

## Introduction

- ▶ We propose an efficient policy search method for preference-based reinforcement learning
- ▶ Evolutionary Direct Policy Search (EDPS) [1]
- ▶ Assume a parametric policy space

$$\Pi = \{\pi_{\Theta} \mid \Theta \in \mathbb{R}^p\}$$

- ▶ Target function is the *expected total reward*

$$\rho_{\pi} = \mathbb{E}_{\mathbf{h} \sim P_{\pi}} [V(\mathbf{h})]$$

that is estimated based on **rollouts**

- ▶ Optimize it by using an *Evolution strategy* (ES), such as CMA-ES [2]
  - ▶ If the number of rollouts is too large, the learning process gets slow
  - ▶ If the number of rollouts is too small, the ranking over the offsprings is not reliable enough
- ▶ Adaptive control of the number of rollouts using **racing algorithms**

### Algorithm 1 EDPS ( $\mathcal{M}, \mu, \lambda, n_{\max}, \delta$ )

Initialization: select an initial parameter vector  $\Omega^{(0)}$  and an initial set of candidate solutions  $\Theta_1^{(0)}, \dots, \Theta_{\mu}^{(0)}$ ,  $\sigma^{(0)}$  is the identity permutation

$t = 0$

**repeat**

$t = t + 1$

**for**  $\ell = 1, \dots, \lambda$  **do**

$$\Theta_{\ell}^{(t)} \sim F(\Omega^{(t-1)}, \Theta_{\sigma^{(t-1)}(1)}^{(t-1)}, \dots, \Theta_{\sigma^{(t-1)}(\mu)}^{(t-1)})$$

**end for**

$$\sigma^{(t)} = \text{Racing}(\mathcal{M}, \pi_{\Theta_1^{(t)}}, \dots, \pi_{\Theta_{\lambda}^{(t)}}; \mu, n_{\max}, \delta)$$

$$\Omega^{(t)} = \text{Update}(\Omega^{(t-1)}, \Theta_{\sigma^{(t)}(1)}^{(t)}, \dots, \Theta_{\sigma^{(t)}(\mu)}^{(t)})$$

