



Peg Insertion with Policy Search

Caris Moses
6.231 Final Project

Motivation

- Policy search provide a **generalizable method** to learning parameters for policies.



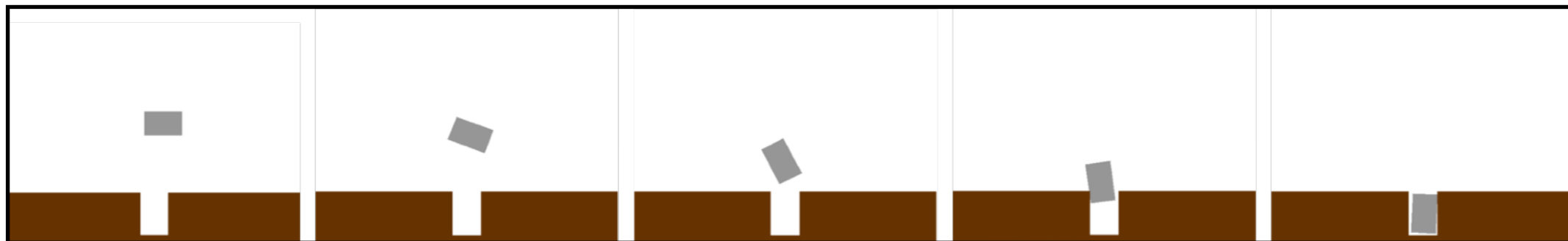
Motivation

- Policy search provide a **generalizable method** to learning parameters for policies.
- The methods considered in this talk are **model-free** and therefore require no knowledge of system dynamics.



Motivation

- Policy search provide a **generalizable method** to learning parameters for polices.
- The methods considered in this talk are **model-free** and therefore require no knowledge of system dynamics.
- Peg insertion is a know **difficult task** for stochastic robot processes.



Policy Parameterizations

$$\pi(a|s;\theta) = ?$$



Policy Parameterizations

Dynamic Movement Primitives

- Attractor Dynamics (PD controller) for 1D state

$$\pi(a|s, v_s) = -K(s - g_s) - Dv_s$$

$$K = 7$$
$$D = .5$$



Policy Parameterizations

Dynamic Movement Primitives^[2]

- Attractor Dynamics (PD controller with parameterized forcing term) for 1D state

$$\pi(a|s, v_s) = -K(s - g_s) - Dv_s + \mathbf{b}^T \boldsymbol{\theta}$$

Policy Parameterizations

Dynamic Movement Primitives^[2]

- Attractor Dynamics (PD controller with parameterized forcing term) for 1D state

$$\pi(a|s, v_s) = -K(s - g_s) - Dv_s + \mathbf{b}^T \boldsymbol{\theta}$$

- Canonical System Dynamics

$$\dot{z}_t = -\alpha z_t$$

$$z_0 = 1$$

Policy Parameterizations

Dynamic Movement Primitives[2]

- Attractor Dynamics (PD controller with parameterized forcing term) for 1D state

$$\pi(a|s, v_s) = -K(s - g_s) - Dv_s + \mathbf{b}^T \boldsymbol{\theta}$$

- Canonical System Dynamics

$$\dot{z}_t = -\alpha z_t \quad z_0 = 1$$

- p Radial Basis Function Kernels

$$[\mathbf{b}_t]_j = \frac{\psi_j}{\sum_{i=1}^p \psi_i} z_t |s_0 - g_s|$$
$$\psi_j = \exp(-h_j(z_t - c_j)^2)$$

Hyperparameters:

