

# A Pontryagin Perspective on Reinforcement Learning

*or: “Open-Loop Reinforcement Learning”*

Onno Eberhard<sup>1,2</sup> · Claire Vernade<sup>2</sup> · Michael Muehlebach<sup>1</sup>

<sup>1</sup>Max Planck Institute for Intelligent Systems

<sup>2</sup>University of Tübingen



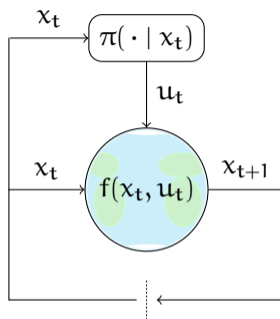
EBERHARD KARLS  
UNIVERSITÄT  
TÜBINGEN



imprs-is

April 19, 2024

## Recap: reinforcement learning

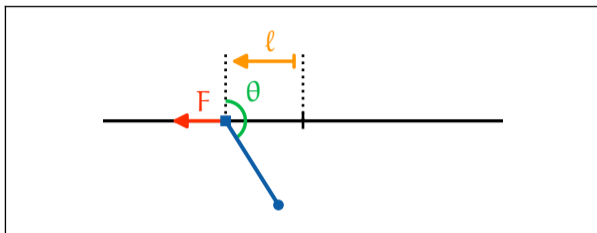


- Goal: learn a policy  $\pi$  that maximizes the sum of rewards

$$\pi^* = \arg \max_{\pi: \mathcal{X} \rightarrow \Delta_{\mathcal{U}}} \mathbb{E}_{\pi} \left[ \sum_{t=0}^{T-1} r(x_t, u_t) + r_T(x_T) \right]$$

- Challenge: the system is unknown, only simulations possible

## Example: inverted pendulum swing-up

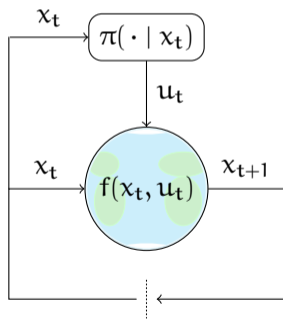


- ▶ State:  $x = (l, \dot{l}, \theta, \dot{\theta})$
- ▶ Action:  $u = F$  (external force)
- ▶ Rewards: positive when upright and at rest, penalty for large force
- ▶ Classical benchmark in control theory

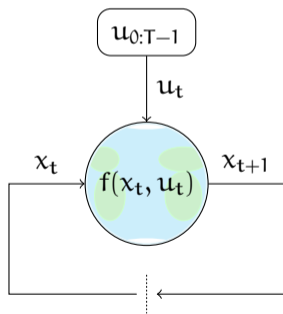
# Open-loop control

- ▶ Pendulum is “easily” solvable with deep RL, e.g. soft actor-critic (SAC)
  - ▶ SAC uses multiple neural networks,  $> 100,000$  parameters in total
  - ▶ Overkill! Why not just learn the actions that are necessary?

Closed-loop control



Open-loop control



# Pontryagin's principle

- ▶ Open-loop optimal control:

$$\mathbf{u}_{0:T-1}^* = \arg \max_{\mathbf{u}_{0:T-1} \in \mathcal{U}^T} \underbrace{\sum_{t=0}^{T-1} r(\mathbf{x}_t, \mathbf{u}_t) + r_T(\mathbf{x}_T)}_{J(\mathbf{u}_{0:T-1})} \quad \text{s.t.} \quad \mathbf{x}_{t+1} = f(\mathbf{x}_t, \mathbf{u}_t)$$

- ▶ Gradient-based optimization: how to compute  $\nabla_{\mathbf{u}_t} J(\mathbf{u}_{0:T-1})$ ?

## Pontryagin's principle for computing $\nabla_{\mathbf{u}_t} J(\mathbf{u}_{0:T-1})$

1. Forward pass:  $\mathbf{x}_{t+1} = f(\mathbf{x}_t, \mathbf{u}_t)$ , where  $\mathbf{x}_0$  is given
2. Backward pass:  $\lambda_t = \nabla_{\mathbf{x}} r(\mathbf{x}_t, \mathbf{u}_t) + \nabla_{\mathbf{x}} f(\mathbf{x}_t, \mathbf{u}_t) \lambda_{t+1}$ , where  $\lambda_T = \nabla r_T(\mathbf{x}_T)$
3. Gradient:  $\nabla_{\mathbf{u}_t} J(\mathbf{u}_{0:T-1}) = \nabla_{\mathbf{u}} r(\mathbf{x}_t, \mathbf{u}_t) + \nabla_{\mathbf{u}} f(\mathbf{x}_t, \mathbf{u}_t) \lambda_{t+1}$

# Open-loop reinforcement learning

## Theorem (informal)

Replace  $\nabla_x f_t$  and  $\nabla_u f_t$  in Pontryagin's equations by estimates  $A_t$  and  $B_t$  with sufficiently small errors  $\|\nabla_x f_t - A_t\|$  and  $\|\nabla_u f_t - B_t\|$  to get an approximate gradient  $g \simeq \nabla J(\mathbf{u}_{0:T-1})$ . Gradient ascent on  $g$  produces iterates  $\mathbf{u}_{0:T-1}^{(0)}, \dots, \mathbf{u}_{0:T-1}^{(N-1)}$  that satisfy

