

---

# Active Contextual Entropy Search

---

**Jan Hendrik Metzen**

Universität Bremen, 28359 Bremen, Germany  
DFKI GmbH, Robotics Innovation Center, 28359 Bremen, Germany  
jhm@informatik.uni-bremen.de

## Abstract

Contextual policy search allows adapting robotic movement primitives to different situations. For instance, a locomotion primitive might be adapted to different terrain inclinations or desired walking speeds. Such an adaptation is often achievable by modifying a small number of hyperparameters. However, learning, when performed on real robotic systems, is typically restricted to a small number of trials. Bayesian optimization has recently been proposed as a sample-efficient means for contextual policy search that is well suited under these conditions. In this work, we extend entropy search, a variant of Bayesian optimization, such that it can be used for *active* contextual policy search where the agent selects those tasks during training in which it expects to learn the most. Empirical results in simulation suggest that this allows learning successful behavior with less trials.

## 1 INTRODUCTION

Contextual policy search (CPS) is a popular means for multi-task reinforcement learning in robotic control [6]. CPS learns a hierarchical policy, in which the lower-level policy is often a domain-specific behavior representation such as dynamical movement primitives (DMPs) [12]. Learning takes place on the upper-level policy, which is a conditional probability density  $\pi(\theta|s)$  that defines a distribution over the parameter vectors  $\theta$  of the lower-level policy for a given context  $s$ . The context  $s$  encodes properties of environment or task such as a desired walking speed for a locomotion behavior or a desired target position for a ball-throw behavior. The objective of CPS is to learn an upper-level policy which maximizes the expected return of the lower-level policy for a given context distribution.

CPS is typically based on local search based approaches such as cost-regularized kernel regression [14] and contextual relative entropy search (C-REPS) [17, 21]. From the field of black-box optimization, it is well-known that local search based approaches are well suited for problems with a moderate dimensionality and no gradient-information. However, for the special case of relatively low-dimensional search spaces combined with an expensive cost function, which limits the number of evaluations of the cost functions, global search approaches like Bayesian optimization [2] are often superior, for instance for selecting hyperparameters [25]. Combining contextual policy search with pre-trained movement primitives<sup>1</sup> can also fall into this category as evaluating the cost function requires an execution of the behavior on the robot while only a small set of hyperparameters might have to be adapted. Bayesian optimization has been used for non-contextual policy search on locomotion tasks [5, 18] and robot grasping [16] and for (passive) contextual policy search on a simulated robotic ball-throwing task [20].

In this work, we focus on problems where the agent can select the task (context) in which it will perform the next trial during learning. This facilitates active learning, which is considered to be a prerequisite for lifelong learning [23]. A core challenge in active multi-task robot control learning is the incommensurability of performance in different tasks, i.e., how a learning system can account for

---

<sup>1</sup>DMPs can be pre-trained for fixed contexts in simulation or via some kind of imitation learning.

the relative (unknown) *difficulty* of a task: for instance, if a relatively small reward is obtained when executing a specific low-level policy in a task, is it because the low-level policy is not well adapted to the task or because the task is inherently more difficult than other tasks? Fabisch et al. [8] presented an approach for estimating the task-difficulty explicitly, which allows defining heuristic intrinsic reward functions based on which a discounted multi-arm bandit selects the next task actively [7].

In this work, we follow a different approach: rather than explicitly addressing the incommensurability of rewards, we propose ACES, an information theoretic approach for active task selection which selects tasks not based on rewards directly but rather based on the expected reduction in uncertainty about the optimal parameters for the contexts. ACES allows selecting task and parameters jointly without requiring a heuristic definition of a task selection criterion. ACES is motivated by entropy search [10], which has been extended in a similar fashion to non-contextual settings [19], and multi-task Bayesian optimization [27], which focuses on problems with discrete context spaces.