$$K > 0$$

$$D > 0$$

$$a > 0$$

$$g_s$$

$$h_j$$

$$c_j \in [0, 1]$$



Policy Parameterizations

Dynamic Movement Primitives

- Example DMP

$$K = 0$$

$$D = 0$$

$$c_0 = .95 \quad h_0 = 100$$

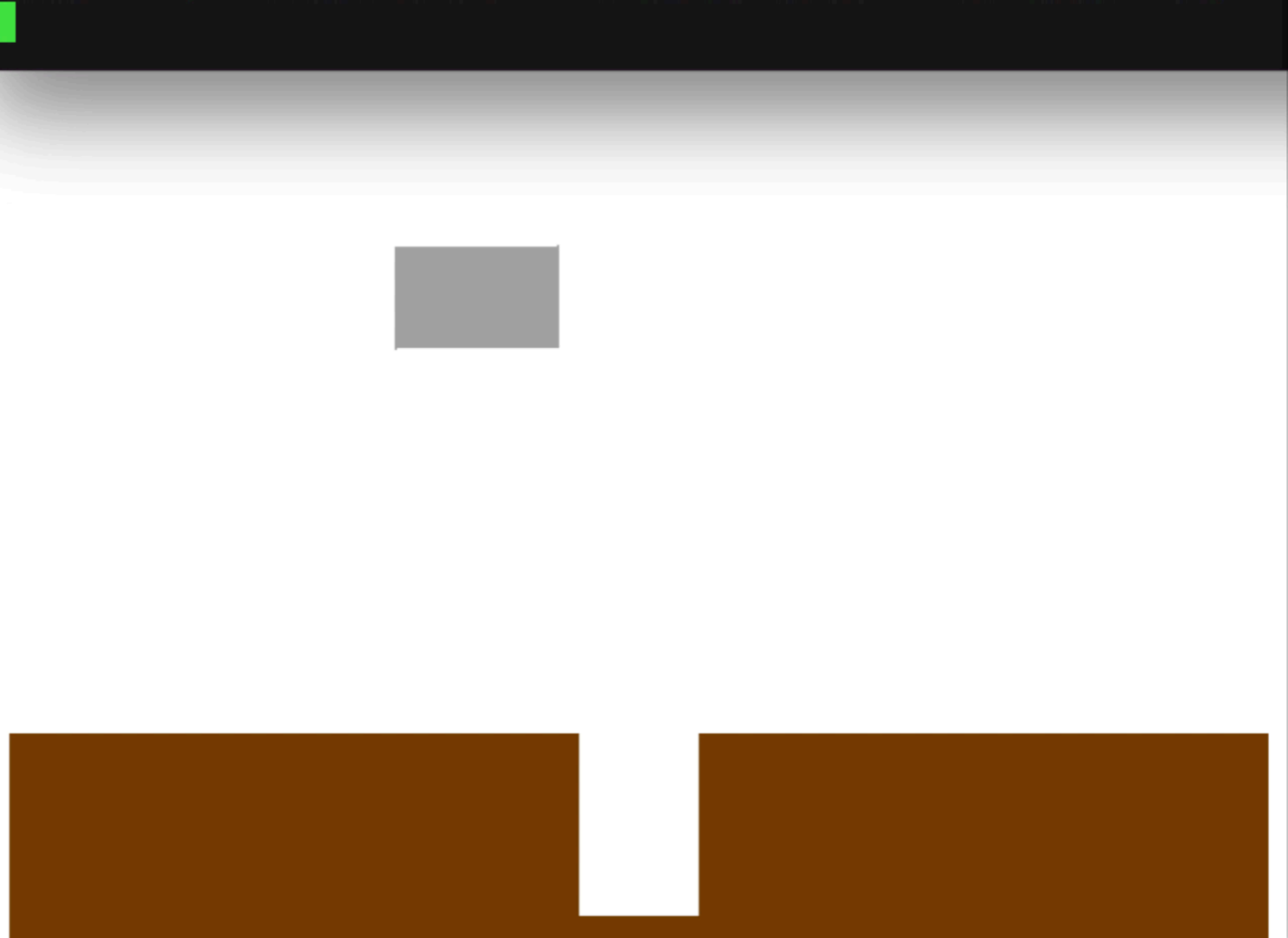
$$c_1 = .6 \quad h_1 = 100$$

$$\theta_x = [3, 0]$$

$$\theta_y = [0, -5]$$

$$\theta_\theta = [2, 0]$$

```
FORCE = [ 4.04625452e+01 -2.77010874e-05 2.69750302e+00]
```



Policy Parameterizations

Dynamic Movement Primitives

- Deterministic Policy

$$\pi(a|s, v_s; \boldsymbol{\theta}) = -K(s - g_s) - Dv_s + \mathbf{b}^T \boldsymbol{\theta}$$

Policy Parameterizations

Stochastic DMPs

- Deterministic Policy

$$\pi(a|s, v_s; \boldsymbol{\theta}) = -K(s - g_s) - Dv_s + \mathbf{b}^T \boldsymbol{\theta}$$

