$, to increase the probability of the data set $p_\theta(\mathbf{A})$.

To minimize the term $\text{KL}(q(\mathbf{Z}) \parallel p_\theta(\mathbf{Z}|\mathbf{A}))$ in the *expectation step* with respect to the old parameters, we have to update the distribution $q(\mathbf{Z})$ which means to equal $q(\mathbf{Z})$ and $p_\theta(\mathbf{Z}|\mathbf{A})$. Therefore the KL divergence is zero and $\ln p_\theta(\mathbf{A}) = \mathcal{L}_\theta(q)$. Now, in the *maximization step*, we can estimate the new parameters with

$$\begin{aligned} \theta_{new} &= \arg \max_{\theta} \mathcal{L}_\theta(q) = \arg \max_{\theta} \int_{\mathbb{Z}} q(\mathbf{Z}) \ln p_\theta(\mathbf{A}, \mathbf{Z}) d\mathbf{Z} - \int_{\mathbb{Z}} q(\mathbf{Z}) \ln q(\mathbf{Z}) d\mathbf{Z} \\ &= \arg \max_{\theta} \mathbb{E}_{q(\mathbf{Z})} [\ln p_\theta(\mathbf{A}, \mathbf{Z})] + \text{const} = \arg \max_{\theta} \Omega_\theta(q) + \text{const} \end{aligned} \quad (2.7)$$

where we can drop the constant term. From the *expectation step* we can now substitute $q(\mathbf{Z})$ with $p_\theta(\mathbf{Z}|\mathbf{A})$ and so the estimation of the new parameters can be done with Maximum Likelihood Estimation (MLE). Because the KL divergence can only be positive or zero, we can be sure that this maximization of the lower bound leads to an increase of the probability $\ln p_\theta(\mathbf{A})$ or at least to the same probability as before. This is due to the fact that we chose the KL divergence to be zero before maximization.

2.3 Dimensionality Reduction with Probabilistic Principal Component Analysis

A standard technique to handle high dimensional spaces is Principle Component Analysis (PCA) [6] which is able to find a low dimensional subspace for a data set. This method has become popular in a broad field of applications where spaces with the above described property have to be handled, i.e., meteorology [15], medicine [16] or face recognition [17]. A good example for the application of PCA is grasping in robotics. The human hand has around 25 DOF, so each hand configuration lies in a 25 dimensional space. But it has been found that about 7 dimension retain enough information to represent each hand configuration with only a small amount of information loss [18]. In this section, a brief overview of the probabilistic version of PCA is given. For the standard PCA we refer the reader to [19].

In contrast to PCA, where we assume to have the full data set, the Probabilistic Principal Component Analysis (PPCA) [7] is able to model the uncertainty in the case of missing data. Furthermore, the PPCA is formulated as a sampling technique from the low dimensional into the high dimensional space.

Assume a d -dimensional data point $\mathbf{x} \in \mathbb{R}^d$ in a high dimensional space, we can find a linear Gaussian model

$$\mathbf{x} = \mathbf{W}\mathbf{z} + \boldsymbol{\mu} + \boldsymbol{\epsilon} \quad (2.8)$$

where the low dimensional latent variable $\mathbf{z} \in \mathbb{R}^n$ is Gaussian distributed according to $p(\mathbf{z}) = \mathcal{N}(\mathbf{0}, \mathbf{I})$. This latent variable is projected with the transformation matrix $\mathbf{W} \in \mathbb{R}^{d \times n}$ into the high dimensional space. The columns of \mathbf{W} span the low-dimensional subspace, but these basis vectors do not have to be orthogonal to each other like in PCA. The mean of this low dimensional subspace is given by the vector $\boldsymbol{\mu} \in \mathbb{R}^d$. To model the uncertainty about the missing data, the isotropic error term $\boldsymbol{\epsilon} \in \mathbb{R}^d$ is used with the Gaussian distribution $p(\boldsymbol{\epsilon}) = \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$. This isotropic error leads to a noise added to the principle components, in contrast to PCA where these projections would lie directly on them. So the resulting Gaussian distribution of Eq. 2.8 is given by

$$p(\mathbf{x}) = \mathcal{N}(\boldsymbol{\mu}, \mathbf{W}\mathbf{W}^T + \sigma^2 \mathbf{I}). \quad (2.9)$$

As our model contains latent variables, the parameters \mathbf{W} , σ^2 and $\boldsymbol{\mu}$ can be determined for the maximization of the probability $p(\mathbf{X})$ with the iterative EM algorithm.

Because the ability to project samples from the manifold to the high dimensional space, PPCA seems to be ideal to be used as underlying theoretical foundation for the exploration method in our desired latent space RL method.

2.4 Policy Search

Most tasks for a robot can be formulated as a policy search problem. The goal of policy search is to find a control policy $\pi(\mathbf{s}_t, t)$ that provides an action $\mathbf{a}_t \in \mathbf{A}$ for every given state $\mathbf{s}_t \in \mathbf{S}$ and time t of the system. The execution of

this action on the robot results in a new state \mathbf{s}_{t+1} and a reward $r_t(\mathbf{s}_t, \mathbf{a}_t)$. Instead of a deterministic policy we will use a stochastic policy $\pi_\theta(\mathbf{a}_t|\mathbf{s}_t, t)$ given as a conditional probability distribution with parameters θ . Furthermore, we assume that the policy is linear in the form of $\mathbf{a}_t = \mathbf{M}_\theta \boldsymbol{\phi}(\mathbf{s}_t, t)$ where the matrix \mathbf{M}_θ depends on the given parameters and the time t . The feature vector $\boldsymbol{\phi}(\mathbf{s}_t, t) = (\phi_1(\mathbf{s}_t, t), \phi_2(\mathbf{s}_t, t), \dots, \phi_p(\mathbf{s}_t, t))^T$ contains basis functions $\phi_1, \phi_2, \dots, \phi_p$ which depend on the current time and the state. Examples for such features are time-dependent Gaussians. The goal of policy search is to determine the parameters θ s.t. the *expected return*

$$J(\theta) = \mathbb{E}_{p_\theta(\boldsymbol{\tau})} [R(\boldsymbol{\tau})] = \int_{\mathbb{T}} p_\theta(\boldsymbol{\tau}) R(\boldsymbol{\tau}) d\boldsymbol{\tau} \quad (2.10)$$

is maximized. This expectation integrates over all possible trajectories $\boldsymbol{\tau}$, given by the set \mathbb{T} . The trajectories are given by the sequence of states and actions $\boldsymbol{\tau} = [\mathbf{s}_{1:T+1}, \mathbf{a}_{1:T}]$. The reward of each trajectory is given by $R(\boldsymbol{\tau})$, which we assume to be the *accumulated immediate rewards* r_t with

$$R(\boldsymbol{\tau}) = \sum_{t=1}^T r_t(\mathbf{s}_t, \mathbf{a}_t) + r_{T+1}(\mathbf{s}_T), \quad (2.11)$$

where the reward r_{T+1} denotes the final reward for reaching the state \mathbf{s}_{T+1} . For the probability of a trajectory, we assume the *Markov property* and so we get

$$p_\theta(\boldsymbol{\tau}) = p(\mathbf{s}_1) \prod_{t=1}^T p(\mathbf{s}_{t+1}|\mathbf{s}_t, \mathbf{a}_t) \pi_\theta(\mathbf{a}_t|\mathbf{s}_t, t). \quad (2.12)$$

On the right side of the term we can find our policy that is to be determined, while the state probabilities on the left side will disappear in Eq. 2.14 later on.

2.5 The Reinforcement Learning Framework of PoWER

Policy Search methods, which are based on EM, formalize the problem statement as an inference problem with latent variables. In order to do so, the rewards have to be transformed into (improper) probability distributions such that each binary reward event is an (unnormalized) probability. In the case of a reward signal where a small number is good and a large bad, this can be easily done by using the exponential function with

$$r_t(\mathbf{s}_t, \mathbf{a}_t) = \exp\{-\lambda r'_t(\mathbf{s}_t, \mathbf{a}_t)\}, \quad (2.13)$$

where $r'_t(\mathbf{s}_t, \mathbf{a}_t)$ is the actual reward and λ a (temperature) parameter to adjust the exponential function. Here the exponential function ensures that the reward r_t is non-negative. From now on we will assume that our reward r_t is such a transformed reward. With EM we can now find a lower bound of the expected return from Eq. (2.10), which is, according to Kober and Peters [20], given by

$$\begin{aligned} L_{\theta_{\text{old}}}(\theta) &= \int_{\mathbb{T}} p_{\theta_{\text{old}}}(\boldsymbol{\tau}) R(\boldsymbol{\tau}) \ln p_\theta(\boldsymbol{\tau}) d\boldsymbol{\tau} \\ &= \mathbb{E}_{p_{\theta_{\text{old}}}(\boldsymbol{\tau})} \left[\left(\sum_{t=1}^T \ln \pi_\theta(\mathbf{a}_t|\mathbf{s}_t, t) \right) R(\boldsymbol{\tau}) \right] \\ &= \mathbb{E}_{p_{\theta_{\text{old}}}(\boldsymbol{\tau})} \left[\sum_{t=1}^T Q^\pi(\mathbf{s}_t, \mathbf{a}_t, t) \ln \pi_\theta(\mathbf{a}_t|\mathbf{s}_t, t) \right], \end{aligned} \quad (2.14)$$

where θ_{old} are the old parameters and θ the new ones, which we have to determine with Maximum Likelihood Estimation (MLE). The trajectories are here the latent variables while the rewards be the observations. The function $Q^\pi(\mathbf{s}_t, \mathbf{a}_t, t)$ is defined as the state-action value function given by

$$Q^\pi(\mathbf{s}_t, \mathbf{a}, t) = \mathbb{E} \left[\sum_{\tilde{t}=t}^T r_{\tilde{t}}(\mathbf{s}_{\tilde{t}}, \mathbf{a}_{\tilde{t}}) \mid \mathbf{s}_t = \mathbf{s}, \mathbf{a}_t = \mathbf{a} \right]. \quad (2.15)$$

We can estimate $Q^\pi(\mathbf{s}_t, \mathbf{a}, t)$ by a single rollout with

$$Q^\pi(\mathbf{s}_t^i, \mathbf{a}_t^i, t) \approx \sum_{\bar{i}=t}^T r_{\bar{i}}^i \quad (2.16)$$

where i denotes the index of the rollout. For the sake of simplicity, we will write Q_t^π instead of $Q^\pi(\mathbf{s}_t, \mathbf{a}, t)$. An advantage of this approach is, that we can obtain the parameters by weighted MLE update and so there is no need for a user-specific *learning rate*, which can be a crucial parameter in policy gradient algorithms [21]. The policy $\pi_\theta(\mathbf{a}_t|\mathbf{s}_t, t)$ can be modelled as a linear Gaussian model. Such policies are, for example, used in the 'Policy learning by Weighting Exploration with the Returns' (PoWER) algorithm [20].

In the stochastic policy of PoWER, given by $\mathbf{a}_t = (\mathbf{m} + \epsilon_t)^\top \boldsymbol{\phi}(\mathbf{s}, t)$ for each dimension of the control action, a Gaussian noise $\epsilon_t \sim \mathcal{N}(\mathbf{0}, \Sigma)$ is added to the current solution $\mathbf{m} \in \mathbb{R}^d$. The use of Eq. (2.14) leads to the update rule for PoWER given by

$$\mathbf{m}_{t+1} = \mathbf{m}_t + \left(\sum_{t=1}^T \mathbf{B}(\mathbf{s}_t, t) Q_t^\pi \right)^{-1} \cdot \left(\sum_{t=1}^T \mathbf{B}(\mathbf{s}_t, t) \epsilon_t Q_t^\pi \right) \quad (2.17)$$

with $\mathbf{B}(\mathbf{s}_t, t) = \boldsymbol{\phi}(\mathbf{s}_t, t) \boldsymbol{\phi}(\mathbf{s}_t, t)^\top \left(\boldsymbol{\phi}(\mathbf{s}_t, t)^\top \Sigma \boldsymbol{\phi}(\mathbf{s}_t, t) \right)^{-1}$. In the case of a non-static Σ , a similar update rule for Σ can be found in [20]. With the current framework we are not able to introduce latent variables in our stochastic policy so far. With these latent variables we would be able to use a (stochastic) linear model of a low dimensional subspace for the exploration in policy search. This extension, which is the main contribution of this thesis, will be discussed in the next chapter.