## 2 BACKGROUND

**Contextual Policy Search** (CPS) denotes a model-free approach to reinforcement learning, in which the (low-level) policy  $\pi_\theta$  is parametrized by a vector  $\theta$ . The choice of  $\theta$  is governed by an upper-level policy  $\pi_\omega$ . For generalizing learned policies to multiple tasks, the task is characterized by a context vector  $s$  and the upper-level policy  $\pi_\omega(\theta|s)$  is conditioned on the respective context. The objective of CPS is to learn  $\pi_\omega$  such that the expected return  $J_\omega$  over all contexts is maximized, with  $J_\omega = \int_s p(s) \int_\theta \pi_\omega(\theta|s) R(\theta, s) d\theta ds$ . Here,  $p(s)$  is the distribution over contexts and  $R(\theta, s)$  is the expected return when executing the low level policy with parameter  $\theta$  in context  $s$ . We refer to Deisenroth et al. [6] for a recent overview of (contextual) policy search approaches in robotics.

**Bayesian optimization for contextual policy search** (BO-CPS) is based on applying ideas from Bayesian optimization to contextual policy search [20]. BO-CPS learns internally a model of the expected return  $R(\theta, s)$  of a parameter vector  $\theta$  in a context  $s$ . The model is based on Gaussian process (GP) regression [22]. It learns from sample returns  $R_i$  obtained in rollouts at query points consisting of a context  $s_i$  determined by the environment and a parameter vector  $\theta_i$  selected by BO-CPS. By learning a joint GP model over the context-parameter space, experience collected in one context is naturally generalized to similar contexts.

The GP model provides both an estimate of the expected return  $\mu_{GP}[R(s, \theta)]$  and its standard deviation  $\sigma_{GP}[R(s, \theta)]$ . Based on this information, the parameter vector for the given context is selected by maximizing an *acquisition function*, which allows controlling the trade-off between exploitation (selecting parameters with maximal estimated return) and exploration (selecting parameters with high uncertainty). Common acquisition functions used in Bayesian optimization such as the probability of improvement (PI) and the expected improvement (EI) [2] are not easily generalized to BO-CPS [20]. In contrast, the acquisition function GP-UCB [26], which defines the acquisition value of a parameter vector in a context as  $\text{GP-UCB}(s, \theta) = \mu_{GP}[R(s, \theta)] + \kappa \sigma_{GP}[R(s, \theta)]$ , where  $\kappa$  controls the exploration-exploitation trade-off, can be applied to BO-CPS straightforwardly resulting in an approach similar to CGP-UCB [15]. BO-CPS selects parameters  $\theta_i$  for a given fixed context  $s_i$  by performing an optimization over the parameter space using the global maximizer DIRECT [13] to find the approximate global maximum, followed by L-BFGS [4] to refine it.

**Entropy search** (ES) is a recently proposed approach to probabilistic global optimization that mainly differs from Bayesian optimization in the choice of the acquisition function [10]. While typical acquisition functions used for Bayesian optimization select query points where they expect the optimum, ES selects query points where it expects to learn most about the optimum. More specifically, ES explicitly represents  $p_{opt}(\theta)$ , the probability that the global optimum (maximum or minimum, depending on the problem) of the unknown function  $f$  is at  $\theta$ . ES estimates  $p_{opt}(\theta)$  at finitely many points  $\{\theta^c\}_{i=1}^{N_\theta}$  on a non-uniform grid that are selected heuristically. Moreover, it approximates  $p_{opt}$  at  $\theta^c$  based on expectation propagation or Monte Carlo integration. To select a query point, ES predicts the change of the GP when drawing a sample at the query point  $\theta^q$  and assuming  $N_y$  different outcomes  $\{y^{(i)}\}$  sampled from the GP’s predictive distribution at  $\theta^q$ . Thereupon, ES selects a query point which minimizes the average loss  $\mathcal{L}(p_{opt}[\theta^q]) = - \int p_{opt}[\theta^q](\theta) \log \frac{p_{opt}[\theta^q](\theta)}{U_I(\theta)} d\theta$ , i.e., which maximizes the relative entropy between  $p_{opt}$  and a uniform measure  $U_I$ , where  $p_{opt}[\theta^q]$  denotes the probability distribution of the global optimum *after* an assumed query at  $\theta^q$ .