- Stochastic Policy $\epsilon \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma})$

Policy Parameterizations

Stochastic DMPs

- Deterministic Policy

$$\pi(a|s, v_s; \boldsymbol{\theta}) = -K(s - g_s) - Dv_s + \mathbf{b}^T \boldsymbol{\theta}$$

- Stochastic Policy $\epsilon \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma})$

- Parameter Space Exploration

$$\pi(a|s, v_s; \boldsymbol{\theta}) = \mathcal{N}(-K(s - g_s) - Dv_s + \mathbf{b}^T \boldsymbol{\theta}, \mathbf{b}^T \boldsymbol{\Sigma} \mathbf{b})$$

Policy Parameterizations

Stochastic DMPs

- Deterministic Policy

$$\pi(a|s, v_s; \boldsymbol{\theta}) = -K(s - g_s) - Dv_s + \mathbf{b}^T \boldsymbol{\theta}$$

- Stochastic Policy $\epsilon \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma})$

- Parameter Space Exploration

$$\pi(a|s, v_s; \boldsymbol{\theta}) = \mathcal{N}(-K(s - g_s) - Dv_s + \mathbf{b}^T \boldsymbol{\theta}, \mathbf{b}^T \boldsymbol{\Sigma} \mathbf{b})$$

- Action Space Exploration

$$\pi(a|s, v_s; \boldsymbol{\theta}) = \mathcal{N}(-K(s - g) - Dv_s + \mathbf{b}^T \boldsymbol{\theta}, \boldsymbol{\Sigma})$$

Policy Parameterizations

Linear Stochastic Policy

- Action Space Exploration

$$\pi(\mathbf{F}|\mathbf{s}; \boldsymbol{\theta}) = \boldsymbol{\theta} \mathbf{s} + \boldsymbol{\epsilon}$$

$$\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma})$$

Policy Search Methods

How do we learn θ ?



Policy Search Methods

Gradient Methods

- Value Function

$$J(\boldsymbol{\theta}) = \mathbf{E}_{\tau}[C(\tau)]$$

- Gradient Method Parameter Update

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \eta \nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta})$$

Policy Search Methods

Gradient Methods: Finite Differences^[1]

- Value Function

$$J(\boldsymbol{\theta}) = \mathbf{E}_{\tau}[C(\tau)]$$

- Gradient Method Parameter Update

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \eta \nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta})$$

- Finite Differences Gradient Estimate

$$\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) \approx \frac{J(\boldsymbol{\theta} + \boldsymbol{\delta}) - J(\boldsymbol{\theta} - \boldsymbol{\delta})}{2\boldsymbol{\delta}}$$

Policy Search Methods

Gradient Methods: REINFORCE[1]

- Gradient Estimate

$$\begin{aligned} J(\boldsymbol{\theta}) &= \mathbf{E}_{\tau} [C(\tau)] \\ &= \int_{\tau} C(\tau) p_{\boldsymbol{\theta}}(\tau) d\tau \end{aligned}$$

Policy Search Methods

Gradient Methods: REINFORCE[1]

- Gradient Estimate

$$\begin{aligned} J(\boldsymbol{\theta}) &= \mathbf{E}_{\tau} [C(\tau)] \\ &= \int_{\tau} C(\tau) p_{\boldsymbol{\theta}}(\tau) d\tau \end{aligned}$$

$$\begin{aligned} \nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) &= \int_{\tau} C(\tau) \nabla_{\boldsymbol{\theta}} p_{\boldsymbol{\theta}}(\tau) d\tau \\ &= \int_{\tau} C(\tau) \nabla_{\boldsymbol{\theta}} \log p_{\boldsymbol{\theta}}(\tau) p_{\boldsymbol{\theta}}(\tau) d\tau \\ &= \mathbf{E}_{\tau} [C(\tau) \nabla_{\boldsymbol{\theta}} \log p_{\boldsymbol{\theta}}(\tau)] \end{aligned}$$

Policy Search Methods

Gradient Methods: REINFORCE[1]