3 Reinforcement Learning in Spaces with Low Intrinsic Dimensionality

3.1 Using the PoWER Framework for Estimating Policies with Latent Variables

The EM framework of PoWER, given in Eq. (2.14), in association with weighted MLE can be exploited for the use of stochastic policies with latent variables. In spirit of the common notations, \mathbf{z} is used for latent variables and $\mathbf{Z} \in \mathbb{R}^{m \times n}$ for latent matrices. Furthermore, for a consistent notation, $\pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t, t)$ is written as $p_{\theta}(\mathbf{a}_t | \mathbf{s}_t, t)$. In order to add latent variables into the policy, we make use of the marginalization rule (see [14]) and define $p_{\theta}(\mathbf{a}_t | \mathbf{s}_t, t) = \int_{\mathbb{Z}} p_{\theta}(\mathbf{a}_t, \mathbf{z} | \mathbf{s}_t, t) d\mathbf{z}$. If EM is applied a new lower bound for Eq. (2.14) can be found with

$$\begin{aligned} \mathcal{L}_{\theta_{\text{old}}}(\theta) &= \mathbb{E}_{p_{\theta_{\text{old}}}(\tau)} \left[\sum_{t=1}^T Q^{\pi}(\mathbf{s}_t, \mathbf{a}_t, t) \ln \int_{\mathbb{Z}} p_{\theta}(\mathbf{a}_t, \mathbf{z} | \mathbf{s}_t, t) d\mathbf{z} \right] \\ &\geq \mathbb{E}_{p_{\theta_{\text{old}}}(\tau)} \left[\sum_{t=1}^T Q^{\pi}(\mathbf{s}_t, \mathbf{a}_t, t) \mathbb{E}_{q(\mathbf{z} | \mathbf{a}_t, \mathbf{s}_t)} [\ln p_{\theta}(\mathbf{a}_t, \mathbf{z} | \mathbf{s}_t, t)] \right] \\ &= \mathcal{L}'_{\theta_{\text{old}}, q}(\theta), \end{aligned} \tag{3.1}$$

where the distribution

$$q(\mathbf{z} | \mathbf{a}_t, \mathbf{s}_t) = \frac{p_{\theta_{\text{old}}}(\mathbf{a}_t, \mathbf{z} | \mathbf{s}_t, t)}{p_{\theta_{\text{old}}}(\mathbf{a}_t | \mathbf{s}_t, t)} \tag{3.2}$$

is given by the posterior distribution $p_{\theta_{\text{old}}}(\mathbf{z} | \mathbf{a}_t, \mathbf{s}_t)$ of the latent variables. It is worth to recognize that we have applied EM twice in this lower bound. First to derive the policy update by weighted MLE and second, to update the joint distribution $p_{\theta}(\mathbf{a}_t, \mathbf{z} | \mathbf{s}_t, t)$. This lower bound can be used for any latent variable model for stochastic policies. In the next sections, a model for the stochastic policy, which uses latent variables, to perform dimensionality reduction will be discussed.

3.2 Low Dimensional and Stochastic Policy Models

If we investigate the PoWER algorithm, we notice that this algorithms handles each dimension of the action space separately. But this seems only reasonable if we have to fulfill a task where the actions, i.e. the angles of a robot arm or hand, are independent for each time step. The fact is, that in reality we have a lot of dependencies between actions for the most common tasks, however this dependencies may not be obvious. In our algorithm, we want to take advantage of this hidden structure in the action space for a faster reinforcement learning in policy search tasks. Since we assume that we start learning without any knowledge or already sampled configurations, we have to learn the structure of this hidden space online. With the described exploitation for latent variables in Eq.(3.1) we are able to learn the properties of the hidden subspace and the optimal policy for the robot with one unified algorithm. Revisiting the model of PPCA given by

$$\mathbf{x} = \mathbf{W}\mathbf{z} + \boldsymbol{\mu} + \boldsymbol{\epsilon}, \tag{3.3}$$

where the axes of the low dimensional subspace are given by the columns of \mathbf{W} , we can find a stochastic policy for uncovering a hidden structure in the action space. This stochastic policy can be found with

$$\mathbf{a} = (\mathbf{W}\mathbf{Z}^T + \mathbf{M} + \mathbf{E}) \boldsymbol{\phi} = \mathbf{W}(\mathbf{Z}^T \boldsymbol{\phi}) + \mathbf{M}\boldsymbol{\phi} + \mathbf{E}\boldsymbol{\phi}, \tag{3.4}$$

and the random matrices $\mathbf{E} \sim \mathcal{N}_{d,p}(\mathbf{0}, \sigma^2 \mathbf{I})$ and $\mathbf{Z}^T \sim \mathcal{N}_{n,p}(\mathbf{0}, \mathbf{I})$, while $\mathbf{M} \in \mathbb{R}^{d \times p}$ is the mean matrix. From now on we will use d for the number of dimensions of the high dimensional action space, n for the number of the low

dimensional space and p for the number of parameters respectively the size of the feature vector $\phi \in \mathbb{R}^p$. As was mentioned before we will write ϕ instead of $\phi(s_t, t)$ for simplification, while keeping the dependencies in mind. Now let us take a closer look at the matrices and their intuition: Like in PPCA the matrix $\mathbf{W} \in \mathbb{R}^{d \times n}$ is a linear transformation from the hidden subspace into the high dimensional action space. The term $(\mathbf{Z}^T \phi)$, which corresponds to a low dimensional action vector, can be seen as the representation of the high dimensional action in the subspace, whose mean action is given by $\mathbf{M}\phi$. Finally, the term $\mathbf{E}\phi$ encodes the uncertainty as an isotropic noise. In a later section, we will introduce a generalization of this model where this term is modelled by a diagonal covariance matrix where we use one parameter per action dimension. In terms of policy search, we can see the term $\mathbf{W}(\mathbf{Z}^T \phi)$ as an directed exploration noise and $\mathbf{E}\phi$ as an undirected exploration noise. The parameters for this model are \mathbf{W}, \mathbf{M} and σ^2 .

3.3 Linear Model Properties

In this section we investigate the properties of the probability distributions which arise as a result from the above given model. The expectation of our model in Eq. (3.3) can be found with

$$\begin{aligned} \mathbb{E}[\mathbf{a}] &= \mathbb{E}[\mathbf{W}(\mathbf{Z}^T \phi) + \mathbf{M}\phi + \mathbf{E}\phi] \\ &= \underbrace{\mathbb{E}[\mathbf{W}(\mathbf{Z}^T \phi)]}_{=0} + \mathbf{M}\phi + \underbrace{\mathbb{E}[\mathbf{E}\phi]}_{=0} \\ &= \mathbf{M}\phi \end{aligned} \quad (3.5)$$

while the covariance of the model can be written as

$$\begin{aligned} \text{cov}(\mathbf{a}) &= \mathbb{E} \left[\left((\mathbf{W}(\mathbf{Z}^T \phi) + \mathbf{M}\phi + \mathbf{E}\phi) - \mathbb{E}[\mathbf{a}] \right) \left((\mathbf{W}(\mathbf{Z}^T \phi) + \mathbf{M}\phi + \mathbf{E}\phi) - \mathbb{E}[\mathbf{a}] \right)^T \right] \\ &= \mathbb{E} \left[(\mathbf{W}(\mathbf{Z}^T \phi) + \mathbf{E}\phi) (\mathbf{W}(\mathbf{Z}^T \phi) + \mathbf{E}\phi)^T \right] \\ &= \mathbb{E} \left[\mathbf{W}\mathbf{Z}^T \phi \phi^T \mathbf{Z}\mathbf{W}^T \right] + \underbrace{\mathbb{E} \left[\mathbf{W}\mathbf{Z}^T \phi \phi^T \mathbf{M}^T \right]}_{=0} + \underbrace{\mathbb{E} \left[\mathbf{W}\mathbf{Z}^T \phi \phi^T \mathbf{E}^T \right]}_{=0} + \underbrace{\mathbb{E} \left[\mathbf{E}\phi \phi^T \mathbf{Z}\mathbf{W}^T \right]}_{=0} + \mathbb{E} \left[\mathbf{E}\phi \phi^T \mathbf{E}^T \right] \\ &= \mathbb{E} \left[\mathbf{W}\mathbf{Z}^T \phi \phi^T \mathbf{Z}\mathbf{W}^T \right] + \mathbb{E} \left[\mathbf{E}\phi \phi^T \mathbf{E}^T \right]. \end{aligned} \quad (3.6)$$

At this point, we have to use Thm. 2 for the expectation of matrix variate distribution. If we apply this theorem we get

$$\begin{aligned} \text{cov}(\mathbf{a}) &= \mathbb{E} \left[\mathbf{W}\mathbf{Z}^T \phi \phi^T \mathbf{Z}\mathbf{W}^T \right] + \mathbb{E} \left[\mathbf{E}\phi \phi^T \mathbf{E}^T \right] \\ &= \mathbf{W} \text{trace}(\phi \phi^T) \mathbf{I} \mathbf{W}^T + \text{trace}(\phi \phi^T) \sigma^2 \mathbf{I} \\ &= \text{trace}(\phi \phi^T) (\mathbf{W}\mathbf{W}^T + \sigma^2 \mathbf{I}) \end{aligned} \quad (3.7)$$

for the covariance of our linear model. Thus the distribution for \mathbf{a} is given as

$$p(\mathbf{a}) = \mathcal{N}(\mathbf{M}\phi, \text{trace}(\phi \phi^T) (\mathbf{W}\mathbf{W}^T + \sigma^2 \mathbf{I})). \quad (3.8)$$

If we revisit our model equation, we can find that our latent variables are given by the random matrix \mathbf{Z} . At a first glance, it seems obvious to use the distribution $p(\mathbf{a}|\mathbf{Z})$ as posterior distribution for EM. But in order to apply Bayes theorem for Gaussian variables (see or Appendix A.2 or [14, p.93]) in an easy way, we need vectors as latent variables instead of matrices. For this simplification, we can use the term $\mathbf{Z}^T \phi$, whose result is a vector, as latent variable. Hence we have to estimate the posterior distribution $p(\mathbf{Z}^T \phi | \mathbf{a})$. For the properties of the distribution $p(\mathbf{Z}^T \phi)$, we can use the fact that $\mathbb{E}[\mathbf{Z}^T \phi] = 0$ holds and so the covariance of this distribution is reduced to the term $\mathbb{E}[\mathbf{Z}^T \phi \phi^T \mathbf{Z}]$, where Thm. 2 can be applied again and we get

$$p(\mathbf{Z}^T \phi) = \mathcal{N}(\mathbf{0}, \text{trace}(\phi \phi^T) \mathbf{I}) \quad (3.9)$$

as a multivariate normal distribution. Given the model Eq. and its distribution in (3.8) we can find that

$$p(\mathbf{a} | \mathbf{Z}^T \phi) = \mathcal{N}(\mathbf{W}(\mathbf{Z}^T \phi) + \mathbf{M}\phi, \sigma^2 \text{trace}(\phi \phi^T) \mathbf{I}) \quad (3.10)$$

holds. With these distributions and Bayes theorem for Gaussian variables we can derivate the posterior distribution with

$$p(\mathbf{Z}^T \phi | \mathbf{a}) = \mathcal{N}(\mathbf{C}\mathbf{W}^T(\mathbf{a} - \mathbf{M}\phi), \mathbf{C}\sigma^2 \text{trace}(\phi \phi^T)) \quad (3.11)$$

where $\mathbf{C} = (\sigma^2 \mathbf{I} + \mathbf{W}^T \mathbf{W})^{-1}$.

3.4 Expectation and Maximization of the Parameters

With the distributions given in the previous section, we can now define the expectation and the maximization steps. Instead of the so far used notation $\mathbb{E}_{p_{\theta_{\text{old}}}(Z^T \phi | \mathbf{a})} [Z^T \phi]$ for the expectation of the latent variables, we will use the notation $\mathbb{E}_{Z^T \phi | \mathbf{a}} [Z^T \phi]$.

3.4.1 Expectation

In the expectation step of the EM algorithm, we have to estimate the parameters of the probability distribution $p_{\theta}(Z^T \phi | \mathbf{a})$. To do so, we can use the old parameters θ_{old} from the previous EM iteration. We can directly estimate this step, because we only need the expectation of this distribution, with

$$\mathbb{E}_{Z^T \phi | \mathbf{a}} [Z^T \phi] = \mathbf{C} \mathbf{W}^T (\mathbf{a} - \mathbf{M} \phi) \quad (3.12)$$

and

$$\mathbb{E}_{Z^T \phi | \mathbf{a}} [Z^T \phi (Z^T \phi)^T] = \mathbf{C} \sigma^2 \text{trace}(\phi \phi^T) + \mathbb{E}_{Z^T \phi | \mathbf{a}} [Z^T \phi] \cdot \mathbb{E}_{Z^T \phi | \mathbf{a}} [Z^T \phi]^T \quad (3.13)$$

with $\mathbf{C} = (\sigma^2 \mathbf{I} + \mathbf{W}^T \mathbf{W})^{-1}$ like above.

3.4.2 Maximization

Now we have to find the parameters for which we maximize the expected return. To do so, we have to maximize our lower bound given in Eq. (3.1) which is equivalent to doing a weighted MLE for our parameters θ . The parameters of our model are \mathbf{W}, \mathbf{M} and σ^2 like mentioned in Sec. 3.2.