**until** Stopping criterion fulfilled

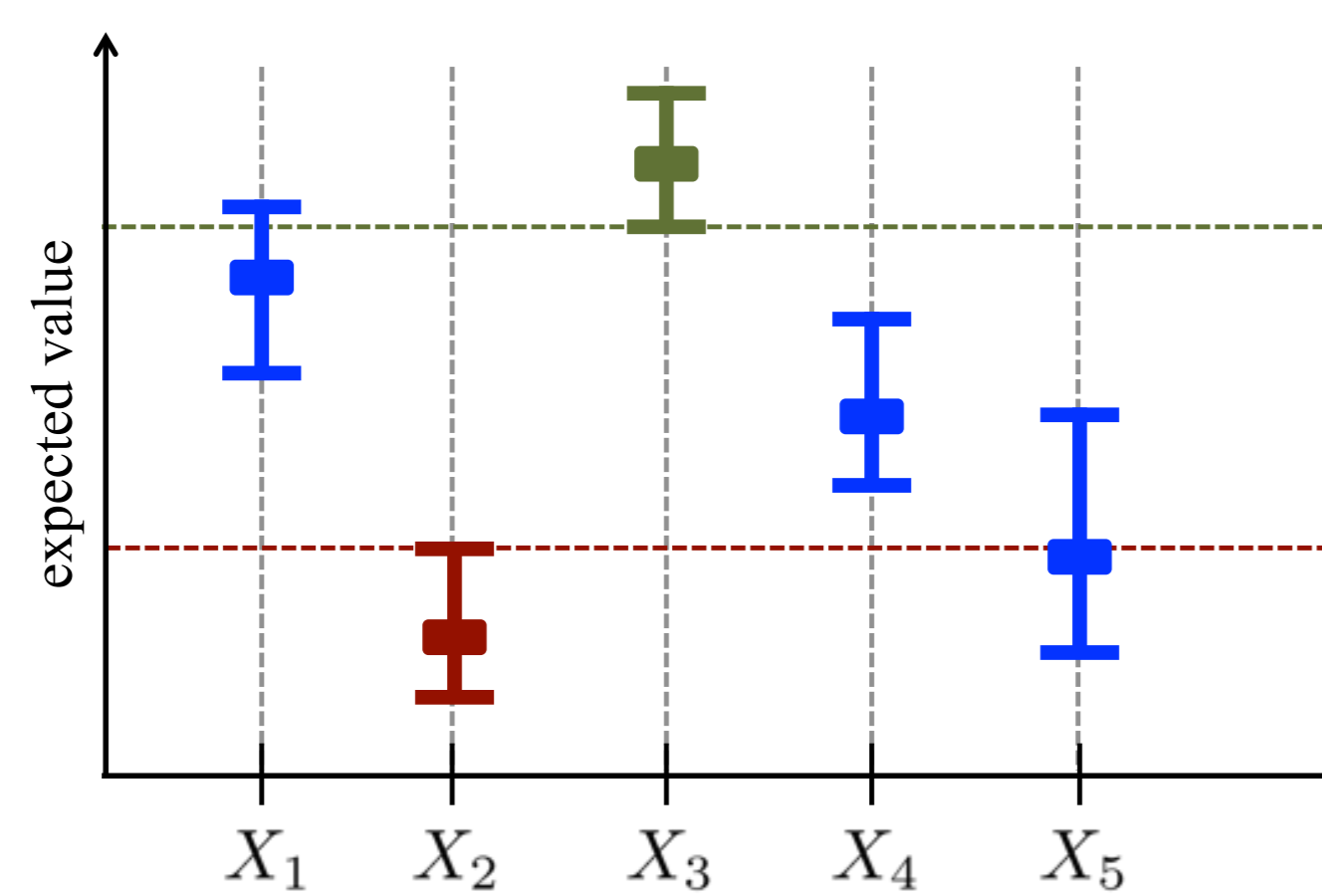
**Return**  $\pi_{\Theta_1^{(t)}}$

## Value-based Racing Algorithm [3]

- ▶  $X_1, \dots, X_K$  are random variables with unknown distribution functions  $P_{X_1}, \dots, P_{X_K}$  and finite expected values  $\mu_i = \int x dP_{X_i}(x)$
- ▶ We solve the following optimization task with high probability

$$\sum_{i \in I} \sum_{j \neq i} \mathbb{I}\{\mu_j < \mu_i\} \longrightarrow \max_{I \subseteq [K]: |I|=\kappa}$$

- ▶ Hoeffding bound:  $\mu_i \in \left[ \hat{\mu}_i - \sqrt{\frac{1}{2n_i} \log \frac{2}{\delta}}, \hat{\mu}_i + \sqrt{\frac{1}{2n_i} \log \frac{2}{\delta}} \right]$  with probability at least  $1 - \delta$



## Preference-based Reinforcement Learning [4,5]

- ▶ **Rollout**: generating a history  $\mathbf{h} \in \mathcal{H}^{(T)}$  by following a policy  $\pi$  for a given MDP  $\mathcal{M} = (\mathcal{S}, \mathcal{A}, P, r)$
- ▶ Assumption: preference relation  $\prec$  on  $\mathcal{H}^{(T)}$
- ▶ Prerequisite: "lifting" of the preference relation  $\prec$  on  $\mathcal{H}^{(T)}$  to a preference relation  $\ll$  on the space of policies  $\Pi$
- ▶ Each policy generates a distribution over the histories  $\mathcal{H}^{(T)}$
- ▶ We can associate policies with random variables  $X$
- ▶ Decision model ( $\ll$ ):

$$X \ll Y \text{ if and only if } P(Y \prec X) < P(X \prec Y)$$

- ▶ Preferential cycles:  $X_1 \ll X_2, X_2 \ll X_3, X_3 \ll X_1$

## Preference-based EDPS

- ▶ Evolution Strategies only need a **ranking** over the candidate solutions to update the parameters of  $F(\cdot)$  (distribution over the search space)
- ▶ **Let's race based on preferences!!**  $\Rightarrow$  Preference-based EDPS
- ▶ Resolving the preferential cycles by using *Copeland relation*:  
 $X_i \ll_C X_j \Leftrightarrow d_i < d_j$ , where  $d_i = \#\{k : X_k \ll X_i, X_k \in \mathcal{X}\}$
- ▶ We solve the following optimization task with high probability

$$\sum_{i \in I} \sum_{j \neq i} \mathbb{I}\{X_j \ll X_i\} \longrightarrow \max_{I \subseteq [K]: |I|=\kappa}$$