- Gradient Estimate

$$\begin{aligned} J(\boldsymbol{\theta}) &= \mathbf{E}_{\tau} [C(\tau)] \\ &= \int_{\tau} C(\tau) p_{\boldsymbol{\theta}}(\tau) d\tau \end{aligned}$$

$$\begin{aligned} \nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) &= \int_{\tau} C(\tau) \nabla_{\boldsymbol{\theta}} p_{\boldsymbol{\theta}}(\tau) d\tau \\ &= \int_{\tau} C(\tau) \nabla_{\boldsymbol{\theta}} \log p_{\boldsymbol{\theta}}(\tau) p_{\boldsymbol{\theta}}(\tau) d\tau \\ &= \mathbf{E}_{\tau} [C(\tau) \nabla_{\boldsymbol{\theta}} \log p_{\boldsymbol{\theta}}(\tau)] \end{aligned}$$

$$= \mathbf{E}_{\tau} [C(\tau) \nabla_{\boldsymbol{\theta}} \log \pi(a|s, v_s; \boldsymbol{\theta})]$$

Policy Search Methods

Gradient Methods: REINFORCE[1]

- Gradient Estimate

$$\begin{aligned} J(\boldsymbol{\theta}) &= \mathbf{E}_{\tau} [C(\tau)] \\ &= \int_{\tau} C(\tau) p_{\boldsymbol{\theta}}(\tau) d\tau \end{aligned}$$

$$\begin{aligned} \nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) &= \int_{\tau} C(\tau) \nabla_{\boldsymbol{\theta}} p_{\boldsymbol{\theta}}(\tau) d\tau \\ &= \int_{\tau} C(\tau) \nabla_{\boldsymbol{\theta}} \log p_{\boldsymbol{\theta}}(\tau) p_{\boldsymbol{\theta}}(\tau) d\tau \\ &= \mathbf{E}_{\tau} [C(\tau) \nabla_{\boldsymbol{\theta}} \log p_{\boldsymbol{\theta}}(\tau)] \end{aligned}$$

$$= \mathbf{E}_{\tau} [C(\tau) \nabla_{\boldsymbol{\theta}} \log \pi(a|s, v_s; \boldsymbol{\theta})]$$

- Only works with **stochastic policies**
- High variance can be mitigated by using **optimal baseline**^[3]

Policy Search Methods

Policy Improvement with Path Integrals^[2]

- (Discrete) Path Integral — path cost with penalty term for unlikely trajectories

$$S(\tau_t) = \sum_{t'=t}^T q_{t'}(s_{t'}, v_{s,t'}) + a_{t'}^T R a_{t'} - \log p(\tau_t | s_{t'}, v_{s,t'})$$

Constant: R (weight on the control cost)

Policy Search Methods

Policy Improvement with Path Integrals^[2]

- Softmax representing the probability of a trajectory given a set of trajectories and their path integrals

$$P(\tau_t) = \frac{\exp(-\frac{1}{\lambda} S(\tau_t))}{\sum_{\tau} \exp(-\frac{1}{\lambda} S(\tau_t))}$$

Constant: λ

Policy Search Methods

Policy Improvement with Path Integrals^[2]

- Softmax representing the probability of a trajectory given a set of trajectories and their path integrals

$$P(\tau_t) = \frac{\exp(-\frac{1}{\lambda} S(\tau_t))}{\sum_{\tau} \exp(-\frac{1}{\lambda} S(\tau_t))}$$

Constant: λ

- Parameter Update

$$\delta\boldsymbol{\theta}_t = \sum_{\tau} P(\tau_t) \mathbf{M}_t \epsilon_t$$

$$\mathbf{M}_t = \frac{R^{-1} \mathbf{b}_t \mathbf{b}_t^T}{\mathbf{b}_t^T R^{-1} \mathbf{b}_t}$$

$$[\delta\boldsymbol{\theta}]_j = \frac{\sum_{t=0}^T (T-t) \psi_{j,t} [\delta\boldsymbol{\theta}_t]_j}{(T-t) \psi_{j,t}}$$

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \delta\boldsymbol{\theta}$$

Policy Search Methods

Policy Improvement with Path Integrals^[2]

- Softmax representing the probability of a trajectory given a set of trajectories and their path integrals

$$P(\tau_t) = \frac{\exp(-\frac{1}{\lambda} S(\tau_t))}{\sum_{\tau} \exp(-\frac{1}{\lambda} S(\tau_t))}$$