A common way to estimate the mean parameter, here \mathbf{M} , is to do a MLE. Thus, \mathbf{M} is given by

$$\begin{aligned} \frac{\partial \ln p(\mathbf{a})}{\partial \mathbf{M}} &= \frac{\partial}{\partial \mathbf{M}} \left(-\frac{1}{2} (\mathbf{a} - \mathbf{M} \phi)^T \mathbf{D}^{-1} (\mathbf{a} - \mathbf{M} \phi) \right) \\ &= \frac{\partial}{\partial \mathbf{M}} \left(-\frac{1}{2} (-2 \mathbf{a}^T \mathbf{D}^{-1} \mathbf{M} \phi + \phi^T \mathbf{M}^T \mathbf{D}^{-1} \mathbf{M} \phi) \right) \\ &\stackrel{A.1}{=} \left(-\frac{1}{2} \left(-2 \mathbf{D}^{-1} \mathbf{a} \phi^T + \frac{\partial \phi^T \mathbf{M}^T \mathbf{D}^{-1} \mathbf{M} \phi}{\partial \mathbf{M}} \right) \right) \\ &\stackrel{A.2}{=} \left(-\frac{1}{2} (-2 \mathbf{D}^{-1} \mathbf{a} \phi^T + 2 \mathbf{D}^{-1} \mathbf{M} \phi \phi^T) \right) \\ &= (\mathbf{D}^{-1} (\mathbf{a} \phi^T - \mathbf{M} \phi \phi^T)) \end{aligned} \quad (3.14)$$

where $\mathbf{D} = (\text{trace}(\phi \phi^T) (\mathbf{W} \mathbf{W}^T + \sigma^2 \mathbf{I})) = \mathbf{D}^T$ holds. In this derivatives we made use of given rules for matrix differentiation at scalar forms, which can be found in the Appendix A.1. Now we can insert this result into the term

$$0 = \mathbb{E}_{\tau} \left[\sum_{t=1}^T \frac{\partial \ln p(\mathbf{a}) Q_t^{\pi}}{\partial \mathbf{M}} \right] \quad (3.15)$$

which is the weighted MLE step and if we solve this equation for \mathbf{M} we find

$$\begin{aligned} 0 &= \mathbb{E}_{\tau} \left[\sum_{t=1}^T \frac{\partial \ln p(\mathbf{a}) Q_t^{\pi}}{\partial \mathbf{M}} \right] \\ \Leftrightarrow 0 &= \mathbb{E}_{\tau} \left[\sum_{t=1}^T \mathbf{D}^{-1} (\mathbf{a} \phi^T - \mathbf{M} \phi \phi^T) Q_t^{\pi} \right] \\ \Leftrightarrow 0 &= (\mathbf{W} \mathbf{W}^T + \sigma^2 \mathbf{I})^{-1} \mathbb{E}_{\tau} \left[\sum_{t=1}^T \frac{1}{\text{trace}(\phi \phi^T)} (\mathbf{a} \phi^T - \mathbf{M} \phi \phi^T) Q_t^{\pi} \right] \\ \Leftrightarrow \mathbf{M} \cdot \mathbb{E}_{\tau} \left[\sum_{t=1}^T \frac{1}{\text{trace}(\phi \phi^T)} \phi \phi^T Q_t^{\pi} \right] &= \mathbb{E}_{\tau} \left[\sum_{t=1}^T \frac{1}{\text{trace}(\phi \phi^T)} \mathbf{a} \phi^T Q_t^{\pi} \right] \\ \Leftrightarrow \mathbf{M} &= \mathbb{E}_{\tau} \left[\sum_{t=1}^T \frac{1}{\text{trace}(\phi \phi^T)} \mathbf{a} \phi^T Q_t^{\pi} \right] \cdot \left(\mathbb{E}_{\tau} \left[\sum_{t=1}^T \frac{1}{\text{trace}(\phi \phi^T)} \phi \phi^T Q_t^{\pi} \right] \right)^{-1}. \end{aligned} \quad (3.16)$$

For the maximization of \mathbf{W} and σ^2 , we have to use our lower bound, where we can find for $\ln p(\mathbf{a}, \mathbf{Z})$ that

$$\begin{aligned}
\frac{\partial \ln p(\mathbf{a}, \mathbf{Z}^T \boldsymbol{\phi})}{\partial \boldsymbol{\theta}_{new}} &= \frac{\partial}{\partial \boldsymbol{\theta}_{new}} (\ln p(\mathbf{a} | \mathbf{Z}^T \boldsymbol{\phi}) + \ln p(\mathbf{Z}^T \boldsymbol{\phi})) \\
&= \frac{\partial}{\partial \boldsymbol{\theta}_{new}} \left(\ln \left(\sqrt{\det(2\pi\sigma^2 \text{trace}(\boldsymbol{\phi}\boldsymbol{\phi}^T) \mathbf{I})}^{-1} \right) + \right. \\
&\quad \left. - \frac{1}{2\sigma^2 \text{trace}(\boldsymbol{\phi}\boldsymbol{\phi}^T)} (\mathbf{a} - \mathbf{WZ}^T \boldsymbol{\phi} - \mathbf{M}\boldsymbol{\phi})^T (\mathbf{a} - \mathbf{WZ}^T \boldsymbol{\phi} - \mathbf{M}\boldsymbol{\phi}) \right) \\
&= \frac{\partial}{\partial \boldsymbol{\theta}_{new}} \left(-\frac{d}{2} \ln(2\pi \text{trace}(\boldsymbol{\phi}\boldsymbol{\phi}^T)) - \frac{d}{2} \ln(\sigma^2) \right. \\
&\quad \left. - \frac{1}{2\sigma^2 \text{trace}(\boldsymbol{\phi}\boldsymbol{\phi}^T)} (\mathbf{a} - \mathbf{WZ}^T \boldsymbol{\phi} - \mathbf{M}\boldsymbol{\phi})^T (\mathbf{a} - \mathbf{WZ}^T \boldsymbol{\phi} - \mathbf{M}\boldsymbol{\phi}) \right)
\end{aligned} \tag{3.17}$$

holds if we differentiate this term w.r.t. one of the parameters. Doing so for the parameter \mathbf{W} and using the derivatives of the scalar forms we get

$$\begin{aligned}
\frac{\partial \ln p(\mathbf{a}, \mathbf{Z}^T \boldsymbol{\phi})}{\partial \mathbf{W}} &\stackrel{3.17}{=} \frac{\partial}{\partial \mathbf{W}} \left(-\frac{d}{2} \ln(2\pi \text{trace}(\boldsymbol{\phi}\boldsymbol{\phi}^T)) - \frac{d}{2} \ln(\sigma^2) + \right. \\
&\quad \left. - \frac{1}{2\sigma^2 \text{trace}(\boldsymbol{\phi}\boldsymbol{\phi}^T)} (\mathbf{a} - \mathbf{WZ}^T \boldsymbol{\phi} - \mathbf{M}\boldsymbol{\phi})^T (\mathbf{a} - \mathbf{WZ}^T \boldsymbol{\phi} - \mathbf{M}\boldsymbol{\phi}) \right) \\
&= \frac{\partial}{\partial \mathbf{W}} \left(-\frac{1}{2\sigma^2 \text{trace}(\boldsymbol{\phi}\boldsymbol{\phi}^T)} (\mathbf{a} - \mathbf{WZ}^T \boldsymbol{\phi} - \mathbf{M}\boldsymbol{\phi})^T (\mathbf{a} - \mathbf{WZ}^T \boldsymbol{\phi} - \mathbf{M}\boldsymbol{\phi}) \right) \\
&= \frac{\partial}{\partial \mathbf{W}} \left(-\frac{1}{2\sigma^2 \text{trace}(\boldsymbol{\phi}\boldsymbol{\phi}^T)} \right. \\
&\quad \left. (-\mathbf{a}^T \mathbf{WZ}^T \boldsymbol{\phi} - \boldsymbol{\phi}^T \mathbf{Z} \mathbf{W}^T \mathbf{a} + \boldsymbol{\phi}^T \mathbf{Z} \mathbf{W}^T \mathbf{WZ}^T \boldsymbol{\phi} + \boldsymbol{\phi}^T \mathbf{Z} \mathbf{W}^T \mathbf{M}\boldsymbol{\phi} + \boldsymbol{\phi}^T \mathbf{M}^T \mathbf{WZ}^T \boldsymbol{\phi}) \right) \\
&= \frac{\partial}{\partial \mathbf{W}} \left(-\frac{1}{2\sigma^2 \text{trace}(\boldsymbol{\phi}\boldsymbol{\phi}^T)} (-2\mathbf{a}^T \mathbf{WZ}^T \boldsymbol{\phi} + \boldsymbol{\phi}^T \mathbf{Z} \mathbf{W}^T \mathbf{WZ}^T \boldsymbol{\phi} + 2\boldsymbol{\phi}^T \mathbf{M}^T \mathbf{WZ}^T \boldsymbol{\phi}) \right) \\
&\stackrel{A.1}{=} -\frac{1}{2\sigma^2 \text{trace}(\boldsymbol{\phi}\boldsymbol{\phi}^T)} \left(-2\mathbf{a} (\mathbf{Z}^T \boldsymbol{\phi})^T + \frac{\partial \boldsymbol{\phi}^T \mathbf{Z} \mathbf{W}^T \mathbf{WZ}^T \boldsymbol{\phi}}{\partial \mathbf{W}} + 2\mathbf{M}\boldsymbol{\phi} (\mathbf{Z}^T \boldsymbol{\phi})^T \right) \\
&\stackrel{A.2}{=} -\frac{1}{\sigma^2 \text{trace}(\boldsymbol{\phi}\boldsymbol{\phi}^T)} \left(-\mathbf{a} (\mathbf{Z}^T \boldsymbol{\phi})^T + \mathbf{WZ}^T \boldsymbol{\phi} (\mathbf{Z}^T \boldsymbol{\phi})^T + \mathbf{M}\boldsymbol{\phi} (\mathbf{Z}^T \boldsymbol{\phi})^T \right).
\end{aligned} \tag{3.18}$$

When we put this result into the weighted MLE for our lower bound we have

$$\begin{aligned}
0 &= \mathbb{E}_\tau \left[\sum_{t=1}^T \frac{\partial}{\partial \mathbf{W}} \mathbb{E}_{\mathbf{Z}^T \boldsymbol{\phi} | \mathbf{a}} [\ln p(\mathbf{a}, \mathbf{Z}^T \boldsymbol{\phi})] Q_t^\pi \right] \\
&= \mathbb{E}_\tau \left[\sum_{t=1}^T \mathbb{E}_{\mathbf{Z}^T \boldsymbol{\phi} | \mathbf{a}} \left[-\frac{1}{\sigma^2 \text{trace}(\boldsymbol{\phi}\boldsymbol{\phi}^T)} \left(-\mathbf{a} (\mathbf{Z}^T \boldsymbol{\phi})^T + \mathbf{WZ}^T \boldsymbol{\phi} (\mathbf{Z}^T \boldsymbol{\phi})^T + \mathbf{M}\boldsymbol{\phi} (\mathbf{Z}^T \boldsymbol{\phi})^T \right) \right] Q_t^\pi \right] \\
&= \mathbb{E}_\tau \left[\sum_{t=1}^T \mathbb{E}_{\mathbf{Z}^T \boldsymbol{\phi} | \mathbf{a}} \left[\frac{1}{\text{trace}(\boldsymbol{\phi}\boldsymbol{\phi}^T)} \left(\mathbf{a} (\mathbf{Z}^T \boldsymbol{\phi})^T - \mathbf{M}\boldsymbol{\phi} (\mathbf{Z}^T \boldsymbol{\phi})^T \right) \right] Q_t^\pi \right] - \mathbb{E}_\tau \left[\sum_{t=1}^T \mathbb{E}_{\mathbf{Z}^T \boldsymbol{\phi} | \mathbf{a}} \left[\frac{1}{\text{trace}(\boldsymbol{\phi}\boldsymbol{\phi}^T)} \mathbf{WZ}^T \boldsymbol{\phi} (\mathbf{Z}^T \boldsymbol{\phi})^T \right] Q_t^\pi \right]
\end{aligned} \tag{3.19}$$

where we can eliminate σ^2 because it does not depend on an index, i.e., it is a constant. Solving this equation for \mathbf{W} , results in the following update equation

$$\begin{aligned}
& \mathbb{E}_\tau \left[\sum_{t=1}^T \mathbb{E}_{\mathbf{Z}^T \phi | \mathbf{a}} \left[\frac{1}{\text{trace}(\phi \phi^T)} \mathbf{W} \mathbf{Z}^T \phi (\mathbf{Z}^T \phi)^T \right] Q_t^\pi \right] = \mathbb{E}_\tau \left[\sum_{t=1}^T \mathbb{E}_{\mathbf{Z}^T \phi | \mathbf{a}} \left[\frac{1}{\text{trace}(\phi \phi^T)} (\mathbf{a} (\mathbf{Z}^T \phi)^T - \mathbf{M} \phi (\mathbf{Z}^T \phi)^T) \right] Q_t^\pi \right] \\
& \Leftrightarrow \mathbf{W} \mathbb{E}_\tau \left[\sum_{t=1}^T \mathbb{E}_{\mathbf{Z}^T \phi | \mathbf{a}} \left[\frac{1}{\text{trace}(\phi \phi^T)} \mathbf{Z}^T \phi (\mathbf{Z}^T \phi)^T \right] Q_t^\pi \right] = \mathbb{E}_\tau \left[\sum_{t=1}^T \mathbb{E}_{\mathbf{Z}^T \phi | \mathbf{a}} \left[\frac{1}{\text{trace}(\phi \phi^T)} (\mathbf{a} (\mathbf{Z}^T \phi)^T - \mathbf{M} \phi (\mathbf{Z}^T \phi)^T) \right] Q_t^\pi \right] \\
& \Leftrightarrow \mathbf{W} = \mathbb{E}_\tau \left[\sum_{t=1}^T \mathbb{E}_{\mathbf{Z}^T \phi | \mathbf{a}} \left[\frac{1}{\text{trace}(\phi \phi^T)} (\mathbf{a} (\mathbf{Z}^T \phi)^T - \mathbf{M} \phi (\mathbf{Z}^T \phi)^T) \right] Q_t^\pi \right] \cdot \\
& \quad \left(\mathbb{E}_\tau \left[\sum_{t=1}^T \mathbb{E}_{\mathbf{Z}^T \phi | \mathbf{a}} \left[\frac{1}{\text{trace}(\phi \phi^T)} \mathbf{Z}^T \phi (\mathbf{Z}^T \phi)^T \right] Q_t^\pi \right] \right)^{-1} \\
& \Leftrightarrow \mathbf{W} = \mathbb{E}_\tau \left[\sum_{t=1}^T \frac{1}{\text{trace}(\phi \phi^T)} (\mathbf{a} \mathbb{E}_{\mathbf{Z}^T \phi | \mathbf{a}} [\mathbf{Z}^T \phi]^T - \mathbf{M} \phi \mathbb{E}_{\mathbf{Z}^T \phi | \mathbf{a}} [\mathbf{Z}^T \phi]^T) Q_t^\pi \right] \cdot \\
& \quad \left(\mathbb{E}_\tau \left[\sum_{t=1}^T \frac{1}{\text{trace}(\phi \phi^T)} \mathbb{E}_{\mathbf{Z}^T \phi | \mathbf{a}} [\mathbf{Z}^T \phi (\mathbf{Z}^T \phi)^T] Q_t^\pi \right] \right)^{-1} \\
& \Leftrightarrow \mathbf{W} = \mathbb{E}_\tau \left[\sum_{t=1}^T \frac{1}{\text{trace}(\phi \phi^T)} (\mathbf{a} - \mathbf{M} \phi) \mathbb{E}_{\mathbf{Z}^T \phi | \mathbf{a}} [\mathbf{Z}^T \phi]^T Q_t^\pi \right] \cdot \left(\mathbb{E}_\tau \left[\sum_{t=1}^T \frac{1}{\text{trace}(\phi \phi^T)} \mathbb{E}_{\mathbf{Z}^T \phi | \mathbf{a}} [\mathbf{Z}^T \phi (\mathbf{Z}^T \phi)^T] Q_t^\pi \right] \right)^{-1}. \tag{3.20}
\end{aligned}$$