### 3 ACTIVE CONTEXTUAL ENTROPY SEARCH

In this section, we present active contextual entropy search (ACES), an extension of ES to CPS which allows selecting both parameters  $\theta_q$  and context  $s_q$  of the next trial. Let  $p_{max}(\theta|s)$  denote the conditional probability distribution of the maximum expected return given context  $s$  and let the loss  $\mathcal{L}^s(s^q, \theta^q) = \mathcal{L}(p_{max}[s^q, \theta^q](\theta|s)) - \mathcal{L}(p_{max}(\theta|s))$  denote the expected change of relative entropy in context  $s$  after performing a trial in context  $s^q$  with parameter  $\theta^q$ . A straightforward extension of ES to active learning in BO-CPS would be selecting  $s^q, \theta^q = \arg \min_{(s^q, \theta^q)} \mathcal{L}^{s^q}(s^q, \theta^q)$ , i.e., select the context  $s^q$  in which the maximum increase of relative entropy is expected. This, however, would not account for information gained about the optima in contexts  $s \neq s^q$  by a query at  $(s^q, \theta^q)$ .

ACES instead averages over the expected change in relative entropy at different points in the context space:  $\text{ACES}(s^q, \theta^q) = \sum_{i=1}^{N_s} \mathcal{L}^{s_i^c}(s^q, \theta^q)$ , where  $\{s_i^c\}_{i=1}^{N_s}$  is a set of contexts which is drawn uniform randomly from the context space. Unfortunately, each evaluation of  $\mathcal{L}^{s_i^c}$  is computationally expensive and thus  $N_s$  would have to be chosen small. On the other hand, GPs have an intrinsic length-scale for many choices of the kernel and thus, a query in context  $s^q$  will only affect  $\mathcal{L}^{s_i^c}$  when  $s_i^c$  is “similar” to  $s^q$ . We define similarity between contexts based on the Mahalanobis distance  $d_M(s_i^c, s^q) = \sqrt{(s_i^c - s^q)^T S^{-1} (s_i^c - s^q)}$  with  $S$  being a diagonal matrix with the (anisotropic) length scales of the GP on the diagonal. Based on this we can approximate  $\text{ACES}(s^q, \theta^q) \approx \sum_{s \in \text{NN}(s^q, \{s^c\}, N_{nn})} \mathcal{L}^s(s^q, \theta^q)$  with NN returning the  $N_{nn}$  nearest neighbors of  $s^q$  in  $\{s^c\}$  according to the Mahalanobis distance. A larger value of  $N_{nn}$  corresponds to a better approximation of  $\text{ACES}(s^q, \theta^q)$  at the cost of a linearly increased computational cost.

Candidate points  $\theta^c(s)$  are selected by performing Thompson sampling on 500 randomly chosen  $\theta$  with  $N_\theta = 20$ . The number of trial contexts  $N_s$  is set to 100 and we compare empirically  $N_{nn} = 1$  and  $N_{nn} = 20$ . The quantity  $p_{max}$  is approximated using Monte-Carlo integration based on drawing 1000 samples from the GP posterior. The number of samples from the GP’s predictive distribution at  $\theta_q$  for approximating the average loss for a query point is set to  $N_y = 10$ . Since there is noise in the Monte Carlo estimates of  $\mathcal{L}^s(s^q, \theta^q)$ , we use CMA-ES [9] as optimizer rather than DIRECT.

### 4 EVALUATION

We present results in a simulated robotic control task, in which the robot arm COMPI [1] is used to throw a ball at a target on the ground encoded in a two-dimensional context vector. The target area is  $[1, 2.5]m \times [-1, 1]m$  and the robot arm is mounted at the origin  $(0, 0)$  of this coordinate system. Thus, contexts can be chosen from a two-dimensional context space:  $s \in [1, 2.5] \times [-1, 1]$ . The low-level policy is a joint-space DMP with preselected start and goal angle for each joint and all DMP weights set to 0. This DMP results in throwing a ball such that it hits the ground close to the center of the target area. Adaptation to different target positions is achieved by modifying a two-dimensional vector  $\theta$ : the first component of  $\theta$  corresponds to the execution time  $\tau$  of the DMP, which determines how far the ball is thrown, and the second component to the final angle  $g_0$  of the first joint, which determines the rotation of the arm around the z-axis.