- ▶ We need an *efficient estimator* of  $S(X_i, X_j) = P(X_i \prec X_j)$
- ▶ A two-sample U-statistic called the *Mann-Whitney U-statistic*

$$\hat{s}_{i,j} = \hat{S}(X_i, X_j) = \frac{1}{n^2} \sum_{\ell=1}^n \sum_{\ell'=1}^n [\mathbb{I}\{x_i^{\ell} \prec x_j^{\ell'}\} + \frac{1}{2} [\mathbb{I}\{x_i^{\ell} \sim x_j^{\ell'}\} + \mathbb{I}\{x_i^{\ell} \perp x_j^{\ell'}\}]]$$

where  $X_i = \{x_i^1, \dots, x_i^n\} \sim X_i$  and  $X_j = \{x_j^1, \dots, x_j^n\} \sim X_j$

- ▶ **Hoeffding, 1963, §5b**: For any  $\epsilon > 0$ , using the notations introduced above,

$$P\left(\left|\hat{S}(X, Y) - S(X, Y)\right| \geq \epsilon\right) \leq 2 \exp(-2n\epsilon^2).$$

## Preference-based Racing Algorithm

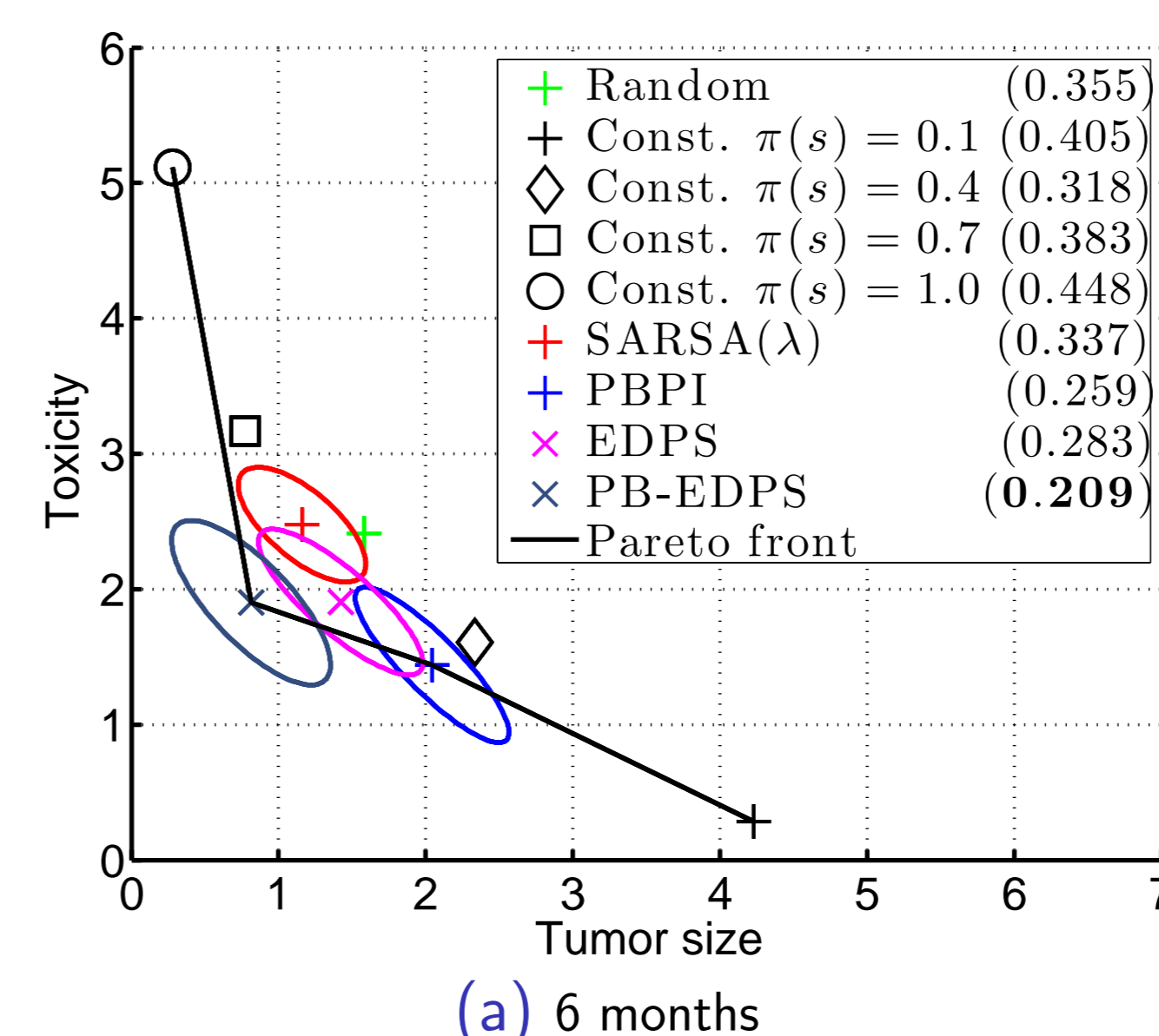
1. Input:  $X_1, \dots, X_K, \kappa, n_{\max}, \delta$
2. Iteratively sample  $X_1, \dots, X_K$
3. Calculate  $\hat{s}_{i,j}$  for all  $1 \leq i, j \leq K$
4. and their confidence intervals as  $[\hat{s}_{i,j} - c_{i,j}, \hat{s}_{i,j} + c_{i,j}]$  where