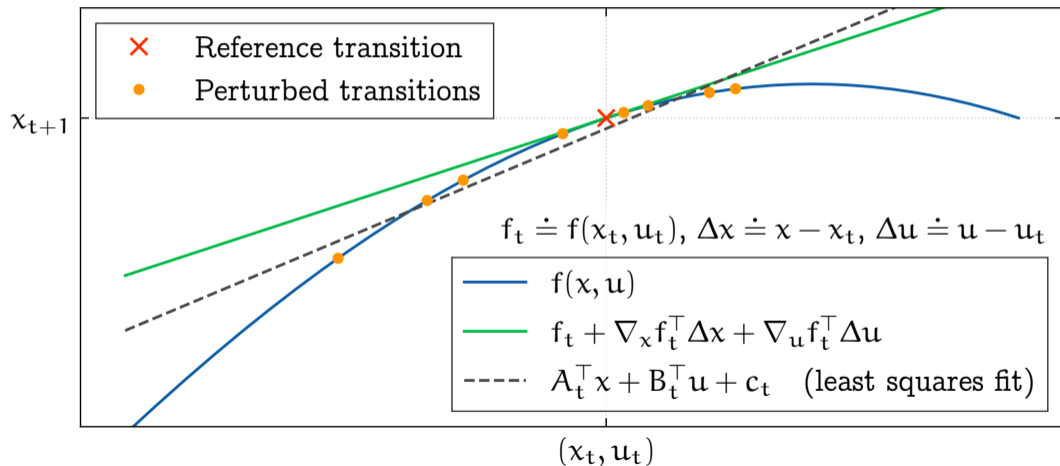
$$\frac{1}{N} \sum_{k=0}^{N-1} \|\nabla_{\mathbf{u}_t} J(\mathbf{u}_{0:T-1}^{(k)})\|^2 \leq \frac{J^* - J(\mathbf{u}_{0:T-1}^{(0)})}{\alpha \eta N}.$$

How should we choose  $A_t$  and  $B_t$ ?

- ▶ Model-based open-loop RL: given a model  $\tilde{f}$ , use  $A_t \doteq \nabla_x \tilde{f}_t$  and  $B_t \doteq \nabla_u \tilde{f}_t$
- ▶ Model-free open-loop RL: estimate  $\nabla_x f_t$  and  $\nabla_u f_t$  directly

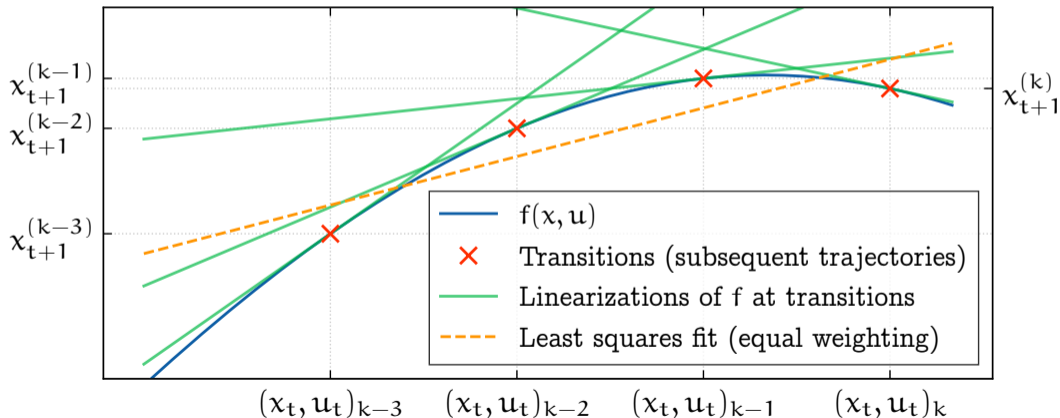
# Model-free open-loop RL

- The Jacobians  $\nabla_x f_t$  and  $\nabla_u f_t$  measure how  $x_{t+1}$  changes if  $(x_t, u_t)$  is perturbed

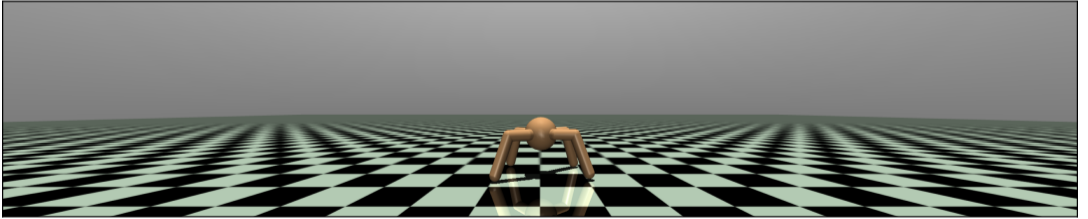
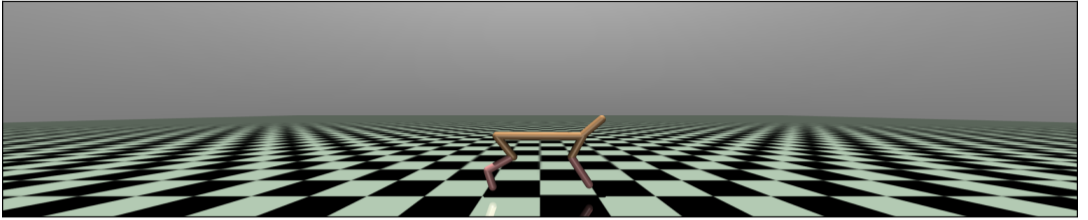


# Off-trajectory open-loop RL

- ▶ Subsequent trajectories are similar, no need to throw all data away!
- ▶ Algorithm: recursive least squares with forgetting







*Fin.*