- Parameter Update

$$\delta\theta_t = \sum_{\tau} P(\tau_t) \mathbf{M}_t \epsilon_t$$

$$[\delta\theta]_j = \frac{\sum_{t=0}^T (T-t) \psi_{j,t} [\delta\theta_t]_j}{(T-t) \psi_{j,t}}$$

$$\theta \leftarrow \theta + \delta\theta$$

- Referred to as **probability weighted averaging**

- Paths with high cost have low weight/probability and vice versa



Results



Problem Description

Stochastic DMP with 2D θ

- Stochastic DMP - each dimension of the state has its own controller

$$\pi(F_x|x, v_x) = -K(x - g_x) - Dv_x$$

$$\pi(F_y|y, v_y; \theta_y) = \mathcal{N}(-K(y - g_y) - Dv_y + \mathbf{b}^T \theta_y, \mathbf{b}^T \Sigma \mathbf{b})$$

$$\pi(F_\theta|\theta, v_\theta; \theta_\theta) = \mathcal{N}(-K(\theta - g_\theta) - Dv_\theta + \mathbf{b}^T \theta_\theta, \mathbf{b}^T \Sigma \mathbf{b})$$

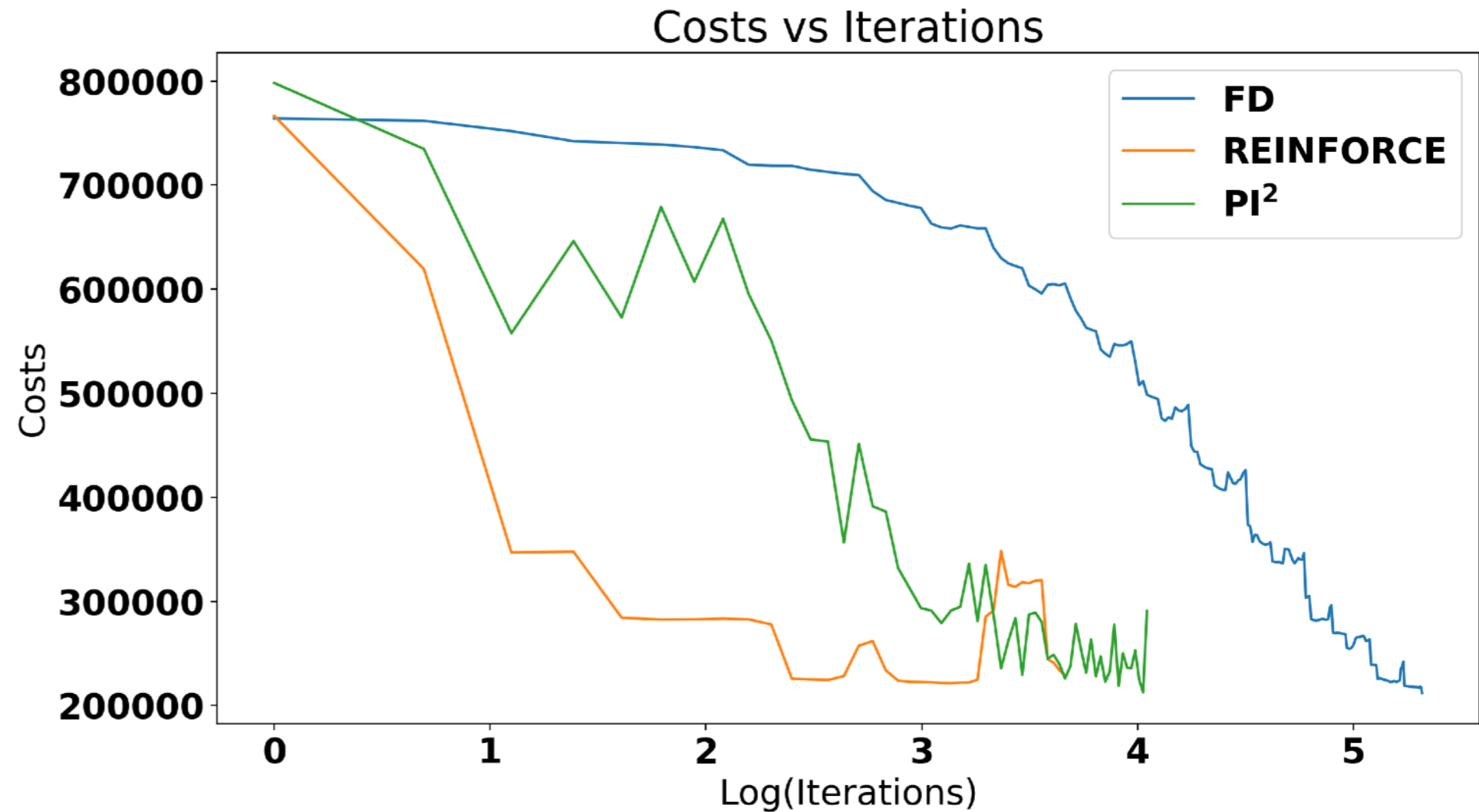
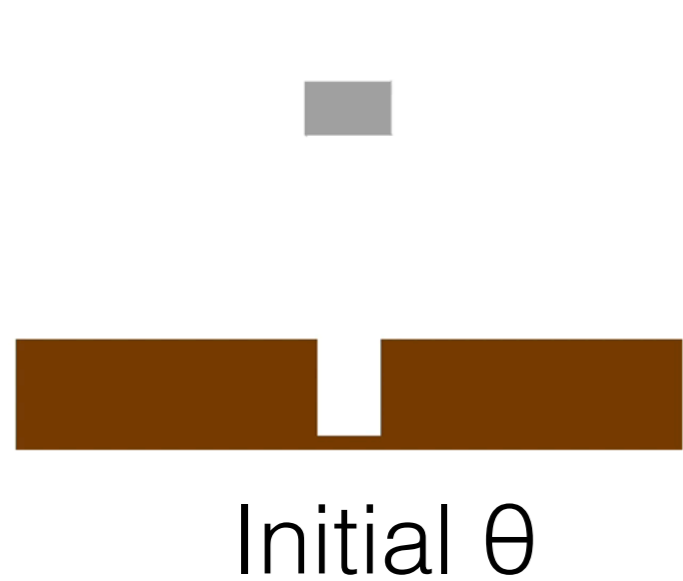
- $p = 1$ RBF Kernel

