



Motivation

- Policy search \equiv inference (via stochastic control)
- Wingate et al. 2011: Bayesian policy search in deterministic domains
- van de Meent et al. 2016: Policy optimization in stochastic domains
- Bayesian** policy search in **stochastic** domains?

- David Wingate, Noah D Goodman, Daniel M Roy, Leslie P Kaelbling, and Joshua B Tenenbaum. Bayesian policy search with policy priors. IJCAI 2011
- Jan-Willem van de Meent, Brooks Paige, David Tolpin, and Frank Wood. Black-box policy search with probabilistic programs. AISTATS 2016

Stochastic Metropolis–Hastings

$$p_{common}(\mathbf{1}|\mathbb{E}(r|\theta)) = \exp(\mathbb{E}(r|\theta) - U_r) = \prod_r (\exp(r - U_r))^{p_S(r|\theta)} = \prod_r p_{common}(\mathbf{1}|r)^{p_S(r|\theta)}$$

$$p_{our}(\mathbf{1}|\mathbb{E}(r|\theta)) = \frac{\mathbb{E}(r|\theta) - L_r}{U_r - L_r} = \sum_r p_S(r|\theta) \left(\frac{r - L_r}{U_r - L_r} \right) = \sum_r p_S(r|\theta) p_{our}(\mathbf{1}|r) - \text{mixture!}$$

Model

$\tau \sim D_\tau$ — simulator trace

$\theta \sim D_\theta$

$\mathbf{1} \sim \text{Bernoulli} \left(\frac{\mathbb{E}_\tau[\mathcal{P}(\theta, \tau)] - L_r}{U_r - L_r} \right)$

Algorithm

- loop**
- $\tau \sim D_\tau$ (* always accept *)
- $\theta' \sim D_\theta$
- $u \sim \text{Uniform}(0, 1)$
- if** $\{u < \min \left(1, \frac{\mathcal{P}(\theta', \tau) - L_r}{\mathcal{P}(\theta, \tau) - L_r} \right)\}$ **then**
- $\theta \leftarrow \theta'$
- end if**
- Output θ
- end loop**

Probabilistic programs for policy search

- $S(\theta)$ — simulator, returns reward r
- θ — parameters of interest

Common conditioning

$$p_{common}(\mathbf{1}|r) = \exp(r - U_r)$$

U_r — an upper bound on r

Generative model

$\theta \sim D_\theta$

$\mathbf{1} \sim \text{Bernoulli}(\exp(\mathbb{E}[S(\theta)] - U_r))$

Our conditioning

$$p_{our}(\mathbf{1}|r) = \frac{r - L_r}{U_r - L_r}$$

U_r, L_r — upper and lower bounds on r

Generative model

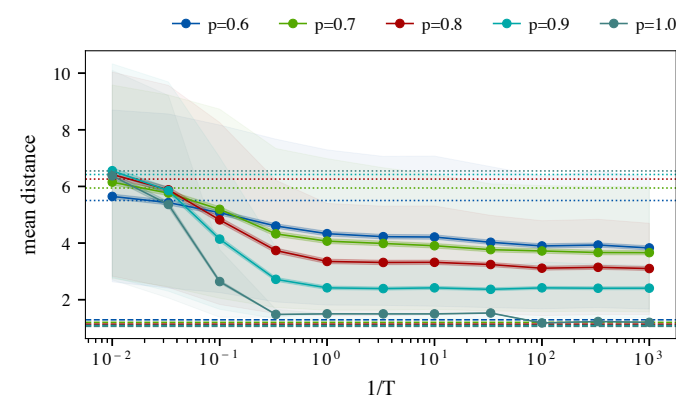
$\theta \sim D_\theta$

$\mathbf{1} \sim \text{Bernoulli} \left(\frac{\mathbb{E}[S(\theta)] - L_r}{U_r - L_r} \right)$

Allows flattening of nested inference!

Case studies

Canadian traveller



RockSample

