

## Stochastic Search In Changing Situations

Abbas Abdolmaleki,<sup>1,2,3,5</sup> David Simões,<sup>1</sup> Nuno Lau,<sup>1</sup>  
Luis Paulo Reis,<sup>2,3</sup> Bob Price,<sup>5</sup> Gerhard Neumann<sup>4</sup>

1: IEETA, DETI, University of Aveiro, Aveiro, Portugal  
2: DSI, University of Minho, Braga, Portugal  
3: LIACC, University of Porto, Porto, Portugal  
4: CLAS, TU Darmstadt, Darmstadt, Germany  
5: PARC, Xerox, Palo Alto, USA

### Abstract

Stochastic search algorithms are black-box optimizer of an objective function. They have recently gained a lot of attention in operations research, machine learning and policy search of robot motor skills due to their ease of use and their generality. However, when the task or objective function slightly changes, many stochastic search algorithms require complete re-learning in order to adapt the solution to the new objective function or the new context. As such, we consider the contextual stochastic search paradigm. Here, we want to find good parameter vectors for multiple related tasks, where each task is described by a continuous context vector. Hence, the objective function might change slightly for each parameter vector evaluation. In this paper, we investigate a contextual stochastic search algorithm known as Contextual Relative Entropy Policy Search (CREPS), an information-theoretic algorithm that can learn from multiple tasks simultaneously. We show the application of CREPS for simulated robotic tasks.

### Introduction

Stochastic search algorithms are gradient-free black-box optimizers of some performance function dependent on a high-dimensional parameter vector. They directly evaluate the execution of a parameter vector by using the return of an episode. Stochastic search algorithms (Hansen, Muller, and Koumoutsakos 2003; Sun et al. 2009; Stulp and Sigaud 2012; Rückstieß, Felder, and Schmidhuber 2008) typically maintain a search distribution over the parameters that we want to optimise, which is used to create samples of the parameter vector. Subsequently, the performance of the sampled parameters is evaluated. Using the samples and their evaluations, a new search distribution is computed by computing gradient based updates (Sun et al. 2009; Rückstieß, Felder, and Schmidhuber 2008), evolutionary strategies (Hansen, Muller, and Koumoutsakos 2003), the cross-entropy method (Mannor, Rubinstein, and Gat 2003), path integrals (Stulp and Sigaud 2012; Theodorou, Buchli, and Schaal 2010), or information-theoretic policy updates (Kupcsik et al. 2013; Abdolmaleki et al. 2015a). However, many of the previously mentioned algorithms cannot be applied to multi-task learning. In other words, if the task setup or objective function changes slightly, re-learning is needed

to adapt the solution to the new situation or the new context. For example, consider optimising the parameters of a humanoid robot controller to kick a ball. Once the characteristics of the ball, such as weight or material, or objective function, such as desired kick distance, change, re-learning is needed. One could independently optimize for several target contexts in order to generalize a task, for example optimizing to kick the ball for different distances(context). Subsequently, when a new unseen context is presented, the optimized contexts can be generalized through regression methods (Niehaus, Röfer, and Laue 2007; Wang, Fleet, and Hertzmann 2009). However now optimizing for different contexts and then generalizing between the optimized parameters for different unseen contexts are two independent processes. Therefore, even though such approaches have been used successfully, they are time consuming as well as inefficient in terms of the number of needed training samples. In other words, we cannot reuse data-points obtained from optimizing a task with context  $s$  to improve and accelerate the optimization of a task with context  $s'$ . As such, it is desirable to learn the selection of the parameters for multiple tasks at once without restarting the learning process once we see a new task. This problem setup is also known as contextual policy search (Kupcsik et al. 2013; Kober, Oztop, and Peters 2010). Recently, such multi-task learning capability was established for information-theoretic policy search algorithms (Peters, Mülling, and Altun 2010), such as the episodic Contextual Relative Entropy Policy Search (CREPS) algorithm (Daniel, Neumann, and Peters 2012; Kupcsik et al. 2013). In (Abdolmaleki et al. 2015c), CREPS was successfully used to optimize a walking controller for different speeds. Despite its advantages, CREPS has a major set-back that does not allow it to perform favourably. Like many other stochastic search algorithms, CREPS maintains a Gaussian search distribution and it updates the mean and covariance matrix of its search distribution iteratively. However due to the covariance matrix update rule of CREPS, we will show that, search distribution might collapse prematurely to a point-estimate before finding a good solution, resulting in a premature convergence which is highly undesirable. Although, this multi-task learning capability is not found in other stochastic search algorithms (Hansen, Muller, and Koumoutsakos 2003; Sun et al. 2009), such as CMA-ES and NES, or commonly used policy search methods

(Stulp and Sigaud 2012; Kober and Peters 2010), they typically don't suffer from premature convergence. Therefore, to solve premature convergence problem of CREPS we use the rank- $\mu$  covariance matrix update rule of CMA-ES (Hansen, Muller, and Koumoutsakos 2003) along with CREPS. Now we combine the old coavarinace matrix with the estimated covariance matrix from samples. We call the resulting algorithm Contextual Relative Entropy Policy Search with Covariance Matrix Adaptation (CREPS-CMA). We evaluate CREPS-CMA on standard functions and simulated robotic tasks including a complex humanoid robot kicking task. We will show that CREPS-CMA works favourably.