$$c_{i,j} = \sqrt{\frac{1}{2n} \log \frac{2K^2 n_{\max}}{\delta}}$$

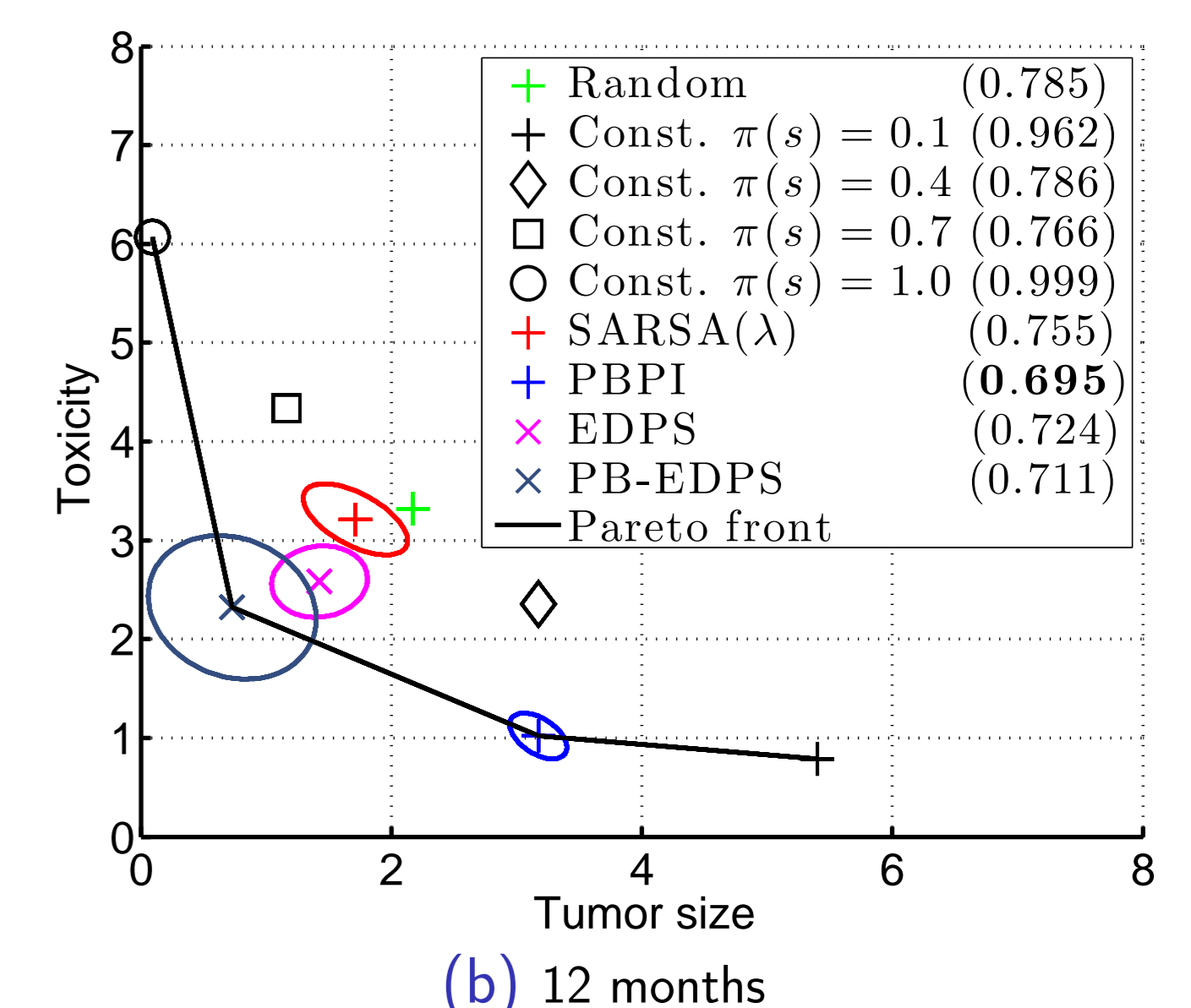
5. When can we stop sampling an option?
  - ▶ Number of options that are beaten by  $i$  so far:  $z_i = |\{j : u_{i,j} < 1/2, j \neq i\}|$
  - ▶ Number of options that beat  $i$  so far:  $o_i = |\{j : \ell_{i,j} > 1/2, j \neq i\}|$
  - ▶  $C = \{i : K - \kappa < |\{j : K - z_j < o_i\}|\}$
  - ▶  $D = \{i : \kappa < |\{j : K - o_j < z_i\}|\}$
6. If  $(i, j \in C \cup D) \vee (1/2 \notin [\ell_{i,j}, u_{i,j}])$  then do not update  $\hat{s}_{i,j}$  any more

