Information Theoretic Regret Bounds for Online Nonlinear Control

<u>Model-based Reinforcement Learning</u> - The Kernelized Nonlinear Regulator Model

This work studies the following nonlinear control problem:

 $x_{h+1} = f(x_h, u_h) + \epsilon$, where $\epsilon \sim N(0, \sigma^2 I)$, (1)

where $h \in \{0, 1, \dots, H-1\}$; the state $x_h \in \mathbb{R}^{d_{\mathcal{X}}}$; the control $u_h \in \mathcal{U}$ where \mathcal{U} may be an arbitrary set; $f : \mathcal{X} \times \mathcal{U} \to \mathcal{X}$.

Assume that f lives in a known RKHS.

• Equivalently: $f(x, u) = W^* \phi(x, u)$. W^* : Unknown linear mapping, $\phi : \mathcal{X} \times \mathcal{U} \to \mathcal{V}$: Known function. \mathcal{V} : Either finite or countably infinite dimensional Euclidean vector space.

▶ Main result: a regret bound $O(\tilde{d}_{\text{effective}}\sqrt{HT})$. (*T*:#episodes) Proof techniques: a novel self-bounding regret lemma.

Experiments

Benchmark Tasks with Random Features

* Benchmark tasks including MuJoCo environments from OpenAI Gym

* We use Random Fourier Features

* It is observed that the Thompson Sampling variant of our proposed algorithm with RFFs quickly increased reward in early stages, indicating low sample complexities empirically

* Our algorithm consistently performs well on simple continuous control tasks





Sham Kakade, Akshay Krishnamurthy,

Kendall Lowrey, Motoya Ohnishi, Wen Sun





<u>Theory & Algorithm</u>

Lower Confidence-based Continuous Control (LC³)

- 1: Initialize BALL⁰ to be any set containing W^*
- 2: for t = 0 ... T do
- $\pi^{t} = \operatorname{argmin}_{\pi \in \Pi} \operatorname{min}_{W \in \operatorname{BALL}^{t}} J^{\pi}(x_{0}; c^{t}, W)$ Execute π^{t} to sample $\tau^{t} := \{x_{h}^{t}, u_{h}^{t}, c_{h}^{t}, x_{h+1}^{t}\}_{h=0}^{H-1}$
- Update BALL^t

6: end for

- Π: Policy class
- \blacktriangleright BALL^t: Confidence ball of weight matrices W at tth episode
- \triangleright c^{t} : Immediate cost function at *t*th episode
- \blacktriangleright $J^{\pi}(x_0; c^t, W)$: Expected total cost of π under the dynamics $W\phi(x, u) + \epsilon$

Assumption 1

We have access to an oracle that implements Line 3 of Algorithm – One may use some effective heuristics such as DDP, MPPI etc.

Assumption 2

For every episode, the cost function is non-negative and the realized cumulative cost has uniformly bounded second moments, i.e.

$$\sup_{\pi \in \Pi} \mathbb{E} \left[\left(\sum_{h=0}^{H-1} c^t(x_h, u_h) \right)^2 \, \middle| \, x_0, \pi \right] \le V_{\max}.$$

Armhand Robotics System

* 33 degree of freedom robotic arm and hand system tasked with picking up an spherical object

* Learning model dynamics for the real world applications such as robotics requires sufficiently complex features

* We test our method with features created by six ensemble of MuJoCo models

* Each element of the ensemble is unable to complete the task in isolation







Cornell University

We do NOT require bounded state space, bounded cost function or bounded feature map

1. We develop a stopping time martingale to handle the unbounded nature of the (realized) cumulative costs

2. We develop a novel way to handle Gaussian smoothing through the chisquared distance function between two distributions

3. We utilize methods developed for the analysis of linear bandits and Gaussian process bandits (e.g. maximum information gain)

We address multi-step extension to RL settings

We prove a "self-bounding" regret bound that relates the instantaneous regret to the second moment of the stochastic process

Self-Bounding, Simulation Lemma

For any policy, model parameterization, and non-negative cost, and for any initial state, we have:

$$J^{\pi}(x_{0}; c, W^{\star}) - J^{\pi}(x_{0}; c, W)$$

$$\leq \sqrt{HV^{\pi}(x_{0}; c, W^{\star})} \sqrt{\mathbb{E}\left[\sum_{h=0}^{H-1} \min\left\{\frac{1}{\sigma^{2}} \|(W^{\star} - W) \phi(x_{h}, u_{h})\|_{2}^{2}, 1\right\}\right]}$$

Maze

* We construct a continuous maze environment, where an agent plans for continuous actions with MPPI

* Exploration is necessary to find the goal

* We compare random walk exploration to our method, with and without posterior sampling