$$c_0 = .95 \quad h_0 = 1.$$



Results

Stochastic DMP with 2D θ



Finite Differences

REINFORCE

PI²



Problem Description

Linear Stochastic Policy

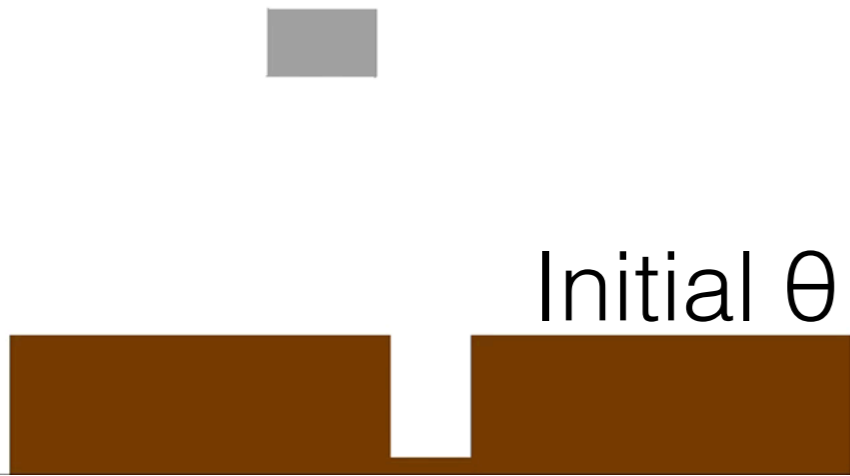
- Stochastic linear policy - the parameter determines all dimensions of the action

$$\pi(\mathbf{F} | \mathbf{s}; \boldsymbol{\theta}) = \boldsymbol{\theta} \mathbf{s} + \epsilon$$