## Problem Statement

Given a context vector  $s$  with  $m$  dimensions that defines a task, we want to find a function  $\mathbf{m}(s) : \mathbb{R}^m \rightarrow \mathbb{R}^n$  that outputs a parameter vector  $\theta$  with  $n$  dimensions such that it maximizes an objective function  $R(\theta, s) : \{\mathbb{R}^n, \mathbb{R}^m\} \rightarrow \mathbb{R}$ . The only accessible information on  $R(\theta, s)$  are evaluations  $\{R^{[k]}\}_{k=1\dots N}$  of samples  $\{s^{[k]}, \theta^{[k]}\}_{k=1\dots N}$ , where  $k$  is the index of the sample and  $N$  is number of samples. We maintain a search distribution  $\pi(\theta|s)$  over the parameter space  $\theta$  of the objective function  $R(\theta, s)$ . The search distribution  $\pi(\theta|s)$  is modeled as linear Gaussian policy, i.e.,

$$\pi(\theta|s) = \mathcal{N}\left(\theta|\mathbf{m}_\pi(s) = \mathbf{A}_\pi^T \varphi(s), \Sigma_\pi\right),$$

where  $\varphi(s)$  is an arbitrary feature function of context  $s$ .

In each iteration, a new coefficient matrix  $\mathbf{A}_\pi$  and a new covariance matrix  $\Sigma_\pi$  is obtained. Typically,  $\varphi(s) = [1 \ s]$ , which results in a linear generalization over contexts. However, one could use non-linear feature functions such as radial basis functions (RBF) (Broomhead and Lowe 1988) for non-linear generalization over contexts. RBF features have been shown to enable algorithms to learn non-linear policies which greatly outperform their linear counterparts on non-linear tasks, such as humanoid walking or humanoid kicking (Abdolmaleki et al. 2015c).

In each iteration, given context samples<sup>1</sup>  $s^{[k]}$ , the current search distribution  $q(\theta|s)$  is used to create samples  $\theta^{[k]}$  of the parameter vector  $\theta$ . Subsequently, the evaluation  $R^{[k]}$  of  $\{s^{[k]}, \theta^{[k]}\}$  is obtained by querying the objective function  $R(\theta, s)$ . In many algorithms, the samples  $\{s^{[k]}, \theta^{[k]}, R^{[k]}\}_{k=1\dots N}$  are used to compute a weight  $d^{[k]}$  for each sample  $k$ . Subsequently, using  $\{s^{[k]}, \theta^{[k]}, d^{[k]}\}_{k=1\dots N}$ , a new Gaussian search distribution  $\pi(\theta|s)$  is estimated. This process will run iteratively until the algorithm converges to a solution. Algorithm 1 shows a compact representation of contextual stochastic search methods.

## Contextual Stochastic Search

The Relative Entropy Policy Search (REPS) algorithm was originally proposed as a step-based policy search algorithm

<sup>1</sup>Please note that the way we sample contexts  $s^{[k]}$  depends on the task. Throughout this paper we use a uniform distribution to sample contexts  $s^{[k]}$ .

---

### Algorithm 1 Contextual stochastic search algorithm

---

#### Repeat

Set  $q(\theta|s)$  to  $\pi(\theta|s)$

Generate context samples  $\{s^{[k]}\}_{k=1\dots N}$

Sample parameters  $\{\theta^{[k]}\}_{k=1\dots N}$  from current search distribution  $q(\theta|s)$  given context samples  $\{s^{[k]}\}_{k=1\dots N}$

Evaluate the reward  $R^{[k]}$  of each sample in the sample set  $\{s^{[k]}, \theta^{[k]}\}_{k=1\dots N}$

Use the data set  $\{s^{[k]}, R^{[k]}\}_{k=1\dots N}$  to compute a weight  $d^{[k]}$  for each sample

Use the data set  $\{s^{[k]}, \theta^{[k]}, d^{[k]}\}_{k=1\dots N}$  to update the new search distribution  $\pi(\theta|s)$

Until search distribution  $\pi(\theta|s)$  converges.

---

(Peters, Mülling, and Altun 2010). Recently, an information-theoretic stochastic search algorithm was presented as a modification of REPS (Kupcsik et al. 2013). However, as we will show, the algorithm suffers from premature convergence. In order to solve this problem, we will extend the algorithm by using the rank- $\mu$  covariance matrix adaptation method of CMA-ES algorithm resulting to CREPS-CMA algorithm. The CMA-ES algorithm doesn't work for contextual setting.

Given a dataset  $\{s^{[k]}, \theta^{[k]}, R^{[k]}\}_{k=1\dots N}$  and the current distribution  $q(\theta|s)$  that has been used to generate the unweighed samples  $\{\theta^{[k]}\}_{k=1\dots N}$ , CREPS-CMA first computes a weight  $d^{[k]}$  for each sample and subsequently, using these weights, a new search distribution  $\pi(\theta|s)$  is computed. Therefore, we start by explaining how the weights are computed and, after that, we explain the search distribution update rules.