In an analogous manner, we can find the derivative w.r.t. σ^2 with

$$\begin{aligned}
& \frac{\partial \ln p(\mathbf{a}, \mathbf{Z}^T \phi)}{\partial \sigma^2} \stackrel{3.17}{=} \frac{\partial}{\partial \sigma^2} \left(-\frac{d}{2} \ln(\sigma^2) - \frac{1}{2\sigma^2 \text{trace}(\phi \phi^T)} (\mathbf{a} - \mathbf{W} \mathbf{Z}^T \phi - \mathbf{M} \phi)^T (\mathbf{a} - \mathbf{W} \mathbf{Z}^T \phi - \mathbf{M} \phi) \right) \\
& = -\frac{d}{2\sigma^2} + \frac{1}{2(\sigma^2)^2 \text{trace}(\phi \phi^T)} (\mathbf{a} - \mathbf{W} \mathbf{Z}^T \phi - \mathbf{M} \phi)^T (\mathbf{a} - \mathbf{W} \mathbf{Z}^T \phi - \mathbf{M} \phi) \tag{3.21}
\end{aligned}$$

and if we plug this result into the lower bound we get

$$\begin{aligned}
0 &= \mathbb{E}_\tau \left[\sum_{t=1}^T \frac{\partial}{\partial \sigma^2} \mathbb{E}_{\mathbf{Z}^T \phi | \mathbf{a}} [\ln p(\mathbf{a}, \mathbf{Z}^T \phi)] Q_t^\pi \right] \\
&= \mathbb{E}_\tau \left[\sum_{t=1}^T \mathbb{E}_{\mathbf{Z}^T \phi | \mathbf{a}} \left[-\frac{d}{2\sigma^2} + \frac{1}{2(\sigma^2)^2 \text{trace}(\phi \phi^T)} (\mathbf{a} - \mathbf{W} \mathbf{Z}^T \phi - \mathbf{M} \phi)^T (\mathbf{a} - \mathbf{W} \mathbf{Z}^T \phi - \mathbf{M} \phi) \right] Q_t^\pi \right] \\
&\Leftrightarrow \mathbb{E}_\tau \left[\sum_{t=1}^T \frac{d}{2\sigma^2} Q_t^\pi \right] = \mathbb{E}_\tau \left[\sum_{t=1}^T \mathbb{E}_{\mathbf{Z}^T \phi | \mathbf{a}} \left[\frac{1}{2(\sigma^2)^2 \text{trace}(\phi \phi^T)} (\mathbf{a} - \mathbf{W} \mathbf{Z}^T \phi - \mathbf{M} \phi)^T (\mathbf{a} - \mathbf{W} \mathbf{Z}^T \phi - \mathbf{M} \phi) \right] Q_t^\pi \right] \\
&\Leftrightarrow \frac{d}{2\sigma^2} \mathbb{E}_\tau \left[\sum_{t=1}^T Q_t^\pi \right] = \frac{1}{2(\sigma^2)^2} \mathbb{E}_\tau \left[\sum_{t=1}^T \frac{1}{\text{trace}(\phi \phi^T)} \mathbb{E}_{\mathbf{Z}^T \phi | \mathbf{a}} [(\mathbf{a} - \mathbf{W} \mathbf{Z}^T \phi - \mathbf{M} \phi)^T (\mathbf{a} - \mathbf{W} \mathbf{Z}^T \phi - \mathbf{M} \phi)] Q_t^\pi \right] \\
&\Leftrightarrow \sigma^2 \mathbb{E}_\tau \left[\sum_{t=1}^T Q_t^\pi \right] = \frac{1}{d} \mathbb{E}_\tau \left[\sum_{t=1}^T \frac{1}{\text{trace}(\phi \phi^T)} \mathbb{E}_{\mathbf{Z}^T \phi | \mathbf{a}} [(\mathbf{a} - \mathbf{W} \mathbf{Z}^T \phi - \mathbf{M} \phi)^T (\mathbf{a} - \mathbf{W} \mathbf{Z}^T \phi - \mathbf{M} \phi)] Q_t^\pi \right] \\
&\Leftrightarrow \sigma^2 \mathbb{E}_\tau \left[\sum_{t=1}^T Q_t^\pi \right] = \frac{1}{d} \mathbb{E}_\tau \left[\sum_{t=1}^T \frac{1}{\text{trace}(\phi \phi^T)} \mathbb{E}_{\mathbf{Z}^T \phi | \mathbf{a}} [(\mathbf{a} - \mathbf{M} \phi)^T (\mathbf{a} - \mathbf{M} \phi) - 2(\mathbf{a} - \mathbf{M} \phi)^T \mathbf{W} \mathbf{Z}^T \phi + \phi^T \mathbf{Z} \mathbf{W}^T \mathbf{W} \mathbf{Z}^T \phi] Q_t^\pi \right]. \tag{3.22}
\end{aligned}$$

The update equation for this parameter is given by

$$\begin{aligned}
\sigma^2 &= \frac{1}{d} \mathbb{E}_\tau \left[\sum_{t=1}^T \frac{1}{\text{trace}(\phi \phi^T)} \left((\mathbf{a} - \mathbf{M} \phi)^T (\mathbf{a} - \mathbf{M} \phi) - 2(\mathbf{a} - \mathbf{M} \phi)^T \mathbf{W} \mathbb{E}_{\mathbf{Z}^T \phi | \mathbf{a}} [\mathbf{Z}^T \phi] \right. \right. \\
& \quad \left. \left. + \text{trace}(\mathbb{E}_{\mathbf{Z}^T \phi | \mathbf{a}} [\mathbf{Z}^T \phi \phi^T \mathbf{Z}]) \mathbf{W}^T \mathbf{W} \right) Q_t^\pi \right] \left(\mathbb{E}_\tau \left[\sum_{t=1}^T Q_t^\pi \right] \right)^{-1}. \tag{3.23}
\end{aligned}$$

3.4.3 Complete Algorithm

With the results given above we can find the 'Policy Search with Probabilistic Principle Component Exploration in the Action Space' (PePPcEr) algorithm which can be found in Alg. 1. The initial parameters σ^2_0, \mathbf{W}_0 and \mathbf{M}_0 can either be chosen randomly or with a pre-analysis with PCA or PPCA. Furthermore, as an additional parameter, the algorithm needs the number of latent dimension n as input. This number depends on the given (robot) system and the task. A pre-analysis on a few samples may help to uncover this property of the hidden subspace. As already mentioned, a feature vector $\phi(\mathbf{s}_t, t) \in \mathbb{R}^p$ is needed. An example for this vector is given later in the experiment section. The result of this algorithm is the matrix \mathbf{M} which contains weights for the basis functions.

Input: Initialized parameters σ^2_0, \mathbf{W}_0 and \mathbf{M}_0 and the number n of the low intrinsic dimension. The function $\phi(\mathbf{s}_t, t)$ serves the basic functions as features.

repeat

Sampling:

for $h=1:H$ **do** # Sample the H rollouts

for $t=1:T$ **do**

$\mathbf{a}_t^h = \mathbf{W}_i \mathbf{Z}^T \phi + \mathbf{M}_i \phi + \mathbf{E} \phi$ with $\mathbf{Z} \sim \mathcal{N}_{p,n}(\mathbf{0}, \mathbf{I})$ and $\mathbf{E} \sim \mathcal{N}_{d,p}(\mathbf{0}, \sigma^2_i \mathbf{I})$

 Execute action \mathbf{a}_t^h

 Observe and store reward $r_t(\mathbf{s}_t, \mathbf{a}_t^h)$

Calculate weights:

$$Q^\pi(\mathbf{s}, \mathbf{a}, t) = \mathbb{E} \left[\sum_{\tilde{t}=t}^T r_{\tilde{t}}(\mathbf{s}_{\tilde{t}}, \mathbf{a}_{\tilde{t}}) \mid \mathbf{s}_t = \mathbf{s}, \mathbf{a}_t = \mathbf{a} \right]$$

Expectation:

$$\mathbf{C} = (\sigma^2_i \mathbf{I} + \mathbf{W}_i^T \mathbf{W}_i)^{-1}$$

foreach \mathbf{a}_t^h **do**

$$\mathbb{E}_{\mathbf{Z}^T \phi | \mathbf{a}_t^h} [\mathbf{Z}^T \phi] = \mathbf{C} \mathbf{W}_i^T (\mathbf{a}_t^h - \mathbf{M}_i \phi)$$

$$\mathbb{E}_{\mathbf{Z}^T \phi | \mathbf{a}_t^h} [\mathbf{Z}^T \phi (\mathbf{Z}^T \phi)^T] = \mathbf{C} \sigma^2_i \text{trace}(\phi \phi^T) + \mathbb{E}_{\mathbf{Z}^T \phi | \mathbf{a}_t^h} [\mathbf{Z}^T \phi] \cdot \mathbb{E}_{\mathbf{Z}^T \phi | \mathbf{a}_t^h} [\mathbf{Z}^T \phi]^T$$

Maximization:

$$\mathbf{M}_{i+1} = \mathbb{E}_\tau \left[\sum_{t=1}^T \text{trace}(\phi \phi^T)^{-1} \mathbf{a}_t^h \phi^T Q_t^\pi \right] \left(\mathbb{E}_\tau \left[\sum_{t=1}^T \text{trace}(\phi \phi^T)^{-1} \phi \phi^T Q_t^\pi \right] \right)^{-1}$$

$$\mathbf{W}_{i+1} = \mathbb{E}_\tau \left[\sum_{t=1}^T \text{trace}(\phi \phi^T)^{-1} (\mathbf{a}_t^h - \mathbf{M}_{i+1} \phi) \mathbb{E}_{\mathbf{Z}^T \phi | \mathbf{a}_t^h} [\mathbf{Z}^T \phi]^T Q_t^\pi \right] \cdot \left(\mathbb{E}_\tau \left[\sum_{t=1}^T \text{trace}(\phi \phi^T)^{-1} \mathbb{E}_{\mathbf{Z}^T \phi | \mathbf{a}_t^h} [\mathbf{Z}^T \phi (\mathbf{Z}^T \phi)^T] Q_t^\pi \right] \right)^{-1}$$

$$\sigma^2_{i+1} = d^{-1} \mathbb{E}_\tau \left[\sum_{t=1}^T \text{trace}(\phi \phi^T)^{-1} \left((\mathbf{a}_t^h - \mathbf{M}_{i+1} \phi)^T (\mathbf{a}_t^h - \mathbf{M}_{i+1} \phi) - 2 (\mathbf{a}_t^h - \mathbf{M}_{i+1} \phi)^T \mathbf{W}_{i+1} \mathbb{E}_{\mathbf{Z}^T \phi | \mathbf{a}_t^h} [\mathbf{Z}^T \phi] + \text{trace}(\mathbb{E}_{\mathbf{Z}^T \phi | \mathbf{a}_t^h} [\mathbf{Z}^T \phi \phi^T \mathbf{Z}] \mathbf{W}_{i+1}^T \mathbf{W}_{i+1}) \right) Q_t^\pi \right] \left(\mathbb{E}_\tau \left[\sum_{t=1}^T Q_t^\pi \right] \right)^{-1}$$

until $\mathbf{M}_i \approx \mathbf{M}_{i+1}$

Output: Policy \mathbf{M} for the feature vector ϕ

Algorithm 1: Policy Search with Probabilistic Principle Component Exploration in the Action Space (PePPcEr)