The upper-level policy<sup>2</sup> is a deterministic policy which selects parameters  $\theta$  based on an affine function of context  $s$ . This policy is trained on the training data  $\{(s_i, \theta_i, R_i)\}_i$  using the C-REPS policy update. The limits on the parameter space are  $g_0 \in [-\frac{\pi}{2}, \frac{\pi}{2}]$  and  $\tau \in [0.4, 2]$ . All approaches use an anisotropic Matérn kernel for the GP surrogate model. Since we focus on a “pure exploration” scenario [3], GP-UCB’s exploration parameter  $\kappa$  is set to a constant value of 5.0. The reward is defined as  $r = -\|s - b_s\|^2 - 0.01 \sum_t v_t^2$ , where  $s$  denotes the goal position,  $b_s$  denotes the position hit by the ball, and  $\sum_t v_t^2$  denotes a penalty term on the sum of squared joint velocities during DMP execution. The maximum achievable reward for different  $b_s$  differs as values of  $b_s$  further away from the origin (where the arm is mounted) require larger joint velocities  $v_t$  and thus incur a larger penalty. Thus, rewards in different contexts are incommensurable.

Figure 1 summarizes the main results of the empirical evaluation. The left graph shows the mean offline performance of the upper-level policy at 16 test contexts on a grid over the context space.

<sup>2</sup>The upper-level policy could in principle be defined directly on the surrogate model. This would, however, require a computationally expensive maximization over the parameter space for each evaluation of the policy.

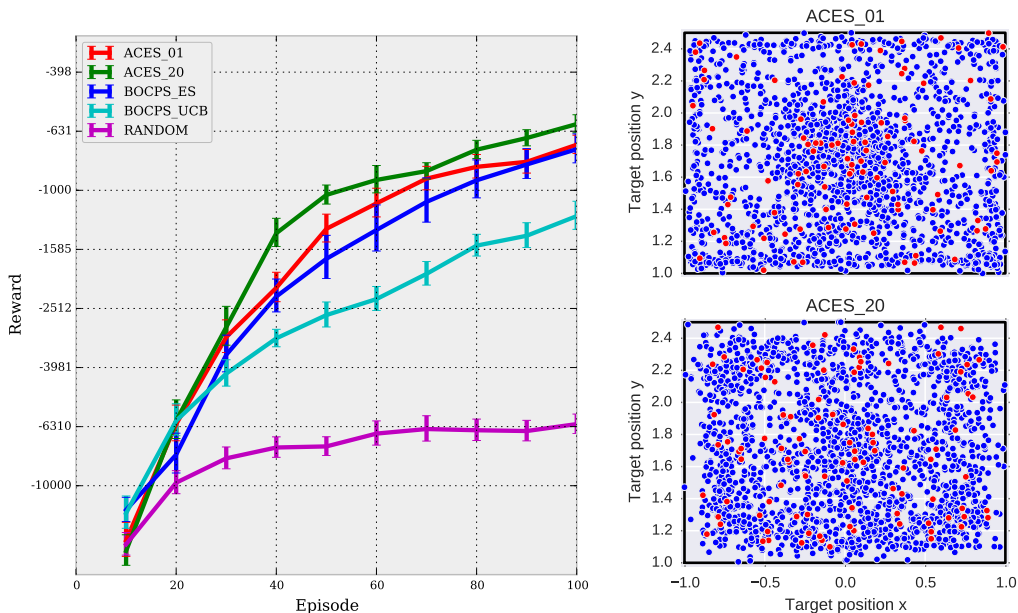


Figure 1: (Left) Learning curves on the simulated robot arm COMPI: the offline performance is evaluated each 10 episodes on 16 test contexts distributed on a grid over the target area. Shown are mean and its standard error over 20 independent runs. (Right) Scatter plot showing the sampled contexts for all (blue) and a single representative run (red).