## Weight Computation

The weight computation for CREPS-CMA is essentially the same as for CREPS (Kupcsik et al. 2013). The key idea behind CREPS is to ensure a smooth and stable learning process by bounding the relative entropy between the old search distribution  $q(\theta|s)$  and the newly estimated policy  $\pi(\theta|s)$ . To do so, we want to optimize the following performance criteria

$$\begin{aligned} & \max_{\pi} \iint \mu(s) \pi(\theta|s) \mathcal{R}_{s\theta} d\theta ds, \\ & \text{s.t.} \int \mu(s) \text{KL}(\pi(\theta|s) || q(\theta|s)) ds \leq \epsilon, \quad (1) \\ & \forall s : 1 = \int \pi(\theta|s) d\theta, \end{aligned}$$

where  $\mu(s)$  denotes the distribution over the context which is specified by the task and  $\mathcal{R}_{s\theta}$  denotes the expected performance when evaluating parameter vector  $\theta$  in context  $s$ . This optimization problem can be solved efficiently by the method of Lagrangian multipliers (Boyd and Vandenberghe

2004). The closed form solution for policy  $\pi(\boldsymbol{\theta}|\mathbf{s})$  is given by

$$\pi(\boldsymbol{\theta}|\mathbf{s}) \propto q(\boldsymbol{\theta}|\mathbf{s}) \exp(\mathcal{R}_{\mathbf{s}\boldsymbol{\theta}}/\eta),$$

where  $\eta$  is a Lagrangian multiplier that sets the temperature of the soft-max distribution given in the previous equation. The temperature parameter  $\eta$  can be found efficiently by optimizing the dual function

$$g(\eta) = \eta\epsilon + \eta \int \mu(\mathbf{s}) \log \left( \int q(\boldsymbol{\theta}|\mathbf{s}) \exp \left( \frac{\mathcal{R}_{\mathbf{s}\boldsymbol{\theta}}}{\eta} \right) d\boldsymbol{\theta} \right) d\mathbf{s}. \quad (2)$$

The optimal value for  $\eta$  can be obtained by minimizing the dual function  $g(\eta)$  such that  $\eta > 0$ , see (Boyd and Vandenberghe 2004). However, approximating the log integral in the dual function (Equation 2) is not feasible as we would need a lot of samples  $\boldsymbol{\theta}^{i,k}$  for each context  $\mathbf{s}^k$ . As the context can often not be directly controlled, we have only access to a single action  $\boldsymbol{\theta}^k$  per context  $\mathbf{s}^k$ .

In order to achieve this requirement, the performance criteria is reformulated. CREP optimizes for the joint probabilities  $p(\mathbf{s}, \boldsymbol{\theta})$  instead of for the policy  $\pi(\boldsymbol{\theta}|\mathbf{s})$ . Additionally, CREPS enforces that  $p(\mathbf{s}) = \int_{\boldsymbol{\theta}} p(\mathbf{s}, \boldsymbol{\theta}) d\boldsymbol{\theta}$  still reproduces the correct context distribution  $\mu(\mathbf{s})$  by using the constraints  $\forall \mathbf{s} : p(\mathbf{s}) = \mu(\mathbf{s})$ . However, there are now an infinite number of constraints. In order to avoid this, these constraints are approximately fulfilled by matching feature expectations, instead of matching single probabilities, i.e.,

$$\int_{\mathbf{s}} p(\mathbf{s}) \boldsymbol{\phi}(\mathbf{s}) d\mathbf{s} = \hat{\boldsymbol{\phi}},$$

where  $\hat{\boldsymbol{\phi}} = \int_{\mathbf{s}} \mu(\mathbf{s}) \boldsymbol{\phi}(\mathbf{s}) d\mathbf{s}$  is the expected feature vector for the given context distribution  $\mu(\mathbf{s})$  and a given feature space  $\boldsymbol{\phi}$ . For example, if the feature vector  $\boldsymbol{\phi}(\mathbf{s})$  contains all linear and squared terms  $[s, s^2]$  of the context vector  $\mathbf{s}$ , these constraints correspond to matching the first and second moment of  $p(\mathbf{s})$  with the moments of  $\mu(\mathbf{s})$ .<sup>2</sup> Note that  $\boldsymbol{\phi}$  is typically different from  $\varphi$ . The resulting optimization program yields

$$\begin{aligned} & \max_p \iint p(\mathbf{s}, \boldsymbol{\theta}) \mathcal{R}_{\mathbf{s}\boldsymbol{\theta}} d\mathbf{s} d\boldsymbol{\theta} \\ & \text{s.t. } \epsilon \geq \text{KL}(p(\mathbf{s}, \boldsymbol{\theta}) || \mu(\mathbf{s}) q(\boldsymbol{\theta}|\mathbf{s})), \\ & \hat{\boldsymbol{\phi}} = \iint p(\mathbf{s}, \boldsymbol{\theta}) \boldsymbol{\phi}(\mathbf{s}) d\mathbf{s} d\boldsymbol{\theta}, \\ & 1 = \iint p(\mathbf{s}, \boldsymbol{\theta}) d\mathbf{s} d\boldsymbol{\theta}. \end{aligned} \quad (3)$$