3.5 Diagonal Matrix Extension

Revisiting the stochastic model of PePPcEr given by

$$\mathbf{a} = (\mathbf{W} \mathbf{Z}^T + \mathbf{M} + \mathbf{E}) \phi = \mathbf{W} (\mathbf{Z}^T \phi) + \mathbf{M} \phi + \mathbf{E} \phi, \quad (3.24)$$

we can find a slightly different model if we choose the covariance matrix of \mathbf{E} as a diagonal matrix with several parameters for each dimension in the action space. Thus we have $\mathbf{E} \sim \mathcal{N}_{d,p}(\mathbf{0}, \mathbf{D} \otimes \mathbf{I})$ with $\mathbf{D} \in \mathbb{R}^{d,d}$ as such a diagonal covariance matrix with respect to the action space and the other distributions as before. For Eq. (3.7) we can now find that

$$\begin{aligned} \text{cov}(\mathbf{a}) &= \mathbb{E}[\mathbf{WZ}^T \boldsymbol{\phi} \boldsymbol{\phi}^T \mathbf{ZW}^T] + \mathbb{E}[\mathbf{E} \boldsymbol{\phi} \boldsymbol{\phi}^T \mathbf{E}^T] \\ &= \mathbf{W} \text{trace}(\boldsymbol{\phi} \boldsymbol{\phi}^T) \mathbf{I} \mathbf{W}^T + \text{trace}(\boldsymbol{\phi} \boldsymbol{\phi}^T) \mathbf{D} \\ &= \text{trace}(\boldsymbol{\phi} \boldsymbol{\phi}^T) (\mathbf{W}\mathbf{W}^T + \mathbf{D}) \end{aligned} \quad (3.25)$$

because of Thm. 2 and $\mathbf{E}^T \sim \mathcal{N}_{p,d}(\mathbf{0}, \mathbf{I} \otimes \mathbf{D})$. Therefore the property

$$\mathbf{p}(\mathbf{a}) = \mathcal{N}(\mathbf{M}\boldsymbol{\phi}, \text{trace}(\boldsymbol{\phi} \boldsymbol{\phi}^T) (\mathbf{W}\mathbf{W}^T + \mathbf{D})) \quad (3.26)$$

exists for the distribution of \mathbf{a} and we can find for $\ln \mathbf{p}(\mathbf{a}, \mathbf{Z})$ that

$$\begin{aligned} \frac{\partial \ln \mathbf{p}(\mathbf{a}, \mathbf{Z}^T \boldsymbol{\phi})}{\partial \boldsymbol{\theta}_{\text{new}}} &= \frac{\partial}{\partial \boldsymbol{\theta}_{\text{new}}} \left(-\frac{d}{2} \ln(2\pi \text{trace}(\boldsymbol{\phi} \boldsymbol{\phi}^T)) - \frac{1}{2} \ln \det(\mathbf{D}) \right. \\ &\quad \left. - \frac{1}{2 \text{trace}(\boldsymbol{\phi} \boldsymbol{\phi}^T)} (\mathbf{a} - \mathbf{WZ}^T \boldsymbol{\phi} - \mathbf{M}\boldsymbol{\phi})^T \mathbf{D}^{-1} (\mathbf{a} - \mathbf{WZ}^T \boldsymbol{\phi} - \mathbf{M}\boldsymbol{\phi}) \right). \end{aligned} \quad (3.27)$$

While the maximization term for \mathbf{M} remains the same it is important to ensure that the maximization for \mathbf{W} does not depend on \mathbf{D} and is still solvable. We can find that this is still true with

$$\begin{aligned} \frac{\partial \ln \mathbf{p}(\mathbf{a}, \mathbf{Z}^T \boldsymbol{\phi})}{\partial \mathbf{W}} &\stackrel{3.27}{=} \frac{\partial}{\partial \mathbf{W}} \left(-\frac{d}{2} \ln(2\pi \text{trace}(\boldsymbol{\phi} \boldsymbol{\phi}^T)) - \frac{1}{2} \ln \det(\mathbf{D}) + \right. \\ &\quad \left. - \frac{1}{2 \text{trace}(\boldsymbol{\phi} \boldsymbol{\phi}^T)} (\mathbf{a} - \mathbf{WZ}^T \boldsymbol{\phi} - \mathbf{M}\boldsymbol{\phi})^T \mathbf{D}^{-1} (\mathbf{a} - \mathbf{WZ}^T \boldsymbol{\phi} - \mathbf{M}\boldsymbol{\phi}) \right) \\ &= \frac{\partial}{\partial \mathbf{W}} \left(-\frac{1}{2\sigma^2 \text{trace}(\boldsymbol{\phi} \boldsymbol{\phi}^T)} (\mathbf{a} - \mathbf{WZ}^T \boldsymbol{\phi} - \mathbf{M}\boldsymbol{\phi})^T \mathbf{D}^{-1} (\mathbf{a} - \mathbf{WZ}^T \boldsymbol{\phi} - \mathbf{M}\boldsymbol{\phi}) \right) \\ &= \frac{\partial}{\partial \mathbf{W}} \left(-\frac{1}{2\sigma^2 \text{trace}(\boldsymbol{\phi} \boldsymbol{\phi}^T)} (-2\mathbf{a}^T \mathbf{D}^{-1} \mathbf{WZ}^T \boldsymbol{\phi} + \boldsymbol{\phi}^T \mathbf{ZW}^T \mathbf{D}^{-1} \mathbf{WZ}^T \boldsymbol{\phi} + 2\boldsymbol{\phi}^T \mathbf{M}^T \mathbf{D}^{-1} \mathbf{WZ}^T \boldsymbol{\phi}) \right) \\ &\stackrel{A.1}{=} -\frac{1}{2\sigma^2 \text{trace}(\boldsymbol{\phi} \boldsymbol{\phi}^T)} \left(-2\mathbf{D}^{-1} \mathbf{a} (\mathbf{Z}^T \boldsymbol{\phi})^T + \frac{\partial \boldsymbol{\phi}^T \mathbf{ZW}^T \mathbf{D}^{-1} \mathbf{WZ}^T \boldsymbol{\phi}}{\partial \mathbf{W}} + 2\mathbf{D}^{-1} \mathbf{M}\boldsymbol{\phi} (\mathbf{Z}^T \boldsymbol{\phi})^T \right) \\ &\stackrel{A.2}{=} -\frac{1}{\sigma^2 \text{trace}(\boldsymbol{\phi} \boldsymbol{\phi}^T)} \left(-\mathbf{D}^{-1} \mathbf{a} (\mathbf{Z}^T \boldsymbol{\phi})^T + \mathbf{D}^{-1} \mathbf{WZ}^T \boldsymbol{\phi} (\mathbf{Z}^T \boldsymbol{\phi})^T + \mathbf{D}^{-1} \mathbf{M}\boldsymbol{\phi} (\mathbf{Z}^T \boldsymbol{\phi})^T \right). \end{aligned} \quad (3.28)$$

because the matrix derivations ensure that the matrix \mathbf{D} can be found on the left side of each term and so Eq. 3.19 still holds.

Now we have to determine the maximization term of \mathbf{D} . First we can find for the differentiation of $\ln \mathbf{p}(\mathbf{a}, \mathbf{Z}^T \boldsymbol{\phi})$ w.r.t. \mathbf{D} that

$$\begin{aligned} \frac{\partial \ln \mathbf{p}(\mathbf{a}, \mathbf{Z}^T \boldsymbol{\phi})}{\partial \mathbf{D}} &\stackrel{3.27}{=} \frac{\partial}{\partial \mathbf{D}} \left(-\frac{1}{2} \ln \det(\mathbf{D}) - \frac{1}{2 \text{trace}(\boldsymbol{\phi} \boldsymbol{\phi}^T)} (\mathbf{a} - \mathbf{WZ}^T \boldsymbol{\phi} - \mathbf{M}\boldsymbol{\phi})^T \mathbf{D}^{-1} (\mathbf{a} - \mathbf{WZ}^T \boldsymbol{\phi} - \mathbf{M}\boldsymbol{\phi}) \right) \\ &\stackrel{A.3}{=} -\frac{1}{2} \mathbf{D}^{-1} - \frac{\partial}{\partial \mathbf{D}} \frac{1}{2 \text{trace}(\boldsymbol{\phi} \boldsymbol{\phi}^T)} (\mathbf{a} - \mathbf{WZ}^T \boldsymbol{\phi} - \mathbf{M}\boldsymbol{\phi})^T \mathbf{D}^{-1} (\mathbf{a} - \mathbf{WZ}^T \boldsymbol{\phi} - \mathbf{M}\boldsymbol{\phi}) \\ &\stackrel{A.2}{=} -\frac{1}{2} \mathbf{D}^{-1} + \frac{1}{2 \text{trace}(\boldsymbol{\phi} \boldsymbol{\phi}^T)} \mathbf{D}^{-1} \text{diag} \left\{ (\mathbf{a} - \mathbf{WZ}^T \boldsymbol{\phi} - \mathbf{M}\boldsymbol{\phi}) (\mathbf{a} - \mathbf{WZ}^T \boldsymbol{\phi} - \mathbf{M}\boldsymbol{\phi})^T \right\} \mathbf{D}^{-1} \end{aligned} \quad (3.29)$$

holds and if we plugging this result into the lower bound we get that

$$\begin{aligned} 0 &= \mathbb{E}_{\boldsymbol{\tau}} \left[\sum_{t=1}^T \frac{\partial}{\partial \mathbf{D}} \mathbb{E}_{\mathbf{Z}^T \boldsymbol{\phi} | \mathbf{a}} [\ln \mathbf{p}(\mathbf{a}, \mathbf{Z}^T \boldsymbol{\phi})] Q_t^\pi \right] \\ \Leftrightarrow 0 &= \mathbb{E}_{\boldsymbol{\tau}} \left[\sum_{t=1}^T -\frac{1}{2} \mathbf{D}^{-1} Q_t^\pi + \sum_{t=1}^T \mathbb{E}_{\mathbf{Z}^T \boldsymbol{\phi} | \mathbf{a}} \left[\frac{1}{2 \text{trace}(\boldsymbol{\phi} \boldsymbol{\phi}^T)} \mathbf{D}^{-1} \text{diag} \left\{ (\mathbf{a} - \mathbf{WZ}^T \boldsymbol{\phi} - \mathbf{M}\boldsymbol{\phi}) (\mathbf{a} - \mathbf{WZ}^T \boldsymbol{\phi} - \mathbf{M}\boldsymbol{\phi})^T \right\} \mathbf{D}^{-1} \right] Q_t^\pi \right] \\ \Leftrightarrow \mathbb{E}_{\boldsymbol{\tau}} \left[\sum_{t=1}^T \frac{1}{2} \mathbf{D}^{-1} Q_t^\pi \right] &= \mathbb{E}_{\boldsymbol{\tau}} \left[\sum_{t=1}^T \mathbb{E}_{\mathbf{Z}^T \boldsymbol{\phi} | \mathbf{a}} \left[\frac{1}{2 \text{trace}(\boldsymbol{\phi} \boldsymbol{\phi}^T)} \mathbf{D}^{-1} \text{diag} \left\{ (\mathbf{a} - \mathbf{WZ}^T \boldsymbol{\phi} - \mathbf{M}\boldsymbol{\phi}) (\mathbf{a} - \mathbf{WZ}^T \boldsymbol{\phi} - \mathbf{M}\boldsymbol{\phi})^T \right\} \mathbf{D}^{-1} \right] Q_t^\pi \right]. \end{aligned} \quad (3.30)$$

If we now multiply \mathbf{D} from the left and the right side we find the maximization with

$$\begin{aligned} \mathbf{D} &= \mathbb{E}_\tau \left[\sum_{t=1}^T \mathbb{E}_{\mathbf{Z}^T \phi | \mathbf{a}} \left[\frac{1}{\text{trace}(\phi \phi^T)} \text{diag} \left\{ (\mathbf{a} - \mathbf{W} \mathbf{Z}^T \phi - \mathbf{M} \phi) (\mathbf{a} - \mathbf{W} \mathbf{Z}^T \phi - \mathbf{M} \phi)^T \right\} \right] Q_t^\pi \right] \cdot \left(\mathbb{E}_\tau \left[\sum_{t=1}^T Q_t^\pi \right] \right)^{-1} \\ \mathbf{D} &= \mathbb{E}_\tau \left[\sum_{t=1}^T \frac{1}{\text{trace}(\phi \phi^T)} \text{diag} \left\{ (\mathbf{a} - \mathbf{M} \phi) (\mathbf{a} - \mathbf{M} \phi)^T - (\mathbf{a} - \mathbf{M} \phi) \mathbb{E}_{\mathbf{Z}^T \phi | \mathbf{a}} [\mathbf{Z}^T \phi]^T \mathbf{W} - \mathbf{W} \mathbb{E}_{\mathbf{Z}^T \phi | \mathbf{a}} [\mathbf{Z}^T \phi] (\mathbf{a} - \mathbf{M} \phi)^T + \right. \right. \\ &\quad \left. \left. \mathbf{W} \mathbb{E}_{\mathbf{Z}^T \phi | \mathbf{a}} [\mathbf{Z}^T \phi (\mathbf{Z}^T \phi)^T] \mathbf{W}^T \right\} Q_t^\pi \right] \cdot \left(\mathbb{E}_\tau \left[\sum_{t=1}^T Q_t^\pi \right] \right)^{-1}. \end{aligned} \quad (3.31)$$

The resulting algorithm can be found in Alg. 2, where we can find the new update for the diagonal covariance matrix and slightly different terms for the expectation step.