Sampling contexts and parameters randomly during learning (“Random”) is shown as a baseline and indicates that generalizing experience using a GP model alone does not suffice for quick learning in this task. Rather, a non-random way of exploration is required. BO-CPS with random context selection and UCB for parameter selection improves considerably over random parameter selection. Using ES for parameter selection further improves the learning speed. Closer inspection (not shown) indicates that ES improves over UCB mainly because UCB samples often at the boundaries of the parameter space since the uncertainty is typically large there. ES samples more often in the inner regions of the parameter space since those regions promise a larger information gain globally.

Active context selection using ACES further improves over BOCPS-ES, in particular when the sum over the context space is approximated using several samples ( $N_{nn} = 20$  in the case of ACES\_20) rather than a single sample ( $N_{nn} = 1$  for ACES\_01). The right graph shows the contexts selected by different variants of ACES. It can be seen that ACES\_20 avoids selecting targets close to the boundary of the context space as those typically reveal less global information about the context-dependent optima as boundary points are far away from most other regions of the context space. We attribute the improved learning performance to this way of selecting targets during learning. In contrast, ACES\_1 samples more often close to the boundaries as it only considers the local information gain and thus has no reason to prefer inner over boundary contexts.

## 5 DISCUSSION AND CONCLUSION

We have presented an active learning approach for contextual policy search based on entropy search. First experimental results indicate that the proposed active learning approach provides considerable speed-ups of the learning of movement primitives compared to a random task selection. Comparison with other active task selection approaches [7] remains future work. Moreover, investigating and enhancing the scalability to higher dimensional problems, potentially by employing a combination with random embedding approaches such as REMBO [28], and combining active task selection with predictive entropy search [11] or portfolio-based approaches [24] would be interesting.

**Acknowledgments** This work was performed as part of the project BesMan<sup>3</sup> and supported through two grants of the German Federal Ministry of Economics and Technology (BMWi, FKZ 50 RA 1216 and FKZ 50 RA 1217).

## References

- [1] V. Bargsten and J. de Gea. COMPI: Development of a 6-DOF compliant robot arm for human-robot cooperation. In *Proceedings of the 8th International Workshop on Human-Friendly Robotics (HFR-2015)*. Technische Universität München (TUM), Oct. 2015.
- [2] E. Brochu, V. M. Cora, and N. De Freitas. A tutorial on Bayesian optimization of expensive cost functions, with application to active user modeling and hierarchical reinforcement learning. Technical report.
- [3] S. Bubeck, R. Munos, and G. Stoltz. Pure Exploration in Multi-armed Bandits Problems. In *Algorithmic Learning Theory*, number 5809 in Lecture Notes in Computer Science, pages 23–37. Springer Berlin Heidelberg, Oct. 2009.
- [4] R. H. Byrd, P. Lu, J. Nocedal, and C. Zhu. A Limited-Memory Algorithm for Bound Constrained Optimization. *SIAM Journal on Scientific Computing*, 16:1190–1208, 1995.
- [5] R. Calandra, A. Seyfarth, J. Peters, and M. P. Deisenroth. Bayesian Optimization for Learning Gaits under Uncertainty. *Annals of Mathematics and Artificial Intelligence (AMAI)*, 2015.
- [6] M. P. Deisenroth, G. Neumann, and J. Peters. A Survey on Policy Search for Robotics. *Foundations and Trends in Robotics*, 2(1-2):1–142, 2013.
- [7] A. Fabisch and J. H. Metzen. Active contextual policy search. *Journal of Machine Learning Research*, 15:3371–3399, 2014.
- [8] A. Fabisch, J. H. Metzen, M. M. Krell, and F. Kirchner. Accounting for Task-Difficulty in Active Multi-Task Robot Control Learning. *KI - Künstliche Intelligenz*, pages 1–9, May 2015.
- [9] N. Hansen and A. Ostermeier. Completely derandomized self-adaptation in evolution strategies. *Evolutionary Computation*, 9:159–195, 2001.
- [10] P. Hennig and C. J. Schuler. Entropy Search for Information-Efficient Global Optimization. *JMLR*, 13:1809–1837, 2012.
- [11] J. M. Hernández-Lobato, M. W. Hoffman, and Z. Ghahramani. Predictive Entropy Search for Efficient Global Optimization of Black-box Functions. In *Advances in Neural Information Processing Systems 27*, pages 918–926, 2014.
- [12] A. J. Ijspeert, J. Nakanishi, H. Hoffmann, P. Pastor, and S. Schaal. Dynamical Movement Primitives: Learning Attractor Models for Motor Behaviors. *Neural Computation*, 25:1–46, 2013.
- [13] D. R. Jones, C. D. Perttunen, and B. E. Stuckman. Lipschitzian optimization without the Lipschitz constant. *Journal of Optimization Theory and Applications*, 79(1):157–181, Oct. 1993.
- [14] J. Kober, A. Wilhelm, E. Oztop, and J. Peters. Reinforcement learning to adjust parametrized motor primitives to new situations. *Autonomous Robots*, 33(4):361–379, 2012.
- [15] A. Krause and C. S. Ong. Contextual Gaussian Process Bandit Optimization. In *Advances in Neural Information Processing Systems 24*, pages 2447–2455, 2011.
- [16] O. B. Kroemer, R. Detry, J. Piater, and J. Peters. Combining active learning and reactive control for robot grasping. *Robot. Auton. Syst.*, 58(9):1105–1116, Sept. 2010.
- [17] A. G. Kupcsik, M. P. Deisenroth, J. Peters, and G. Neumann. Data-Efficient Generalization of Robot Skills with Contextual Policy Search. In *27th AAAI Conference on Artificial Intelligence*, June 2013.
- [18] D. Lizotte, T. Wang, M. Bowling, and D. Schuurmans. Automatic gait optimization with gaussian process regression. pages 944–949, 2007.

