

W

Autonomous Robotics

Winter 2024

Abhishek Gupta

TAs: Karthikeya Vemuri, Arnav Thareja

Marius Memmel, Yunchu Zhang



Class Outline

State Estimation

Robotic System Design

Filtering

Localization

SLAM

Control

Feedback Control

PID Control

MPC

LQR

Planning

Search

Heuristic Search

Motion Planning

Lazy Search

Learning

Imitation Learning

Policy Gradient

Actor-Critic

Model-Based RL

Logistics

- One paper presentation Feb 26
- Guest lecture 1 moved to March 1
- Project 5 (Final Project) will be released Feb 26.

Lecture Outline

Imitation Learning Recap + Compounding Error



Imitation Learning: Improvements – Multimodality



Policy Gradient

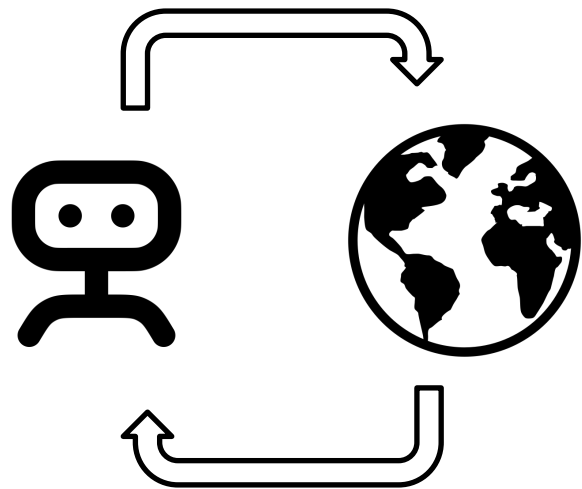


Improving Policy Gradient

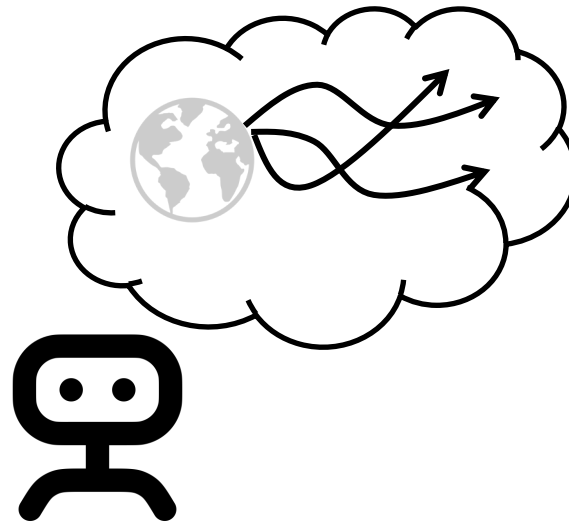
Ok so how can we learn policies?

$$\max_{\theta} \mathbb{E}_{\tau \sim \pi_{\theta}} \left[\sum_{t=0}^T r(s_t, a_t) \right]$$