The solution for  $p(\mathbf{s}, \boldsymbol{\theta})$  is now given by

$$p(\mathbf{s}, \boldsymbol{\theta}) \propto q(\boldsymbol{\theta}|\mathbf{s}) \mu(\mathbf{s}) \exp((\mathcal{R}_{\mathbf{s}\boldsymbol{\theta}} - V(\mathbf{s}))/\eta),$$

where  $V(\mathbf{s}) = \boldsymbol{\phi}(\mathbf{s})^T \mathbf{w}$  is a context dependent baseline which is subtracted from the return  $\mathcal{R}_{\mathbf{s}\boldsymbol{\theta}}$ . The parameters  $\mathbf{w}$  and  $\eta$  are again Lagrangian multipliers that can be obtained by optimizing the dual function, given as

<sup>2</sup>In this paper, in all experiments  $\boldsymbol{\phi}(\mathbf{s}) = [s, s^2]$ , therefore we will always match the first and second moment of  $p(\mathbf{s})$  with the moments of  $\mu(\mathbf{s})$ .

$$\begin{aligned} g(\eta, \mathbf{w}) = & \eta\epsilon + \hat{\boldsymbol{\phi}}^T \mathbf{w} \\ & + \eta \log \left( \iint \mu(\mathbf{s}) q(\boldsymbol{\theta}|\mathbf{s}) \exp \left( \frac{\mathcal{R}_{\mathbf{s}\boldsymbol{\theta}} - \boldsymbol{\phi}(\mathbf{s})^T \mathbf{w}}{\eta} \right) d\boldsymbol{\theta} d\mathbf{s} \right). \end{aligned} \quad (4)$$

This policy update results in a weight

$$d^{[k]} = \exp((\mathcal{R}_{\mathbf{s}\boldsymbol{\theta}} - V(\mathbf{s}))/\eta)$$

for each sample  $[\mathbf{s}^{[k]}, \boldsymbol{\theta}^{[k]}]$ , which we can use to estimate a new search distribution  $\pi_{\pi}(\boldsymbol{\theta}|\mathbf{s})$ .

### Search Distribution Update Rule

Given dataset  $\{\mathbf{s}^{[k]}, \boldsymbol{\theta}^{[k]}, d^{[k]}\}_{k=1\dots N}$  and the old Gaussian search distribution

$$q(\boldsymbol{\theta}|\mathbf{s}) = \mathcal{N}(\boldsymbol{\theta} | \mathbf{m}_q(\mathbf{s}) = \mathbf{A}_q^T \boldsymbol{\varphi}(\mathbf{s}), \Sigma_q),$$

we want to find the new search distribution

$$\pi(\boldsymbol{\theta}|\mathbf{s}) = \mathcal{N}(\boldsymbol{\theta} | \mathbf{m}_{\pi}(\mathbf{s}) = \mathbf{A}_{\pi}^T \boldsymbol{\varphi}(\mathbf{s}), \Sigma_{\pi}),$$

by finding  $\mathbf{A}_{\pi}$  and  $\Sigma_{\pi}$ . Therefore we need two update rules, one for updating the context-dependent mean function  $\mathbf{m}_{\pi}$  of the search distribution and another one for updating the covariance matrix of the distribution  $\Sigma_{\pi}$ .

**Context-Dependent Mean-Function Update Rule** In order to find  $\mathbf{m}_{\pi}$ , the parameters  $\mathbf{A}$  can be obtained by the weighted maximum likelihood

$$\mathbf{A} = (\boldsymbol{\Phi}^T \mathbf{D} \boldsymbol{\Phi} + \lambda \mathbf{I})^{-1} \boldsymbol{\Phi}^T \mathbf{D} \mathbf{U}, \quad (5)$$

where  $\boldsymbol{\Phi}^T = [\boldsymbol{\varphi}^{[1]}, \dots, \boldsymbol{\varphi}^{[N]}]$  contains the feature vector for all samples,  $\mathbf{U} = [\boldsymbol{\theta}^{[1]}, \dots, \boldsymbol{\theta}^{[N]}]$  contains all the sample parameters and  $\mathbf{D}$  is the diagonal weighting matrix containing the weightings  $d^{[k]}$ .

**Covariance Matrix Update Rule** Standard CREPS directly uses the weighted sample covariance matrix  $\mathbf{S}$  as  $\Sigma_{\pi}$  which is obtained by

$$\begin{aligned} \mathbf{S} &= \frac{\sum_{k=1}^N d^{[k]} (\boldsymbol{\theta}^{[k]} - \mathbf{A}^T \boldsymbol{\varphi}(\mathbf{s}^{[k]})) (\boldsymbol{\theta}^{[k]} - \mathbf{A}^T \boldsymbol{\varphi}(\mathbf{s}^{[k]}))^T}{Z}, \\ Z &= \frac{(\sum_{k=1}^N d^{[k]})^2 - \sum_{k=1}^N (d^{[k]})^2}{\sum_{k=1}^N (d^{[k]})}. \end{aligned} \quad (6)$$