Input: Initialized parameters $\mathbf{D}_0, \mathbf{W}_0$ and \mathbf{M}_0 and the number n of the low intrinsic dimension. The function $\phi(\mathbf{s}_t, t)$ serves the basic functions as features.

repeat

Sampling:

for $h=1:H$ **do** # Sample the H rollouts

for $t=1:T$ **do**

$\mathbf{a}_t^h = \mathbf{W}_i \mathbf{Z}^T \phi + \mathbf{M}_i \phi + \mathbf{E} \phi$ with $\mathbf{Z} \sim \mathcal{N}_{p,n}(\mathbf{0}, \mathbf{I})$ and $\mathbf{E} \sim \mathcal{N}_{d,p}(\mathbf{0}, \mathbf{D}_i \otimes \mathbf{I})$

 Execute action \mathbf{a}_t^h

 Observe and store reward $r_t(\mathbf{s}_t, \mathbf{a}_t^h)$

Calculate weights:

$$Q^\pi(\mathbf{s}, \mathbf{a}, t) = \mathbb{E} \left[\sum_{\tilde{t}=t}^T r_{\tilde{t}}(\mathbf{s}_{\tilde{t}}, \mathbf{a}_{\tilde{t}}) \mid \mathbf{s}_t = \mathbf{s}, \mathbf{a}_t = \mathbf{a} \right]$$

Expectation:

$$\mathbf{C} = (\mathbf{I} + \mathbf{W}_i^T \mathbf{D}_i^{-1} \mathbf{W}_i)^{-1}$$

foreach \mathbf{a}_t^h **do**

$$\mathbb{E}_{\mathbf{Z}^T \phi | \mathbf{a}_t^h} [\mathbf{Z}^T \phi] = \mathbf{C} \mathbf{W}_i^T \mathbf{D}_i^{-1} (\mathbf{a}_t^h - \mathbf{M}_i \phi)$$

$$\mathbb{E}_{\mathbf{Z}^T \phi | \mathbf{a}_t^h} [\mathbf{Z}^T \phi (\mathbf{Z}^T \phi)^T] = \mathbf{C} \cdot \text{trace}(\phi \phi^T) + \mathbb{E}_{\mathbf{Z}^T \phi | \mathbf{a}_t^h} [\mathbf{Z}^T \phi] \cdot \mathbb{E}_{\mathbf{Z}^T \phi | \mathbf{a}_t^h} [\mathbf{Z}^T \phi]^T$$

Maximization:

$$\mathbf{M}_{i+1} = \mathbb{E}_\tau \left[\sum_{t=1}^T \text{trace}(\phi \phi^T)^{-1} \mathbf{a}_t^h \phi^T Q_t^\pi \right] \left(\mathbb{E}_\tau \left[\sum_{t=1}^T \text{trace}(\phi \phi^T)^{-1} \phi \phi^T Q_t^\pi \right] \right)^{-1}$$

$$\mathbf{W}_{i+1} = \mathbb{E}_\tau \left[\sum_{t=1}^T \text{trace}(\phi \phi^T)^{-1} (\mathbf{a}_t^h - \mathbf{M}_{i+1} \phi) \mathbb{E}_{\mathbf{Z}^T \phi | \mathbf{a}_t^h} [\mathbf{Z}^T \phi]^T Q_t^\pi \right] \cdot \left(\mathbb{E}_\tau \left[\sum_{t=1}^T \text{trace}(\phi \phi^T)^{-1} \mathbb{E}_{\mathbf{Z}^T \phi | \mathbf{a}_t^h} [\mathbf{Z}^T \phi (\mathbf{Z}^T \phi)^T] Q_t^\pi \right] \right)^{-1}$$

$$\mathbf{D}_{i+1} = \mathbb{E}_\tau \left[\sum_{t=1}^T \text{trace}(\phi \phi^T)^{-1} \text{diag} \left\{ (\mathbf{a} - \mathbf{M}_{i+1} \phi) (\mathbf{a} - \mathbf{M}_{i+1} \phi)^T - 2 \mathbf{W}_{i+1} \mathbb{E}_{\mathbf{Z}^T \phi | \mathbf{a}} [\mathbf{Z}^T \phi] (\mathbf{a} - \mathbf{M}_{i+1} \phi)^T + \right. \right. \\ \left. \left. \mathbf{W}_{i+1} \mathbb{E}_{\mathbf{Z}^T \phi | \mathbf{a}} [\mathbf{Z}^T \phi (\mathbf{Z}^T \phi)^T] \mathbf{W}_{i+1}^T \right\} Q_t^\pi \right] \cdot \left(\mathbb{E}_\tau \left[\sum_{t=1}^T Q_t^\pi \right] \right)^{-1}$$

until $\mathbf{M}_i \approx \mathbf{M}_{i+1}$

Output: Policy \mathbf{M} for the feature vector ϕ

Algorithm 2: Policy Search with Probabilistic Principle Component Exploration in the Action Space (PePP_cEr) with Diagonal Matrix Extension.

4 Experiments

We conducted a set of experiments in order to compare PePPeEr to other existing algorithms for policy search. In this chapter we will present the setups under which the experiments were carried out. Both results on simulation and with a physical robot will be reported. All experiments were done with the Diagonal Matrix Extension for the PePPeEr Algorithm.

4.1 Learning Inverse Kinematics

For the comparison of PePPeEr with other algorithms we set up an experiment with a simulated robot arm. The arm has 26 hinge-joints and 27 segments, an example of such an arm for 8 joints can be found in Fig. 4.1. The goal is to track a trajectory using the end-effector of the robot arm, which we choose to be an ellipse. This ellipse, which is visible as a red circle in Fig. 4.1, follows the equation

$$f(t) = \begin{pmatrix} \cos(t \cdot 0.15) + \text{DOF} - 2.5 \\ \sin(t \cdot 0.15) \cdot 1.5 + 2 \\ 0 \end{pmatrix}, \quad (4.1)$$

where $1 \leq t \leq 50$ is the current time step and DOF the degrees-of-freedom of the robot arm. This forms an ellipse at the upper end of the robot arm. In the start configuration all joints were set to zero degree and the position of the end-effector was $(\text{DOF}, 0, 0)^T$ since each segment is one unit long.

For each of the 50 time steps, we have to learn parameters to bring the end-effector to touch the current point of this trajectory. In our case, the parameters are weights for the summation of Gaussian features, which depends on time. The activations of these Gaussians can be found in Fig. 4.2, where we can see that each Gaussian has a variance of three and the distance between two means is three time steps. To learn the inverse kinematics, we set the reward function to

$$r_t(\mathbf{s}_t, \mathbf{a}_t) = \exp(-D), \quad (4.2)$$

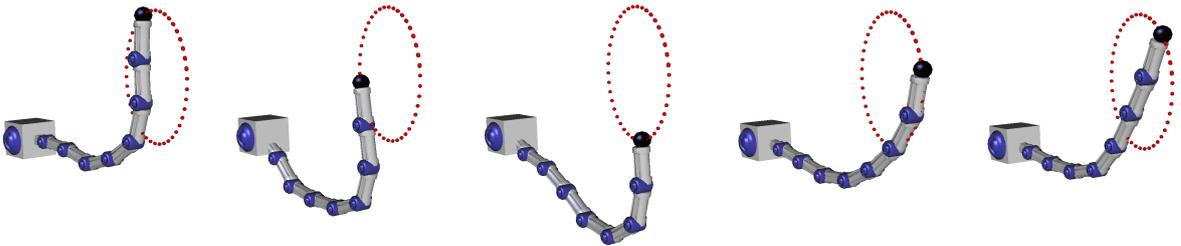


Figure 4.1.: A robot arm with 8 DOF in a 2D Tracking Task. These movements were learned with PePPeEr within 1000 iterations and have an error near zero. The red circle is the desired trajectory of the end-effector.

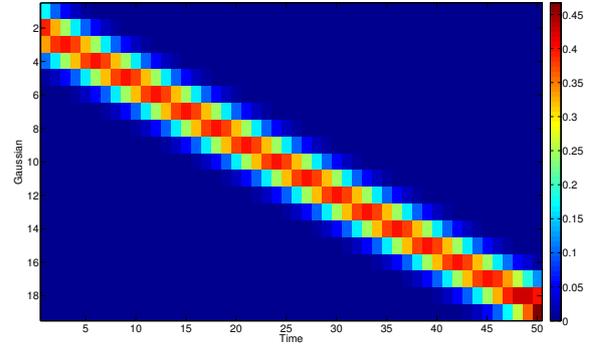


Figure 4.2.: The rows represent the Gaussians and the columns the specific time step. The values of the Gaussians are normalized for each time step.

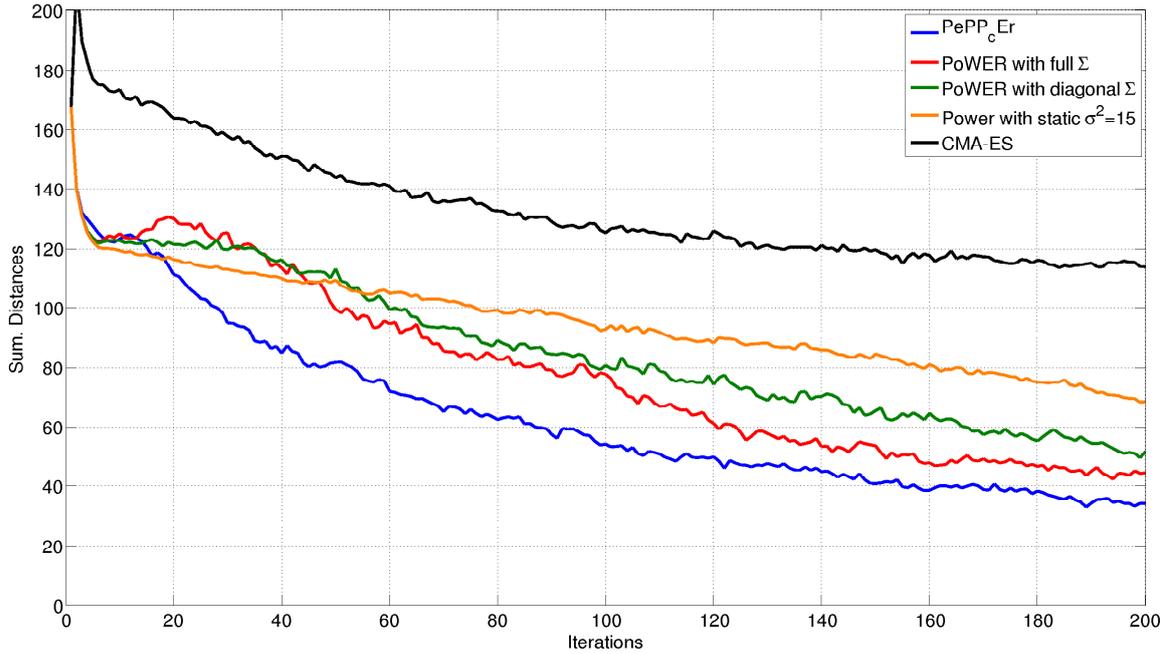


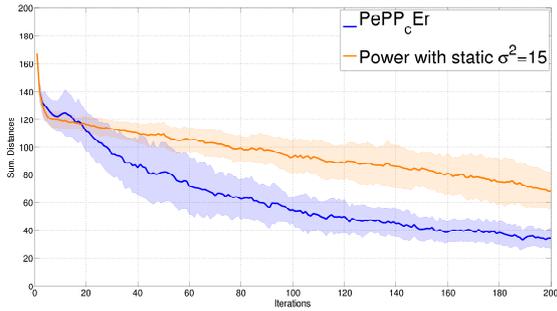
Figure 4.3.: Comparison between PePP_cEr , PoWER and CMA-ES on the inverse kinematic task with a 26-linked robot arm. In each iteration we generated 50 parameterized joint configurations and used the 25 best samples. Every Algorithm was performed with 30 trials and plotted with his mean.

where D is the distance between the end-effector and the target at time step t and the executed action \mathbf{a}_t . The state of the system is not relevant in this function.

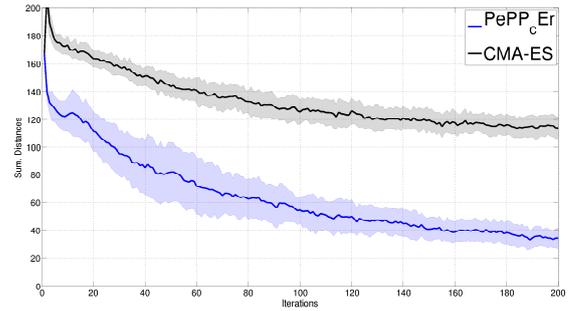
While our task space is only two dimensional, the robot arm meanwhile has 26 DOF which implies that we have a high redundancy in our system. PePP_cEr tries to uncover the latent structure of this 26-dimensional space by determining the mean and the principal axes of this subspace, where each axis represents the co-articulation of different links. Algorithms like PoWER are not able to identify these dependencies between the dimensions of the actions. The Covariance Matrix Adaption Evolution Strategy (CMA-ES) Algorithm can find a full covariance matrix for all parameters, but it will be slower than an Algorithm that has a strong exploration in the principal components concerning tasks where only a few dimensions are relevant. We ran the explained setup with different policy search algorithms, the results can be found in Fig. 4.3. We can see that after 19 iterations PePP_cEr significantly outperforms the different PoWER implementations and the CMA-ES Algorithm. For a better comparison between PePP_cEr and the different algorithms we show in Fig. 4.4 only PePP_cEr and one other algorithm per graph. As mentioned above we used 19 time-dependent Gaussians over 50 time steps (see Fig. 4.2 for an intuition) so we have to estimate 19 weights for one action dimension to approximate the desired trajectory. For 26 dimensions this results in a total number of 494 parameters which have to be estimated for this task. The number of assumed latent dimensions for the PePP_cEr Algorithm was set to $n = 11$ in this evaluation. Furthermore we used the Diagonal Matrix Extension. The algorithm runs 200 iterations and in each iteration 50 samples were created. From this set the 25 samples with highest rewards were chosen for updating the specific parameters of the algorithms. This was performed to have a better distinction between good and bad samples.