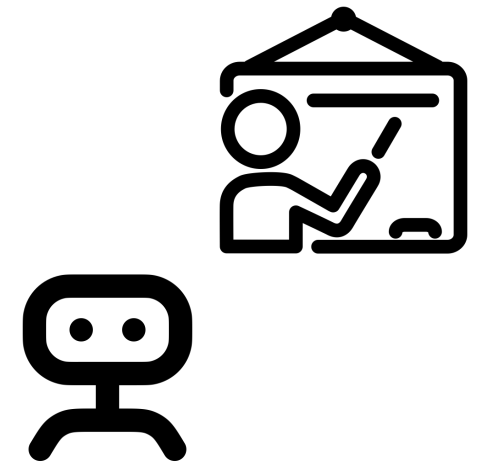
Model-free RL



Model-based RL



Imitation Learning



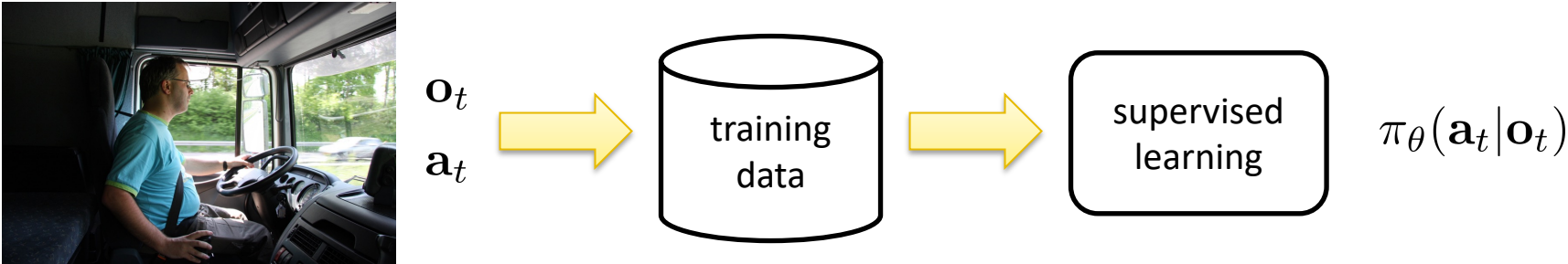
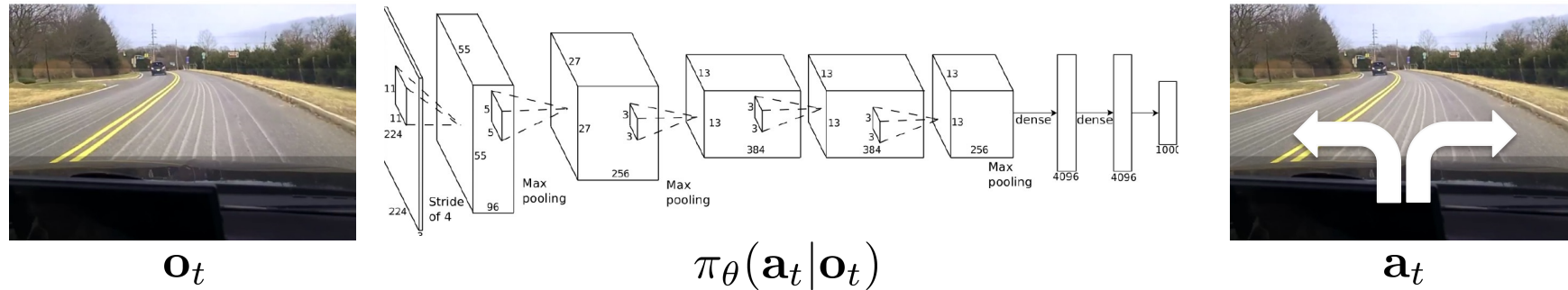
Idea 1: Imitation Learning via Behavior Cloning

Given: Demonstrations of optimal behavior

$$\arg \max_{\theta} \mathbb{E}_{(s^*, a^*) \sim \mathcal{D}} [\log \pi_{\theta}(a^* | s^*)]$$

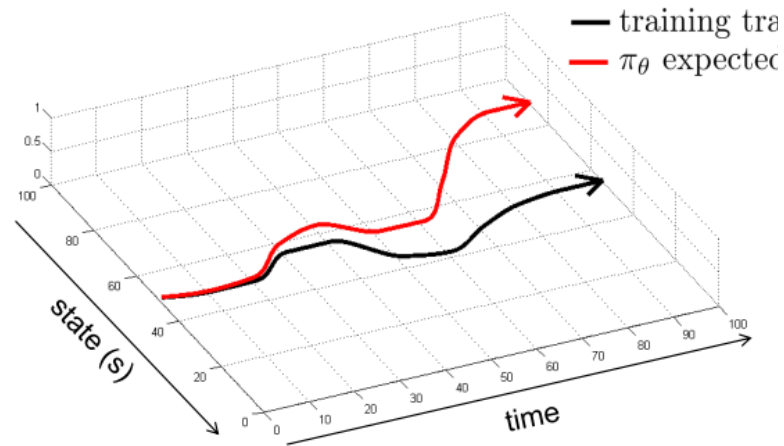
Goal: Train a policy to mimic the demonstrator

Idea: Treat imitation learning as a supervised learning problem!