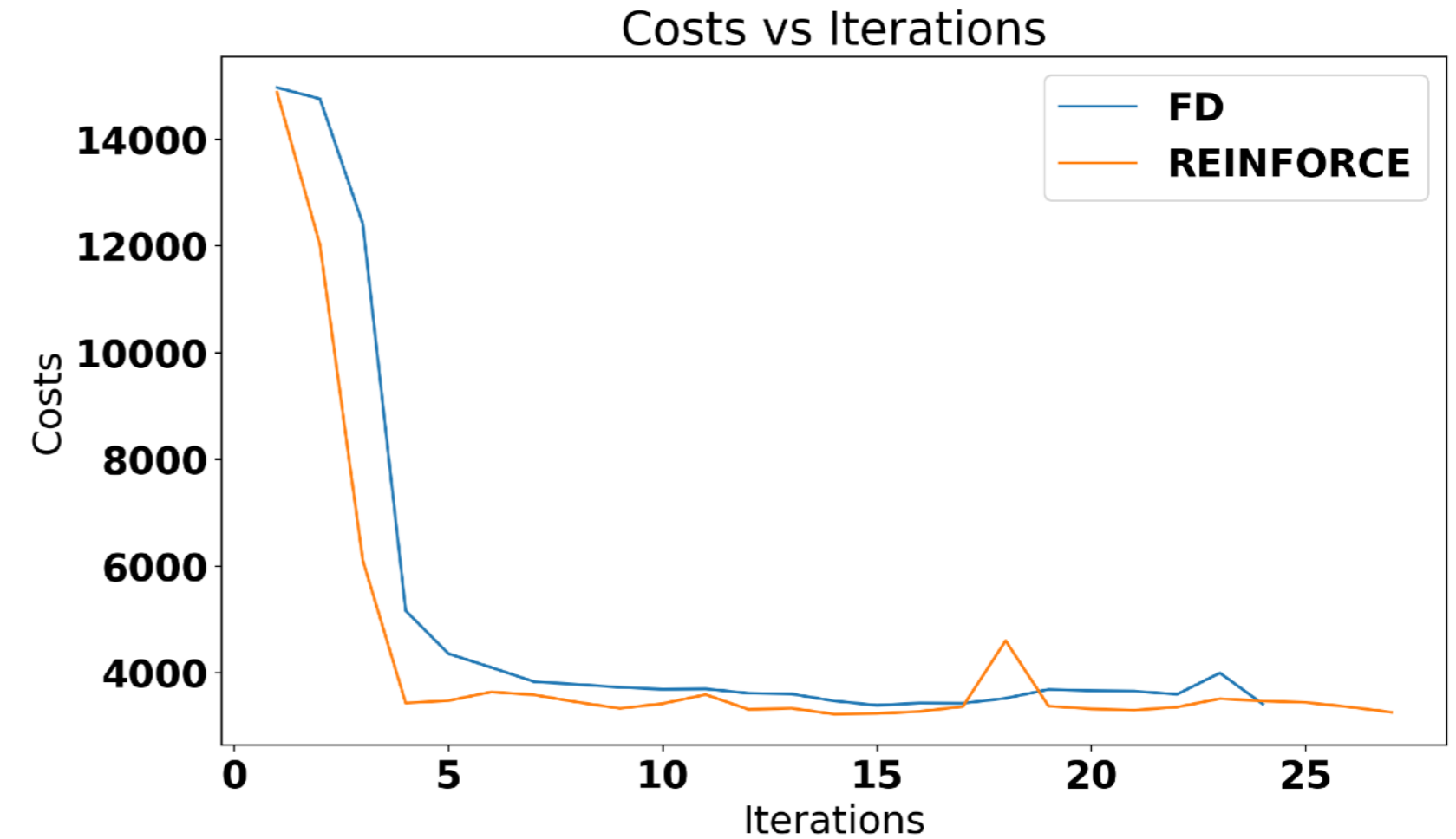
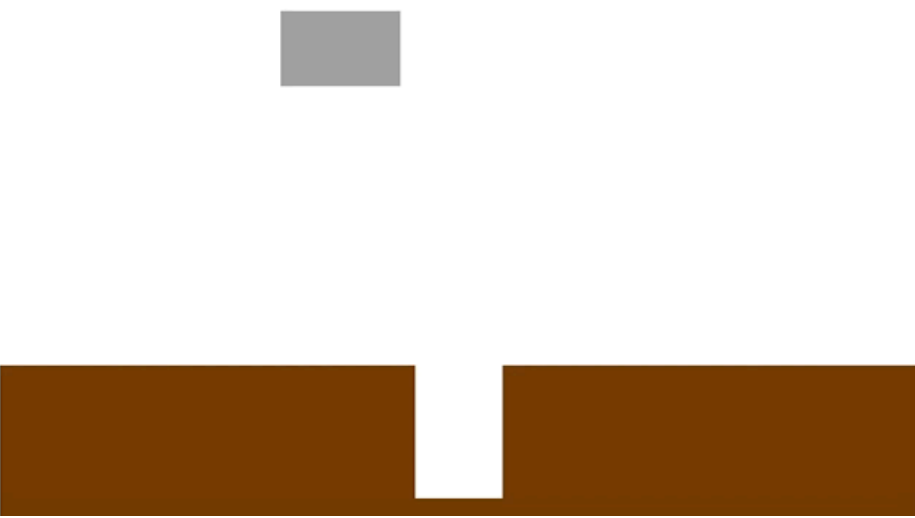
$$\mathbf{F} = \begin{bmatrix} F_x \\ F_y \\ F_\theta \end{bmatrix} \quad \mathbf{s} = \begin{bmatrix} x \\ y \\ \theta \\ v_x \\ v_y \\ v_\theta \end{bmatrix} \quad \mathbf{g} = \begin{bmatrix} g_x \\ g_y \\ g_\theta \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