4.2 Learning to Stand on One Leg

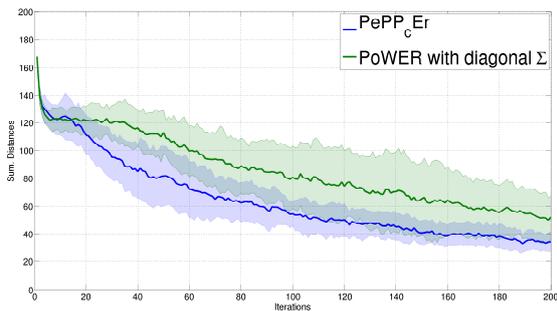
To show that the PePP_cEr Algorithm can be used in real-world tasks, we performed a learning task on the NAO robot. Our reward function measured the height of the left leg and the head of the robot, in order to learn to stand on one leg. This task requires the co-articulation of different body parts and is often used in biomechanical studies on synergies and low-dimensional control in humans, such as in [22]. Learning was performed in a physics-based simulator instead of the real robot, in order to prevent damages of the robot and to automate the process. Because normally sampling on real robots is



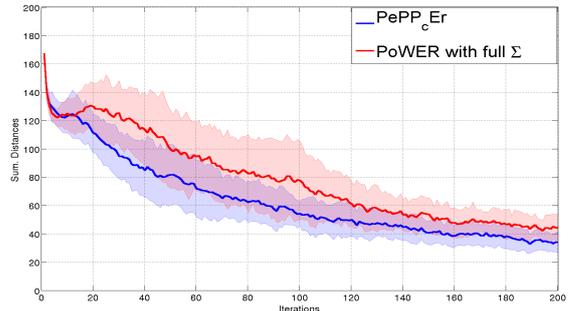
(a) Comparison between PePP_cEr and PoWER with a static diagonal covariance matrix $I \cdot 15$.



(b) Comparison between PePP_cEr and CMA-ES, where the parameters of CMA-ES are automatically chosen.



(c) Comparison between PePP_cEr and PoWER with a diagonal covariance matrix, which is adapted in each iteration.



(d) Comparison between PePP_cEr and the PoWER with an adaptation of a full covariance matrix in each iteration.

Figure 4.4.: Different comparisons between PePP_cEr and one algorithm from Fig. 4.3, each algorithm result plotted with its mean and variance.

cost and time intensive, the goal was to use only a few samples in each iteration to learn a stable and good policy. Thus, it is still possible to do learning directly on the robot, if no simulation is available.

In each iteration we executed 20 movements and used the 5 best executions for learning. Each movement consists of 5 time steps and was created by two time-dependent Gaussian features, so we have 52 parameters to estimate due to the 26 DOF of the NAO robot. As can be seen in Fig. 4.6 only a small number of samples is needed to learn a stable and smooth solution. Policy 1 was learned within 5 iterations, so 100 executions on the simulated robot in total.

In Policy 2 we used 600 samples in 30 iterations, which results in a movement where the robot can lift the leg even higher. Because of the stochastic nature of the PePP_cEr Algorithm, we can find very broad variations of leg lifting movements, two extreme examples can be found in Fig. 4.5. The shown movements seems extreme, but they are still stable and smooth, so that we can execute them on the real robot. Fig. 4.7 shows the evolution of a policy during the policy search process.

Initially the leg is not lifted. However, after several iterations the algorithm identifies a stable lifting-up movement. If we compare the results of the first iterations for different policies in Fig. 4.8 we can see that the found solutions have a strong dependency on the result of the first iteration. This behaviour can be traced back to the fact that the EM Algorithm is a local search algorithm. The results of Policy 4 are particularly insightful since in the solution, which is found in the first iteration of this trial, lifts the right leg instead of the left one. Even if this was the only

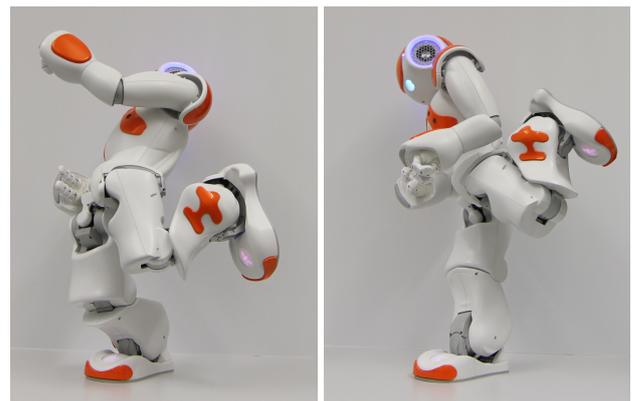


Figure 4.5.: The final poses of two extreme solutions learned by the PePP_cEr Algorithm, which are still stable.

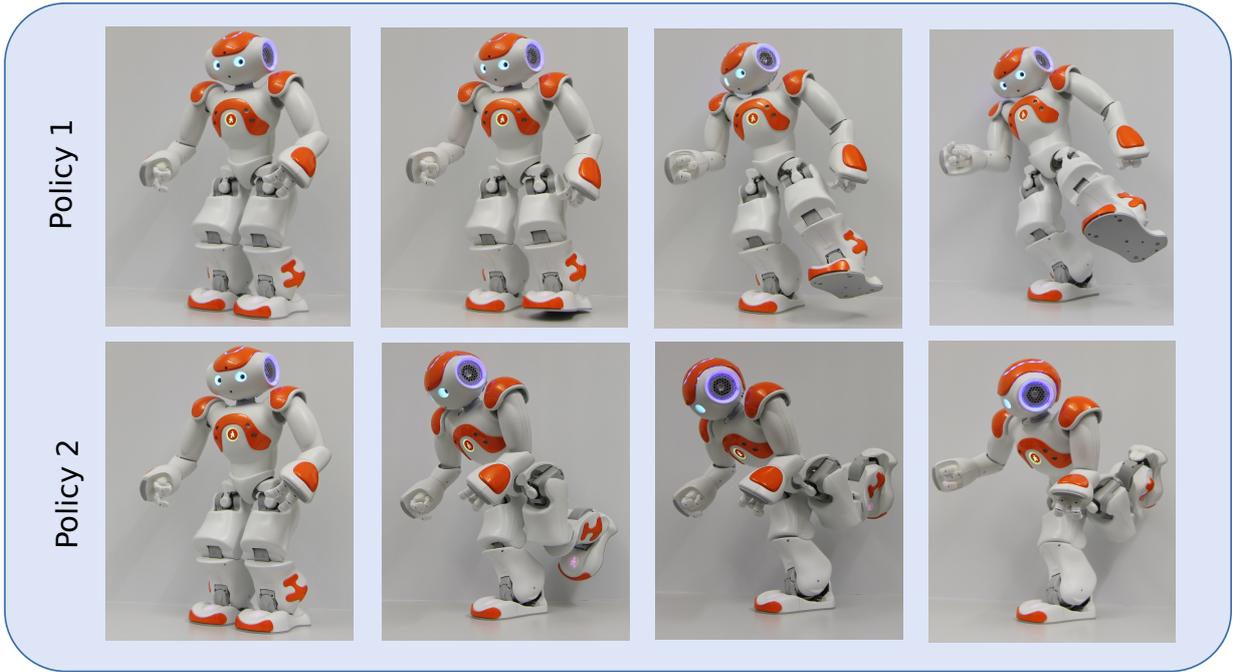


Figure 4.6.: Two different policies for standing on one leg learned using latent space policy search. The PePPeR algorithm needed only 100 executions for the first policy.

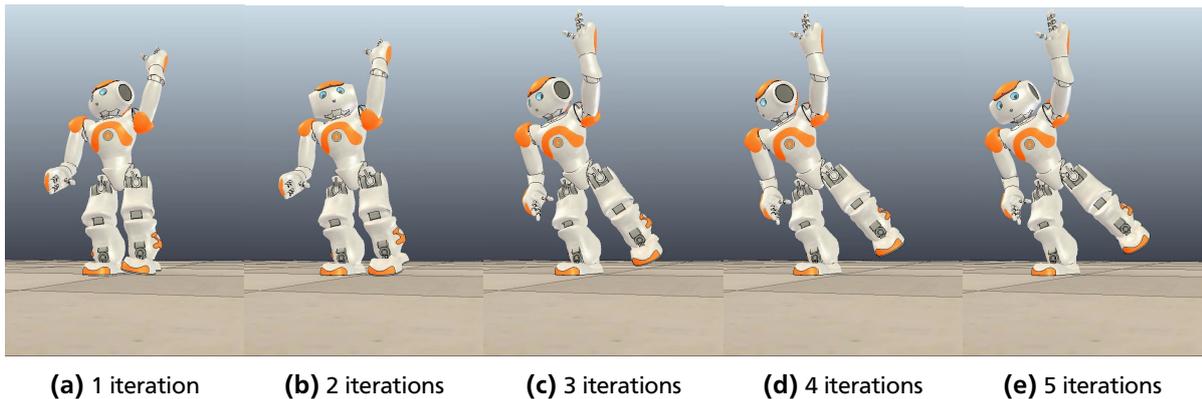


Figure 4.7.: The results of the first five iterations of Policy 5 learned by the PePPeR Algorithm

stable position found in the first iteration, the PePPeR Algorithm is able to find a smooth and stable movement for lifting the left leg within five iterations, therefore performing 100 executions.

4.3 Experimental Intuition

To give an intuition for the PePPeR Algorithm, we performed and visualized a trial on an optimization task. The test function was the sixth Bukin function [23], given by

$$f(x, y) = 100\sqrt{|y - 0.01x^2|} + 0.01|x + 10|, \min f(x, y) = f(-10, 1) = 0 \quad (4.3)$$

as a minimization problem. This function has one global minima, but many local minima that lie all in a ridge, which can be seen in Fig. 4.9, due to this the function models global optimization problem.

In the given plots in Fig. 4.10 we can see the solution and principal component for each iteration of the PePPeR Algorithm. For the sake of visibility, we show the principal component ten times larger. We can see that after a few iterations the mean point reaches the ridge and the principal component lies in the direction of the ridge as expected, because the five best samples lie along the ridge. Especially in the last 10 iterations, before the global minima is reached, the principal component leads to a sampling along the ridge to reach the global minima.

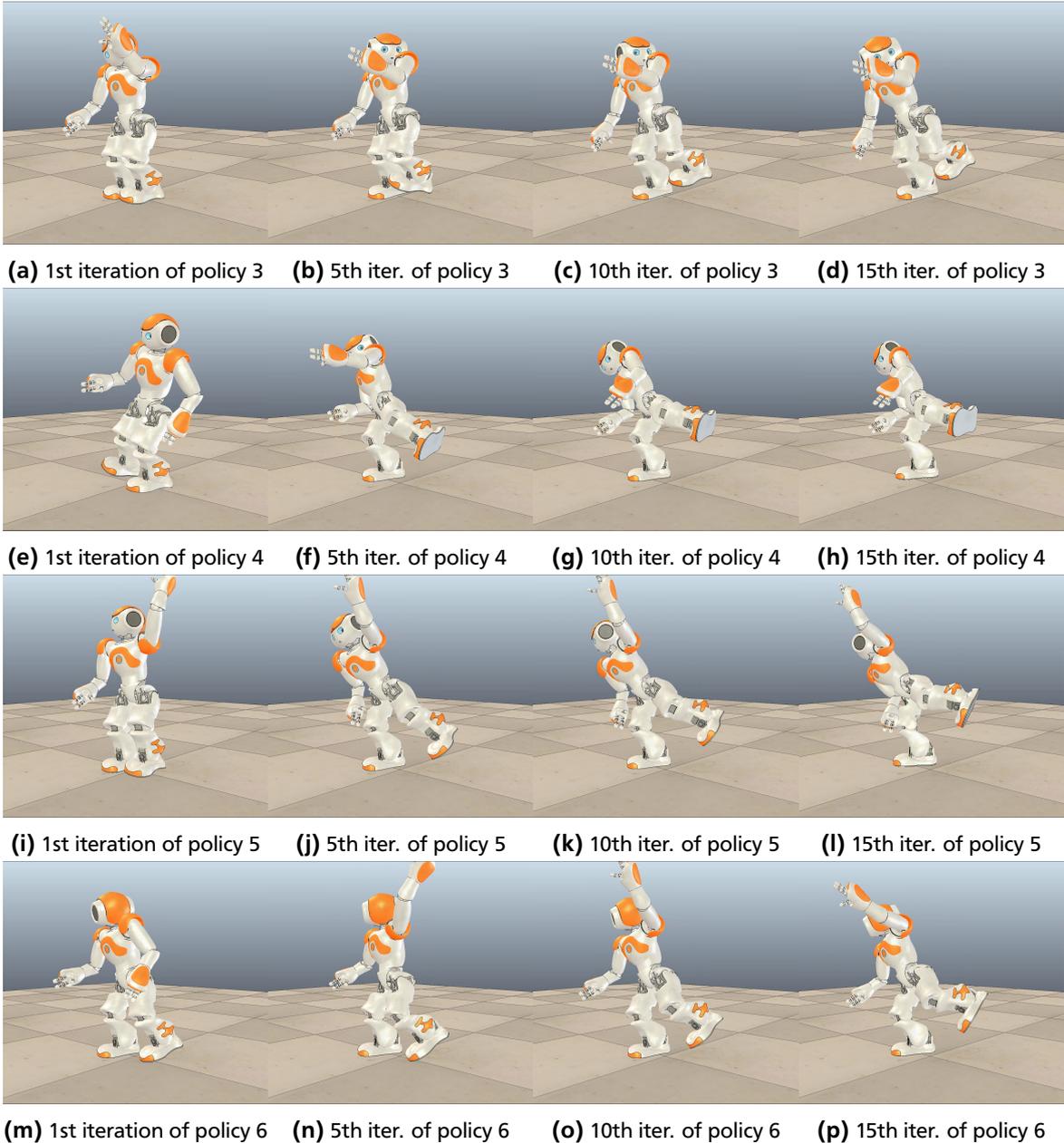


Figure 4.8.: The resulting poses after 1, 5, 10 and 15 iterations with the PePPeR Algorithm. Each row is one independent trial.

However, we can see that in the first two iterations it may be beneficial to use, in this case, both principal components for the first few iterations. When the ridge is reached, it is sufficient to use only one principal component for the exploration because of the almost linear structure of the ridge. While the adding and removing of columns to the transformation matrix W is already possible between two iterations, we need here a reasonable decision criterion for the number of dimensions. This possibility of using different numbers of principal components for various stages during an optimization process may be a crucial feature for a faster convergence and needs further investigations in the future.