So does behavior cloning really work?

- Imitation Learning \neq Supervised Learning



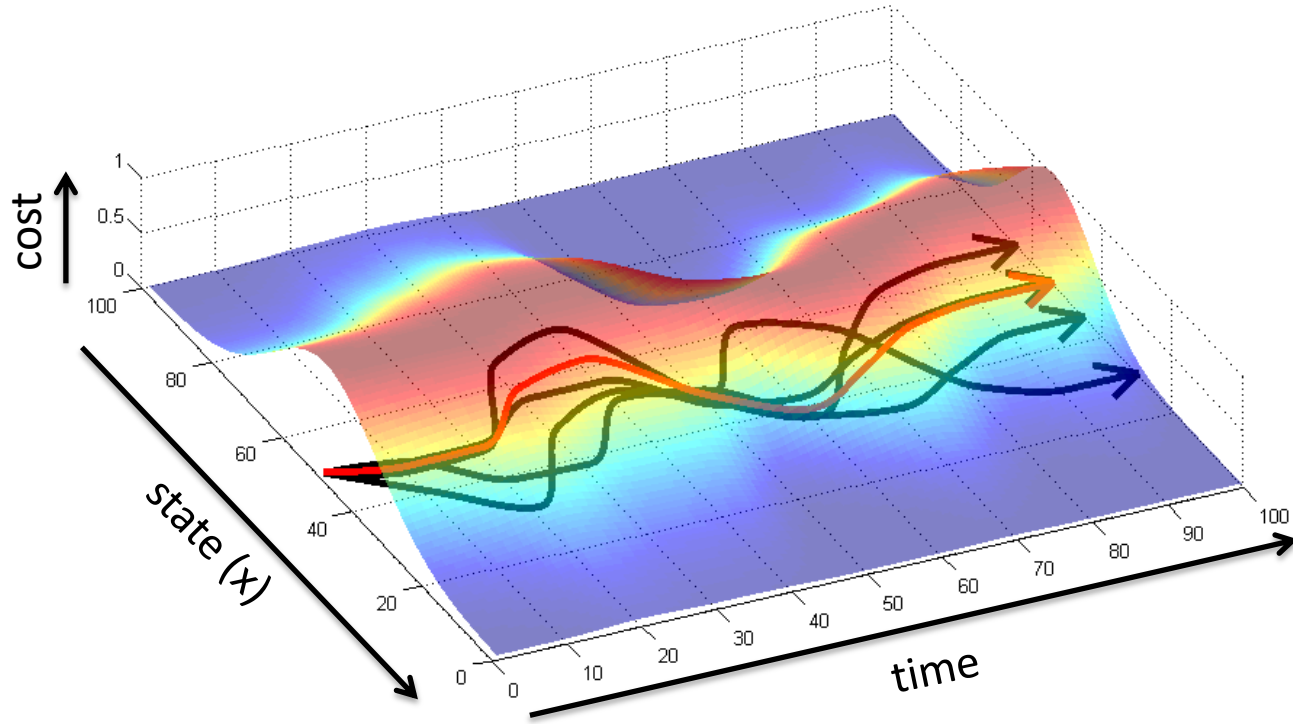
$$\arg \max_{\theta} \mathbb{E}_{(s^*, a^*) \sim \mathcal{D}} [\log \pi_{\theta}(a^* | s^*)] \quad \mathbb{E}_{(s, a) \sim \rho(\pi)} [\mathbf{1}(a = a^*)]$$



Not the same!

