
Active Contextual Entropy Search

Jan Hendrik Metzen

Universität Bremen, 28359 Bremen, Germany
jhm@informatik.uni-bremen.de

Coauthor

Affiliation, Address
email

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 actual 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, which is well suited under these conditions. In this work, we extend entropy search, a particular kind of Bayesian optimization, such that it can be used for *active* contextual policy search, where the learning systems selects those tasks during training in which it expects to learn the most.

1 INTRODUCTION

Contextual policy search (CPS) is a popular means for multi-task reinforcement learning in robotic control [5]. CPS learns a hierarchical policy, in which the lower-level policy is often a domain-specific behavior representation such as dynamical movement primitives (DMPs) [10]. Learning takes place on the upper-level policy, which is typically a conditional probability density $\pi(\theta|s)$ that defines a distribution over the parameter vectors θ of the lower-level policy for a given context s . This context vector s encodes properties of the environment or the 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 [12] and contextual relative entropy search (C-REPS) [14, 17]. 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 [20]. Combining contextual policy search with pre-trained movement primitives¹ 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 [4, 15] and robot grasping [13] and for contextual policy search on a simulated robotic ball-throwing task [16].

In this work, we focus on the problem of actively selecting the task (context), in which the agent performs the next trial during learning. This constitutes an active learning approach, which is considered to be a prerequisite for lifelong learning [19]. 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 the relative (unknown) *difficulty* of a task: for instance, if a relatively small reward is

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

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? Recently, Fabisch et al. [7] presented one approach for estimating the task-difficulty explicitly, which allows defining heuristic intrinsic reward functions based on which a discounted multi-arm bandit can select the next task actively [6].

In this work, we follow a different approach: rather than explicitly addressing the incommensurability of rewards, we propose an information theoretic approach for active task selection which selects the task not based on rewards directly but rather based on the expected reduction in uncertainty about the optimal parameters conditional on the context. This approach is motivated by entropy search [9], which implements a similar extension for non-contextual Bayesian optimization.

2 BACKGROUND

Contextual Policy Search (CPS) refers to a model-free approach to reinforcement learning (RL), in which the (low-level) policy π_θ is parametrized by a vector θ . The choice of θ is governed by an upper-level policy $\pi_\omega(\theta)$. 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 π_ω such that the expected return J_ω 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 θ in context s . CPS is typically based on local search based approaches such as cost-regularized kernel regression [12] and contextual relative entropy search (C-REPS) [14, 17]. We refer to Deisenroth et al. [5] for a recent overview of (contextual) policy search approaches in the robotics domain.

Bayesian optimization for contextual policy search (BO-CPS) [16] is based on applying ideas from Bayesian optimization [2, 20] to contextual policy search. BO-CPS learns internally a model of the expected return $R(\theta, s)$ of a parameter vector θ in a context s . This model is learned by means of Gaussian process (GP) regression [18] 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 θ_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 the uncertainty $\sigma_{GP}[R(s, \theta)]$ of this estimate. Based on this information, the parameter vector for the given context is selected by maximizing an *acquisition function*. These acquisition functions allow 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 [16]. In contrast, the acquisition function GP-UCB, which defines the acquirability 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 κ controls the exploration-exploitation trade-off, can be applied to BO-CPS straightforwardly. BO-CPS selects parameters θ_i for a given fixed context s_i by performing an optimization over the parameter space using the global maximizer DIRECT [11] to find the approximate global maximum, followed by L-BFGS [3] to refine it.

Entropy search (ES) [9] is a recently proposed approach to probabilistic global optimization that mainly differs from Bayesian optimization in the choice of the acquisition function. 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 $p_{opt}(f, \theta)$, the probability that the global optimum (maximum or minimum, depending on the problem) of the unknown function f is at θ . ES estimates $p_{opt}(f)$ at finitely many points $\{\theta^c\}_{i=1}^{N_\theta}$ on a non-uniform grid that are selected heuristically and approximates p_{opt} at θ^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 θ^q and assuming N_y different outcomes $\{y^{(i)}\}$ sampled from the GP’s predictive distribution at θ^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* a query at θ^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 θ_q and context s_q of the next trial. Let $p_{max}^{(s)}(\theta)$ denote the probability distribution of the maximum expected return in context s and let the loss $\mathcal{L}^s(s^q, \theta^q) = \mathcal{L}(p_{max}^{(s)}[s^q, \theta^q]) - \mathcal{L}(p_{max}^{(s)})$ denote the expected change of relative entropy in context s after performing a trial in context s^q with parameter θ^q . A straightforward extension of ES to active learning in BO-CPS would be selecting $s^q, \theta^q = \arg \min_{(s^q, \theta^q)} \mathcal{L}^s(s^q, \theta^q)$. This, however, would not account for information gained by a query at (s^q, θ^q) about the optima in contexts $s \neq s^q$.

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 $\{s^c\}_{i=1}^{N_s}$ 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 $d_M(s_i^c, s^q) = \sqrt{(s_i^c - s^q)S^{-1}(s_i^c - s^q)}$ 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.