Typically, the number of samples used for the estimated sample covariance matrix  $\mathbf{S}$  is much smaller than number of free parameters of the covariance matrix. In this case, it has been shown that the sample covariance matrix from Equation 6 is not a good estimate of the true covariance matrix (Abdolmaleki et al. 2015b) and biases the search distribution towards a specific region of the search space. Due to this effect, the search distribution uncontrollably loses its exploration entropy along many dimensions of the parameter space which causes premature convergence. Such loss of exploration is a highly unwanted effect in stochastic

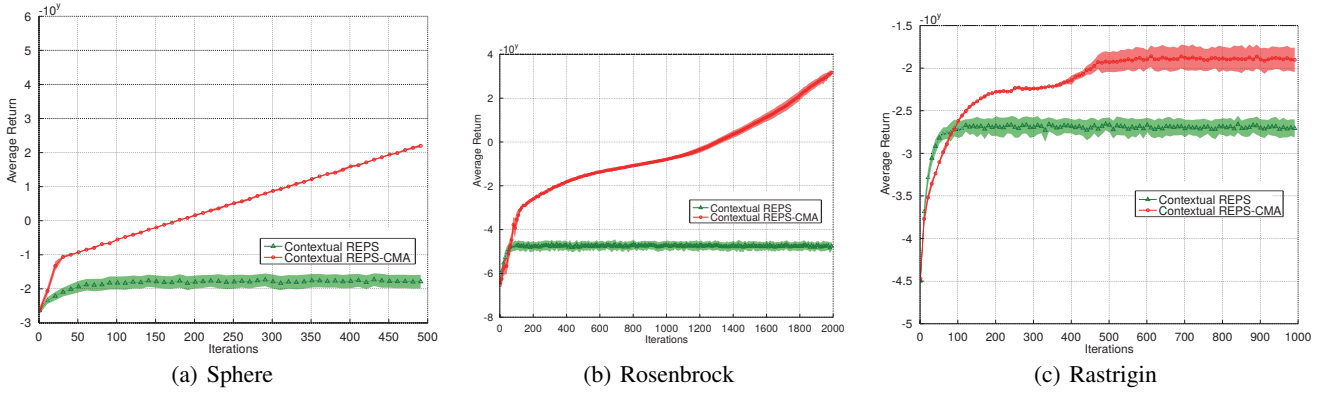


Figure 1: The performance comparison of stochastic search methods for optimising contextual version of standard functions (a) Sphere, (b) Rosenbrock and (c) Rastrigin. The results show that while CREPS-CMA performs well, Contextual REPS suffers from premature convergence.

search. In order to alleviate this problem, inspired by CMA-ES (Hansen, Muller, and Koumoutsakos 2003), which is not a contextual algorithm, we combine the old covariance matrix and the sample covariance matrix from Equation 6, i.e.,

$$\Sigma_{\pi} = (1 - \lambda)\Sigma_q + (\lambda)\mathbf{S}.$$

There are different ways to determine the interpolation factor  $\lambda \in [0, 1]$  between the sample covariance matrix  $\mathbf{S}$  and the old covariance matrix  $\Sigma_q$ . For example, in (Abdolmaleki et al. 2015b), the factor  $\lambda \in [0, 1]$  is chosen in such a way that the entropy of the new search distribution is reduced by a certain amount, while also being scaled with the number of effective samples. For CREPS-CMA, we use the rank- $\mu$  update in CMA-ES algorithm (Hansen, Muller, and Koumoutsakos 2003), i.e.,

$$\phi_{\text{eff}} = \frac{1}{\sum_{k=1}^N (d^{[k]})^2}, \quad \lambda = \min\left(1, \frac{\phi_{\text{eff}}}{p^2}\right)$$

where  $\phi_{\text{eff}}$  is the number of effective samples and  $p$  is the dimension of the parameter space  $\theta$ . See Algorithm 2 for a compact representation of the CREPS-CMA algorithm.

## Experiments

In this section, We evaluate CREPS-CMA algorithm.<sup>3</sup> We chose three series of optimization tasks for comparisons. In the first series, we use standard optimization test functions (Molga and Smutnicki 2005), such as the Sphere, the Rosenbrock and the Rastrigin (multi-modal) functions. We extend these functions to be applicable on the contextual setting. The task is to find the optimum 15 dimensional parameter vector  $\theta$  for a given 1 dimensional context  $s$ .

For the second series of optimization tasks, we use a 5-link planar robot that has to reach a given point, which is varied by the context. We used dynamic movement primitives (DMPs) (Ijspeert and Schaal 2003) as the underlying policy representation with 25 parameters (5 basis functions

<sup>3</sup>Matlab source-code is available on-line at <https://dl.dropboxusercontent.com/u/16387578/ContextualREPS-CMA.zip>

per dimension). The context is the two dimensional position of the via-point to reach with the end-effector. Figure 2 shows the setup of the robot.

For the third series of optimization tasks, we used a simulated humanoid robot that has to kick a ball with a certain distance, which is varied by the context. We used a simple linear model as a controller, which takes as input the starting and end joint positions of the robot’s legs and the time the movement is intended to take. The model linearly interpolates the starting and end positions through time. The context is the distance that the ball should reach. Figure 3 shows an example of the robot movement.