---

<sup>3</sup>More information are available at <http://robotik.dfki-bremen.de/en/research/projects/besman.html>.

- [19] A. Marco, P. Hennig, J. Bohg, S. Schaal, and S. Trimpe. Automatic LQR Tuning based on Gaussian Process Optimization: Early Experimental Results. In *Proceedings of the Second Machine Learning in Planning and Control of Robot Motion Workshop*, Hamburg, 2015. IROS.
- [20] J. H. Metzen, A. Fabisch, and J. Hansen. Bayesian Optimization for Contextual Policy Search. In *Proceedings of the Second Machine Learning in Planning and Control of Robot Motion Workshop*, Hamburg, 2015. IROS.
- [21] J. Peters, K. Mülling, and Y. Altun. Relative Entropy Policy Search. In *Proceedings of the Twenty-Fourth AAAI Conference on Artificial Intelligence*, Atlanta, Georgia, USA, 2010. AAAI Press.
- [22] C. Rasmussen and C. Williams. *Gaussian Processes for Machine Learning*. MIT Press, 2006.
- [23] P. Ruvolo and E. Eaton. Active Task Selection for Lifelong Machine Learning. In *Twenty-Seventh AAAI Conference on Artificial Intelligence*, June 2013.
- [24] B. Shahriari, Z. Wang, M. W. Hoffman, A. Bouchard-Côté, and N. de Freitas. An Entropy Search Portfolio for Bayesian Optimization. In *NIPS workshop on Bayesian optimization*, 2014.
- [25] J. Snoek, H. Larochelle, and R. P. Adams. Practical Bayesian Optimization of Machine Learning Algorithms. In *Advances in Neural Information Processing Systems 25*, pages 2951–2959, 2012.
- [26] N. Srinivas, A. Krause, and M. Seeger. Gaussian process optimization in the bandit setting: No regret and experimental design. In *Proceedings of the 27th International Conference on Machine Learning*, 2010.
- [27] K. Swersky, J. Snoek, and R. P. Adams. Multi-Task Bayesian Optimization. In C. J. C. Burges, L. Bottou, M. Welling, Z. Ghahramani, and K. Q. Weinberger, editors, *Advances in Neural Information Processing Systems 26*, pages 2004–2012. Curran Associates, Inc., 2013.
- [28] Z. Wang, M. Zoghi, F. Hutter, D. Matheson, and N. d. Freitas. Bayesian Optimization in High Dimensions via Random Embeddings. In *International Joint Conferences on Artificial Intelligence (IJCAI)*, 2013.