In the experiments, candidate points $\theta^c(s)$ are selected using Thompson sampling with $N_\theta = 20$. The number of test context is chosen to as $N_s = 100$ 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 θ_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 [8] as global 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 $[1, 2.5]m \times [-1, 1]m$ for the arm mounted at the origin of this coordinate system. 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 two meta-parameters: the execution time τ of the DMP, which determines how far the ball is thrown, and the final angle g_0 of the first joint, which determines the rotation of the arm around the z-axis.

The upper-level policy, which defines a distribution over meta-parameters conditioned on a context, is a Gaussian policy with the mean being an affine function of context, initial mean $(g_0, \tau)^T = (0, 1)^T$, initial variance $((0.1\pi)^2, 1^2)$, and parameter space $g_0 \in [-\frac{\pi}{2}, \frac{\pi}{2}]$ and $\tau \in [0.2, 5]$. All approaches use an anisotropic Matérn kernel for the GP surrogate model. GP-UCB’s exploration parameter is set to $\kappa = 5.0$. 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.

Figure 1 summarizes the main results of the empirical evaluation. The left graph shows the mean performance of the greedy policy, which was evaluated every 10 episodes on 16 test contexts. 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.

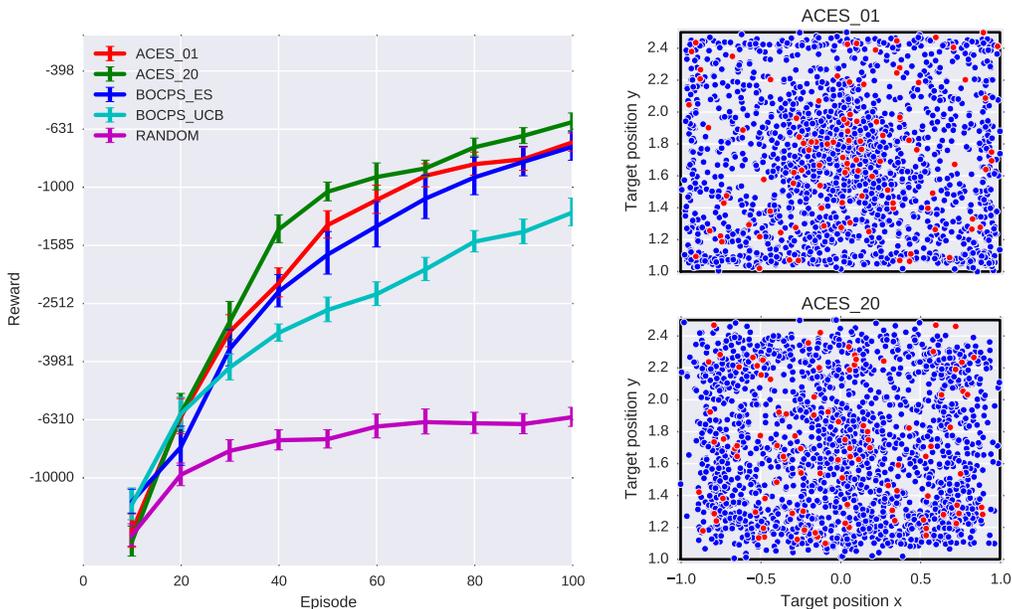


Figure 1: (Left) Learning curves on the simulated robot arm COMPI: the offline performance is evaluated each 10 episodes on 16 test contexts distributed equally in the target area. Shown are mean and its standard error over 20 independent runs. (Right) Scatter plot showing the sampled contexts for all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

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. An open question for future work is if the proposed approach can be scaled to higher dimensional problems. A combination with random embedding approaches such as REMBO [21] could be an interesting starting point for this.

Acknowledgements This work was 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] COMPI - compliant robot arm. <http://robotik.dfki-bremen.de/en/research/robot-systems/comp.html>.
- [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] 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.
- [4] R. Calandra, N. Gopalan, A. Seyfarth, J. Peters, and M. P. Deisenroth. Bayesian gait optimization for bipedal locomotion. In *Proceedings of Learning and Intelligent Optimization Conference (LION8)*, 2014.
- [5] 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.
- [6] A. Fabisch and J. H. Metzen. Active contextual policy search. *Journal of Machine Learning Research*, 15:3371–3399, 2014.
- [7] 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.
- [8] N. Hansen and A. Ostermeier. Completely derandomized self-adaptation in evolution strategies. *Evolutionary Computation*, 9:159–195, 2001.
- [9] P. Hennig and C. J. Schuler. Entropy Search for Information-Efficient Global Optimization. *JMLR*, 13:1809–1837, 2012.
- [10] 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.
- [11] 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.
- [12] J. Kober, A. Wilhelm, E. Oztog, and J. Peters. Reinforcement learning to adjust parametrized motor primitives to new situations. *Autonomous Robots*, 33(4):361–379, 2012.
- [13] 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.
- [14] 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.
- [15] D. Lizotte, T. Wang, M. Bowling, and D. Schuurmans. Automatic gait optimization with gaussian process regression. pages 944–949, 2007.
- [16] 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.
- [17] 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.
- [18] C. Rasmussen and C. Williams. *Gaussian Processes for Machine Learning*. MIT Press, 2006.
- [19] P. Ruvolo and E. Eaton. Active Task Selection for Lifelong Machine Learning. In *Twenty-Seventh AAAI Conference on Artificial Intelligence*, June 2013.
- [20] 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.
- [21] 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.