Results

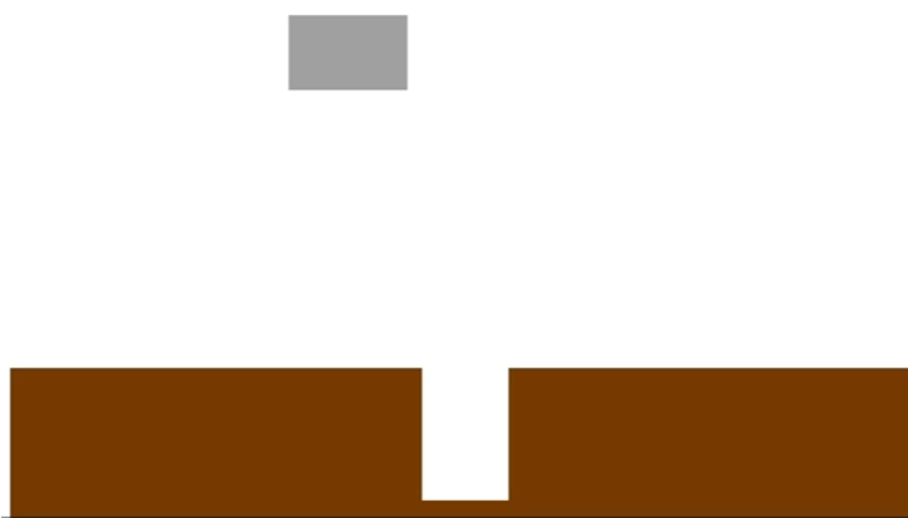
Linear Stochastic Policy



Finite Differences



REINFORCE



Analysis

- **Difficult** to manipulate a peg into a hole using simply **one-dimensional RBF** forcing term.
- **Linear stochastic policies are better** suited for the insertion task.
- **Non-smooth cost function** makes peg insertion task difficult to learn in general.



Novelty and Contribution

- Compared three methods of policy search: REINFORCE, PI^2 , and Finite Differences on a peg insertion task.
- Found that **REINFORCE outperforms both PI^2 and Finite Differences** for a 2D parameter with a stochastic **DMP** policy.
- **REINFORCE** and Finite Differences both perform **better** with a **linear** stochastic policy than with a **DMP** stochastic policy.



Thank you!

- Questions?
- Contact Information: carism@mit.edu

- References

[1] Peters, Jan (2010) Policy gradient methods. Scholarpedia, 5(11):3698.

[2] Theodorou, Evangelos, Jonas Buchli, and Stefan Schaal. "A generalized path integral control approach to reinforcement learning." Journal of Machine Learning Research 11.Nov (2010): 3137-3181.

[3] Zhao, Tingting, et al. "Analysis and improvement of policy gradient estimation." Advances in Neural Information Processing Systems. 2011.