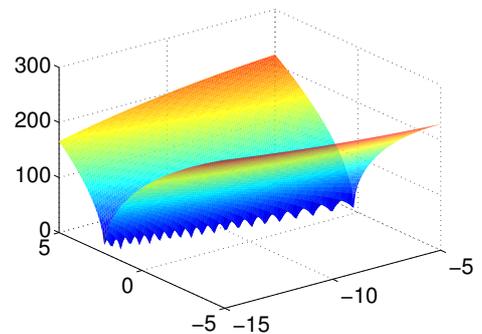


Figure 4.9.: The sixth Bukin function.

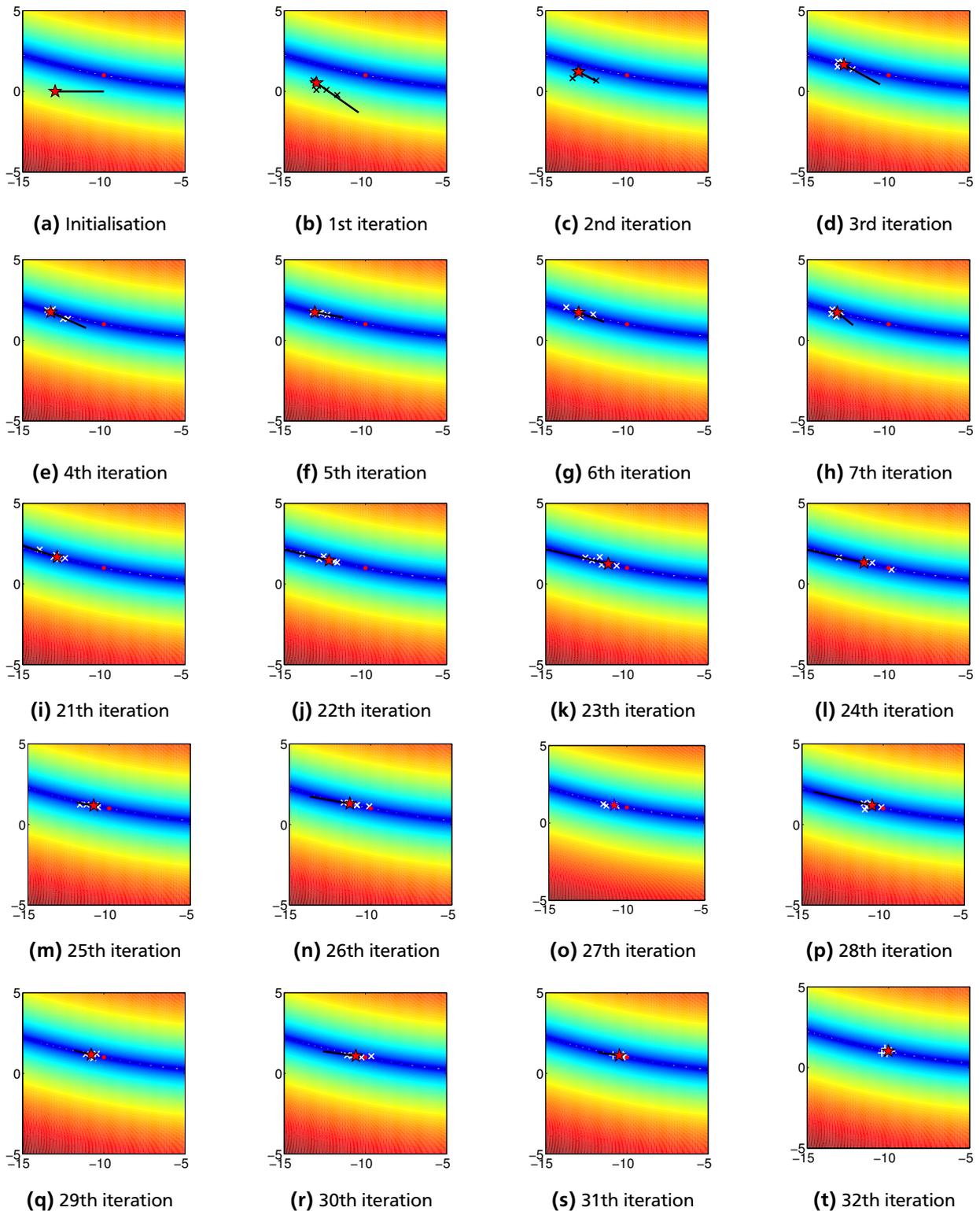


Figure 4.10.: The PePPER Algorithm performed on the sixth Bukin function. The red star is the current mean in the given iteration and the black line the principal component. Both are calculated from five samples, given by the white crosses. The global minima of the Bukin function is given by the red dot.

5 Conclusion and Future Work

In this thesis we presented a novel policy search algorithm for exploiting latent spaces in high dimensional action spaces. We combined the techniques of policy search and dimensionality reduction in an algorithm instead of using the dimensionality reduction as a pre-processing step.

The PePPeEr Algorithm was introduced with an isotropic undirected exploration and a diagonal matrix extension, where we can determine a parameter for this exploration for each dimension. We showed with a evaluation on a simulated and a real-world task that the property of the PePPeEr Algorithm, to uncover the hidden structure of a latent subspace, is able to speed up the convergence to an optimal solution in the given robot tasks. This indicates that our approach is suitable to do policy search for robots where we have a (high) redundancy in the actions and can only use a small set of samples due to costs and time.

However, so far the PePPeEr algorithm is only able to uncover the latent space in the action space. Although the algorithm is still better as the PoWER Algorithm with a full covariance matrix per action dimension, we want to investigate the possibility to uncover latent spaces in the feature space too. This may not be reasonable for time-dependent Gaussian features, but if we have to handle Dynamic Motor Primitives (i.e. see [4]) or feedback features it is possible that with a growing number of features like this we can find such a low dimensional space. This leads to the possibility to use a larger number of features, which can be a crucial feature for an online learning or improving task.

Furthermore the problem to guess the number of dimensions of the latent space right for a maximal efficiency of the PePPeEr Algorithm still exists. While it is possible to do this with a pre-analysis, we want to avoid this and instead guess the number of needed dimensions while learning, because the algorithm allows to change the dimensionality of the latent space while running.

A problem of the PoWER framework is to choose a right parameter for the transformed reward functions if we use the exponential function. So we have to adjust this reward function $\exp(-\lambda \cdot r_t)$ with a parameter $\lambda > 0$ to get a better distinction between good and bad samples. In contrast, the REPS [24] framework offers the possibility to adjust this parameter dynamically in each iteration. Because the PePPeEr Algorithm can be used in any policy search framework where we have a stochastic policy $\pi(\mathbf{a}_t, \mathbf{s}_t)$, we want to implement PePPeEr in the REPS Algorithm in the future.

Bibliography

- [1] 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 *Proceedings of the International Symposium on Experimental Robotics*, 2004, pp. 363–372.
- [2] K. Muelling, J. Kober, O. Kroemer, and J. Peters, "Learning to select and generalize striking movements in robot table tennis," *International Journal of Robotics Research*, no. 3, pp. 263–279, 2013.
- [3] J. Zico Kolter and A. Y. Ng, "The stanford littledog: A learning and rapid replanning approach to quadruped locomotion," *Int. J. Rob. Res.*, vol. 30, no. 2, pp. 150–174, Feb. 2011.
- [4] S. Schaal, J. Peters, J. Nakanishi, and A. Ijspeert, "Learning movement primitives," in *Robotics Research*. Springer, 2005, pp. 561–572.
- [5] R. E. Bellman, *Adaptive Control Processes: A guided tour*. Princeton University Press, 1961, ch. 5.16.
- [6] K. Pearson, "Liii. on lines and planes of closest fit to systems of points in space," *The London, Edinburgh, and Dublin Philosophical Magazine and Journal of Science*, vol. 2, no. 11, pp. 559–572, 1901.
- [7] M. E. Tipping and C. M. Bishop, "Probabilistic principal component analysis," *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, vol. 61, no. 3, pp. 611–622, 1999.
- [8] M. P. Deisenroth, R. Calandra, A. Seyfarth, and J. Peters, "Toward fast policy search for learning legged locomotion," in *IEEE/RSJ International Conference on Intelligent Systems and Robots*, 2012.
- [9] J. Z. Kolter and A. Y. Ng, "Learning omnidirectional path following using dimensionality reduction," in *Proceedings of Robotics: Science and Systems*, 2007.
- [10] S. Bitzer, M. Howard, and S. Vijayakumar, "Using dimensionality reduction to exploit constraints in reinforcement learning," in *Intelligent Robots and Systems (IROS), 2010 IEEE/RSJ International Conference on*, Oct 2010, pp. 3219–3225.
- [11] C. L. Nehaniv and K. Dautenhahn, "Imitation in animals and artifacts," K. Dautenhahn and C. L. Nehaniv, Eds. Cambridge, MA, USA: MIT Press, 2002, ch. The Correspondence Problem, pp. 41–61.
- [12] K. S. Luck, G. Neumann, E. Berger, J. Peters, and H. B. Amor, "Latent space policy search for robotics," 2014, *in review*.
- [13] A. K. Gupta and D. K. Nagar, *Matrix variate distributions*. CRC Press, 2000, vol. 104.
- [14] C. M. Bishop and N. M. Nasrabadi, *Pattern recognition and machine learning*. springer New York, 2006, vol. 1.
- [15] A. Hannachi, I. T. Jolliffe, and D. B. Stephenson, "Empirical orthogonal functions and related techniques in atmospheric science: A review," *International Journal of Climatology*, vol. 27, no. 9, pp. 1119–1152, 2007. [Online]. Available: <http://dx.doi.org/10.1002/joc.1499>
- [16] M. Khezri and M. Jahed, "An exploratory study to design a novel hand movement identification system," *Computers in Biology and Medicine*, vol. 39, no. 5, pp. 433 – 442, 2009. [Online]. Available: <http://www.sciencedirect.com/science/article/pii/S0010482509000420>
- [17] B. Moghaddam and A. Pentland, "Probabilistic visual learning for object representation," *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, vol. 19, no. 7, pp. 696–710, Jul 1997.

-
- [18] J. Lin, Y. Wu, and T. Huang, "Articulate hand motion capturing based on a monte carlo nelder-mead simplex tracker," in *Pattern Recognition, 2004. ICPR 2004. Proceedings of the 17th International Conference on*, vol. 4, 2004, pp. 975–978 Vol.4.
- [19] I. Jolliffe, *Principal component analysis*. Wiley Online Library, 2005.
- [20] J. Kober and J. Peters, "Policy search for motor primitives in robotics," *Machine Learning*, vol. 84, no. 1-2, pp. 171–203, 2011.
- [21] J. Peters and S. Schaal, "Reinforcement learning of motor skills with policy gradients," *Neural networks*, vol. 21, no. 4, pp. 682–697, 2008.
- [22] G. Torres-Oviedo and L. H. Ting, "Subject-specific muscle synergies in human balance control are consistent across different biomechanical contexts," *Journal of neurophysiology*, vol. 103, no. 6, pp. 3084–3098, 2010.
- [23] A. Bukin, *New Minimization Strategy For Non-Smooth Functions*. Budker Institute of Nuclear Physics preprint BUDKER-INP-1997-79, Novosibirsk, 1997.
- [24] C. Daniel, G. Neumann, and J. Peters, "Hierarchical relative entropy policy search," in *Proceedings of the International Conference on Artificial Intelligence and Statistics (AISTATS 2012)*, 2012.



A Appendix

A.1 Derivatives of Scalar Forms

The following derivatives with respect to matrices are given without any proof.

$$\frac{\partial \mathbf{a}^T \mathbf{X} \mathbf{b}}{\partial \mathbf{X}} = \mathbf{a} \mathbf{b}^T \quad (\text{A.1})$$

$$\frac{\partial \mathbf{b}^T \mathbf{X}^T \mathbf{D} \mathbf{X} \mathbf{c}}{\partial \mathbf{X}} = \mathbf{D}^T \mathbf{X} \mathbf{b} \mathbf{c}^T + \mathbf{D} \mathbf{X} \mathbf{c} \mathbf{b}^T \quad (\text{A.2})$$

If \mathbf{X} is a diagonal matrix then we have

$$\frac{\partial \ln \det(\mathbf{X})}{\partial \mathbf{X}} = \mathbf{X}^{-1} \quad (\text{A.3})$$

A.2 Marginalization Rule

This result is taken from [14, p. 93]:

If we have \mathbf{x} given as a marginal normal distribution and \mathbf{y} given \mathbf{x} as a conditional normal distribution defined by

$$p(\mathbf{x}) = \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Lambda}^{-1}) \quad (\text{A.4})$$

$$p(\mathbf{y}|\mathbf{x}) = \mathcal{N}(\mathbf{A}\mathbf{x} + \mathbf{b}, \mathbf{L}^{-1}) \quad (\text{A.5})$$

we are able to determine the normal distribution \mathbf{x} given \mathbf{y} with

$$p(\mathbf{x}|\mathbf{y}) = \mathcal{N}(\boldsymbol{\Sigma}(\mathbf{A}^T \mathbf{L}(\mathbf{y} - \mathbf{b}) + \boldsymbol{\Lambda} \boldsymbol{\mu}), \boldsymbol{\Sigma}) \quad (\text{A.6})$$

where $\boldsymbol{\Sigma}$ is defined by $\boldsymbol{\Sigma} = (\boldsymbol{\Lambda} + \mathbf{A}^T \mathbf{L} \mathbf{A})^{-1}$.