In Figures 1 and 2, we show the average, as well as two times the standard deviation of the results, over 10 trials for each experiment. Note that the y-axis of all plots is in a logarithmic scale. Figure 3(c) shows the average and two times the standard deviation of the results over 10 trials for several contexts of the humanoid kick.

## Standard Optimization Test Functions

We chose three standard optimization functions, which are the Sphere function

$$f(s, \theta) = \sum_{i=1}^p x_i^2,$$

the Rosenbrock function

$$f(s, \theta) = \sum_{i=1}^{p-1} [100(x_{i+1} - x_i^2)^2 + (1 - x_i)^2],$$

and also a multi-modal function, known as the Rastrigin function

$$f(s, \theta) = 10p + \sum_{i=1}^p [x_i^2 - 10 \cos(2\pi x_i)],$$

where  $x = \theta + \mathbf{A}s$ . The matrix  $A$  is a constant matrix that was chosen randomly. In our case, because the context  $s$  is 1 dimensional,  $A$  is a  $p \times 1$  dimensional vector. Our definition

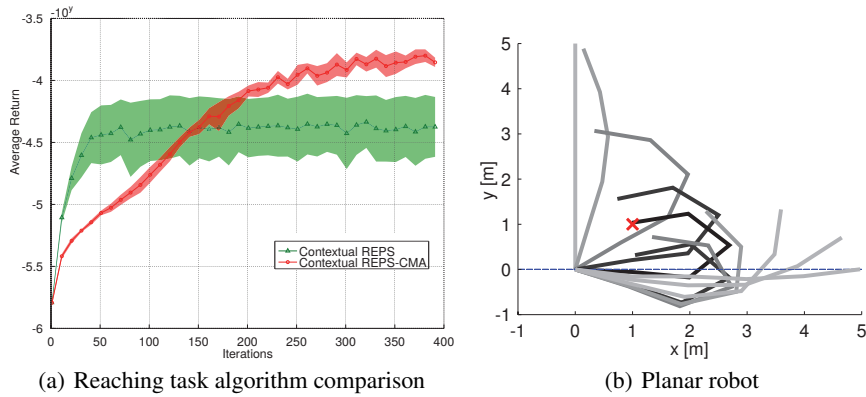


Figure 2: (a) Algorithmic comparison for a planar reaching task (5 joints, 25 parameters). In this task, CREPS-CMA has converged faster and learned the task well. Contextual REPS suffers from premature convergence and cannot learn the task. (b) The planar reaching task used for our comparisons. A 5-link planar robot has to reach a via-point  $v_{50} = [1, 1]$  in task space. The via-point position is the 2 dimensional context vector and is given. The via-point is indicated by the red cross. The postures of the resulting motion are shown as overlay, where darker postures indicate a posture which is close in time to the via-point.

for  $x$  means the optimum  $\theta$  for these functions is linearly dependent on the given context  $s$ . The initial search area of  $\theta$  for all experiments is restricted to the hypercube  $-5 \leq \theta_i \leq 5, i = 1, \dots, p$  and contexts are uniformly sampled from the interval  $0 \leq s_i \leq 3, i = 1, \dots, z$  where  $z$  is the dimension of the context space  $s$ . In our experiments, the mean of the initial distributions has been chosen randomly in the defined search area.

We compared CREPS-CMA with the standard Contextual REPS. In each iteration, we generated 50 new samples. The results in Figure 1 show that CREPS-CMA could successfully learn the contextual tasks while standard Contextual REPS suffers from premature convergence.

### Planar Reaching

In this task, we used a 5-link planar robot with DMPs (Ijspeert and Schaal 2003) as the underlying control policy. Each link had a length of  $1m$ . The robot is modelled as a decoupled linear dynamical system. The robot has to reach a via-point  $v_{50} = [1, 1]$  at time step 50 with its end effector and, at the final time step  $T = 100$ , the point  $v_{100} = [5, 0]$ . We varied the first via-point from 0 to 2 in both dimensions.

The reward was given by a context-dependent quadratic cost term for the two via-points as well as quadratic costs for high accelerations. Note that this performance function is highly non-linear in the parameters as the via-points are defined in end effector space. We used 5 basis functions per degree of freedom for the DMPs while the goal attractor for reaching the final state was assumed to be known. Hence, our parameter vector had 25 dimensions.

We generated 100 new samples in each iteration. The results in Figure 2(a) shows that CREPS-CMA successfully learned the task while contextual REPS cannot complete it.

### Humanoid Kick

In this task, robot should learn to kick the ball to different distances. The context here is the distance that ball travels

which varies from  $3m$  to  $12m$ . The robot kick controller is a simple linear model with three input groups: the action time  $t$ , the initial position of the robot, represented as a vector of joint angles with dimension  $l$ , and the final robot position, also represented as a vector of joint angles with dimension  $l$ . The action time is the amount of time the robot takes to move from the initial to the final position. The joint angles are linearly interpolated across  $t$  to create the corresponding movement.

Our humanoid robot has  $l = 6$  joints in each leg, and other remaining joints (arms and head) are ignored. In other words, the joint angles are 12-dimensional vectors and the controller receives a 25-dimensional parameter vector, which is then interpolated and coded into motor commands. Figures 3(a) and 3(b) show the initial and final positions of an exemplary kick.