## Medical Experiments

- ▶ Medical treatment design for cancer clinical trials
- ▶ State  $s = (\mathcal{S}, X)$  describes the health condition of the patient:  $\mathcal{S}$  is the tumor size and  $X$  the level of toxicity
- ▶ Action is the dosage level  $a \in [0, 1]$
- ▶ A history  $\mathbf{h}$  represents a treatment of a virtual patient
  1.  $\mathbf{h}' \preceq \mathbf{h}$  if the patient survives in  $\mathbf{h}$  but not in  $\mathbf{h}'$ , and both histories are incomparable ( $\mathbf{h}' \perp \mathbf{h}$ ) if the patient does neither survive in  $\mathbf{h}'$  nor in  $\mathbf{h}$ .
  2. Otherwise, preference depends on the worst wellness of the patient and the final tumor size:  $\mathbf{h}' \preceq \mathbf{h}$  if (and only if)  $C_X \leq C'_X$  and  $C_S \leq C'_S$  where  $C_X$  and  $C'_X$  denote the *maximal* toxicity during the whole treatment
  3. Pareto dominance



(a) 6 months



(b) 12 months

Illustration of patient status under different treatment policies. On the x-axis is the tumor size after 6 (a) and 12 (b) months, on the y-axis the highest toxicity during the treatment. The death rates are shown in parentheses at the upper right corner.

- [1] Heidrich-Meisner, V., Igel, C.: Hoeffding and Bernstein races for selecting policies in evolutionary direct policy search. ICML, pp. 401–408 (2009)
- [2] Hansen, N., Kern, S.: Evaluating the CMA evolution strategy on multimodal test functions. PPSN VIII, pp. 282–291 (2004)
- [3] Maron, O., Moore, A.: Hoeffding races: accelerating model selection search for classification and function approximation. NIPS, pp. 59–66 (1994)
- [4] Fürnkranz, J., Hüllermeier, E., Cheng, W., Park, S.: Preference-based reinforcement learning: a formal framework and a policy iteration algorithm. Machine Learning **89**(1-2), 123–156 (2012)
- [5] Akrou, R., Schoenauer, M., Sebag, M.: Preference-based policy learning. ECML/PKDD, pp. 12–27 (2011)