The reward was based on the reward function

$$R(s, \theta) = -(x - s)^2 - y^2,$$

where  $x$  and  $y$  are the distances between the target point and the ball along the x and y axes. We use CREPS-CMA to optimise this reward function and similar to (Abdolmaleki et al. 2015c) we use radial basis feature function for non-linear generalization over contexts. After 1000 iterations with 20 new samples per iteration, the robot could kick the ball with a good accuracy shown on Figure 3(c).<sup>4</sup> The average error distance of the ball to the target was  $0.34 \pm 0.11m$ .

### Conclusion

Many optimization algorithms have been proposed by the scientific community. However, these algorithms usually optimize for a single context of a task, like the ideal gait for the highest speed, the lowest energy consumption, or both.

<sup>4</sup>Demonstration video is available on-line at <https://dl.dropboxusercontent.com/u/16387578/ICARSC16kick.mp4>.

The robot is placed at several points and kicks the ball to the centre of the field using the learned policy.

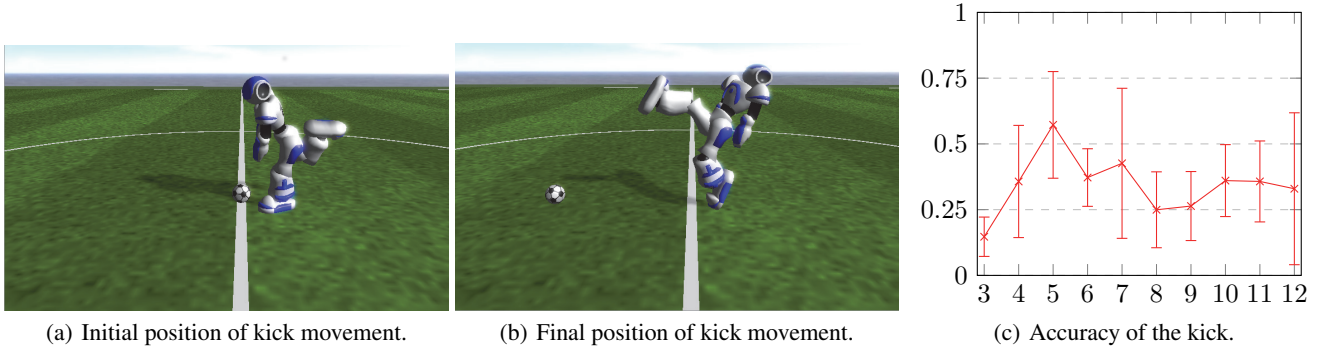


Figure 3: (a) The initial position of an exemplary humanoid kick. (b) The final position of an exemplary humanoid kick. (c) The results of the task. The y-axis represents the distance at which the ball was from the intended target, in meters, while the x-axis represents the distance at which the ball was being kicked from, also in meters.

---

### Algorithm 2 CREPS-CMA

---

**Input :** Data Set  $\mathcal{D}\{\mathbf{s}^{[k]}, \boldsymbol{\theta}^{[k]}, R^{[k]}\}_{k=1\dots N}$  and the old covariance matrix  $\Sigma_q$

**Compute the weights  $d^{[k]}$  for each sample:**

1- Optimize the dual function  $g$  and find optimum  $\eta$  and  $\mathbf{w}$

$$g(\eta, \mathbf{w}) = \eta\epsilon + \hat{\phi}^T \mathbf{w} + \eta \log \left( \sum_{K=1}^N \frac{1}{N} \exp \left( \frac{R^{[k]} - \phi(\mathbf{s}^{[k]})^T \mathbf{w}}{\eta} \right) \right).$$

2- Compute weights  $d^{[k]} = \exp \left( \frac{R^{[k]} - \phi(\mathbf{s}^{[k]})^T \mathbf{w}}{\eta} \right)$ .

3- Normalize  $d^{[k]}$  such that  $\sum_{k=1}^N d^{[k]} = 1$ .

**Compute the new mean function  $m_\pi(\mathbf{s})$ :**

Use weighted maximum likelihood to estimate parameters  $\mathbf{A}_\pi$  of the new mean function

$$\mathbf{A}_\pi = (\Phi^T \mathbf{D} \Phi + \lambda \mathbf{I})^{-1} \Phi^T \mathbf{D} \mathbf{U},$$

where  $\Phi^T = [\varphi^{[1]}, \dots, \varphi^{[N]}]$  contains the feature vector for all samples,  $\mathbf{U} = [\boldsymbol{\theta}^{[1]}, \dots, \boldsymbol{\theta}^{[N]}]$  contains all the sample parameters and  $\mathbf{D}$  is the diagonal weighting matrix containing the weights  $d^{[k]}$ .

**Compute the sample covariance  $\mathbf{S}$ :**

$$\mathbf{S} = \sum_{k=1}^N d^{[k]} (\boldsymbol{\theta}^{[k]} - \mathbf{A}_\pi^T \varphi(\mathbf{s}^{[k]})) (\boldsymbol{\theta}^{[k]} - \mathbf{A}_\pi^T \varphi(\mathbf{s}^{[k]}))^T / Z,$$

**Compute the number of effective samples  $\phi_{\text{eff}}$  and  $\lambda$ :**

$$\phi_{\text{eff}} = \frac{1}{\sum_{k=1}^N (d^{[k]})^2}, \quad \lambda = \min \left( 1, \frac{\phi_{\text{eff}}}{p^2} \right)$$

**Compute the new covariance matrix  $\Sigma$ :**

$$\Sigma_\pi = \lambda \Sigma_q + (1 - \lambda) \mathbf{S}.$$


---

Therefore, in this paper, we investigated contextual stochastic search methods for multi-task learning. We alleviated the premature convergence problem of contextual REPS, which resulted in the CREPS-CMA algorithm. The results show that the algorithm performs favourably and solves the premature convergence issue. We also show its applicability in practical situations, such as a humanoid robot kick task. Regarding future work, it is worthwhile to incorporate the step-size control feature of CMA-ES into CREPS-CMA. This step-size control aims to make consecutive movements of the distribution mean orthogonal in expectation. It effectively prevents premature convergence while allowing fast convergence to an optimum.

### References

- Abdolmaleki, A.; Lioutikov, R.; Peters, J.; Lua, N.; Reis, L.; and Neumann, G. 2015a. Regularized covariance estimation for weighted maximum likelihood policy search methods. In *Advances in Neural Information Processing Systems (NIPS)*, MIT Press.
- Abdolmaleki, A.; Lua, N.; Reis, L.; and Neumann, G. 2015b. Regularized covariance estimation for weighted maximum likelihood policy search methods. In *Proceedings of the International Conference on Humanoid Robots (HUMANOIDS)*.
- Abdolmaleki, A.; Lua, N.; Reis, L.; Peters, J.; and Neumann, G. 2015c. Contextual Policy Search for Generalizing a Parameterized Biped Walking Controller. In *IEEE International Conference on Autonomous Robot Systems and Competitions (ICARSC)*.
- Boyd, S., and Vandenberghe, L. 2004. *Convex optimization*. Cambridge university press.
- Broomhead, D. S., and Lowe, D. 1988. Radial basis functions, multi-variable functional interpolation and adaptive networks. Technical report, DTIC Document.
- Daniel, C.; Neumann, G.; and Peters, J. 2012. Hierarchical Relative Entropy Policy Search. In *International Conference on Artificial Intelligence and Statistics (AISTATS)*.
- Hansen, N.; Muller, S.; and Koumoutsakos, P. 2003. Re-

ducing the Time Complexity of the Derandomized Evolution Strategy with Covariance Matrix Adaptation (CMA-ES). *Evolutionary Computation*.

Ijspeert, A., and Schaal, S. 2003. Learning Attractor Landscapes for Learning Motor Primitives. In *Advances in Neural Information Processing Systems 15(NIPS)*.

Kober, J., and Peters, J. 2010. Policy Search for Motor Primitives in Robotics. *Machine Learning* 1–33.

Kober, J.; Oztop, E.; and Peters, J. 2010. Reinforcement Learning to adjust Robot Movements to New Situations. In *Proceedings of the Robotics: Science and Systems Conference (RSS)*.

Kupcsik, A.; Deisenroth, M. P.; Peters, J.; and Neumann, G. 2013. Data-Efficient Contextual Policy Search for Robot Movement Skills. In *Proceedings of the National Conference on Artificial Intelligence (AAAI)*.

Mannor, S.; Rubinstein, R.; and Gat, Y. 2003. The Cross Entropy method for Fast Policy Search. In *Proceedings of the 20th International Conference on Machine Learning(ICML)*.

Molga, M., and Smutnicki, C. 2005. Test Functions for Optimization Needs. In <http://www.zsd.ict.pwr.wroc.pl/files/docs/functions.pdf>.

Niehaus, C.; Röfer, T.; and Laue, T. 2007. Gait optimization on a humanoid robot using particle swarm optimization. In *Proceedings of the Second Workshop on Humanoid Soccer Robots in conjunction with the*, 1–7.

Peters, J.; Mülling, K.; and Altun, Y. 2010. Relative Entropy Policy Search. In *Proceedings of the 24th National Conference on Artificial Intelligence (AAAI)*. AAAI Press.

Rückstieß, T.; Felder, M.; and Schmidhuber, J. 2008. State-dependent Exploration for Policy Gradient Methods. In *Proceedings of the European Conference on Machine Learning (ECML)*.

Stulp, F., and Sigaud, O. 2012. Path Integral Policy Improvement with Covariance Matrix Adaptation. In *International Conference on Machine Learning (ICML)*.

Sun, Y.; Wierstra, D.; Schaul, T.; and Schmidhuber, J. 2009. Efficient Natural Evolution Strategies. In *Proceedings of the 11th Annual conference on Genetic and evolutionary computation(GECCO)*.

Theodorou, E.; Buchli, J.; and Schaal, S. 2010. A Generalized Path Integral Control Approach to Reinforcement Learning. *The Journal of Machine Learning Research*.

Wang, J. M.; Fleet, D. J.; and Hertzmann, A. 2009. Optimizing walking controllers. *ACM Transactions on Graphics (TOG)* 28(5):168.