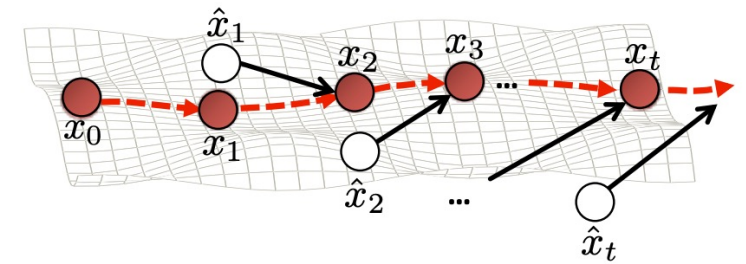
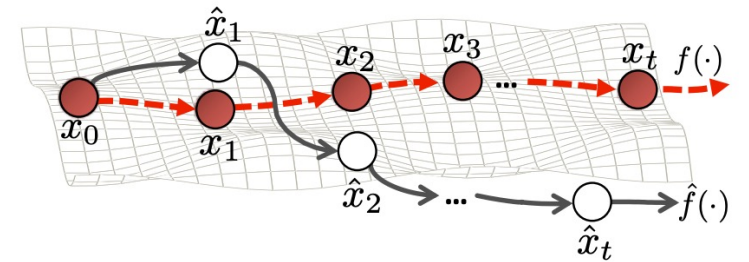
What is the general principle?

- training trajectory
- π_θ expected trajectory



stability

Corrective labels that bring you back to the data



Concrete Instantiation: DAgger

can we make $p_{\text{data}}(\mathbf{o}_t) = p_{\pi_\theta}(\mathbf{o}_t)$?

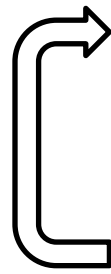
idea: instead of being clever about $p_{\pi_\theta}(\mathbf{o}_t)$, be clever about $p_{\text{data}}(\mathbf{o}_t)$!

DAgger: Dataset Aggregation

goal: collect training data from $p_{\pi_\theta}(\mathbf{o}_t)$ instead of $p_{\text{data}}(\mathbf{o}_t)$

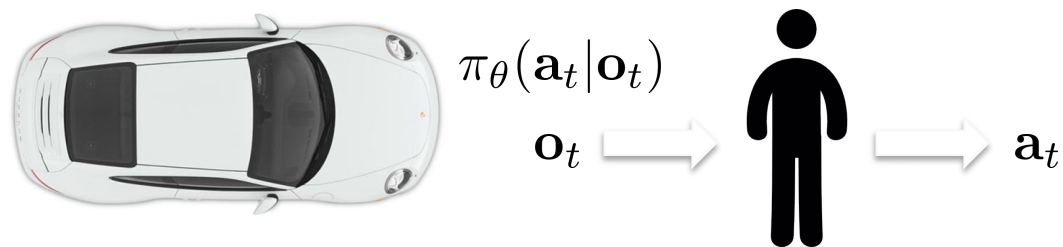
how? just run $\pi_\theta(\mathbf{a}_t|\mathbf{o}_t)$

but need labels \mathbf{a}_t !

- 
1. train $\pi_\theta(\mathbf{a}_t|\mathbf{o}_t)$ from human data $\mathcal{D} = \{\mathbf{o}_1, \mathbf{a}_1, \dots, \mathbf{o}_N, \mathbf{a}_N\}$
 2. run $\pi_\theta(\mathbf{a}_t|\mathbf{o}_t)$ to get dataset $\mathcal{D}_\pi = \{\mathbf{o}_1, \dots, \mathbf{o}_M\}$
 3. Ask human to label \mathcal{D}_π with actions \mathbf{a}_t
 4. Aggregate: $\mathcal{D} \leftarrow \mathcal{D} \cup \mathcal{D}_\pi$

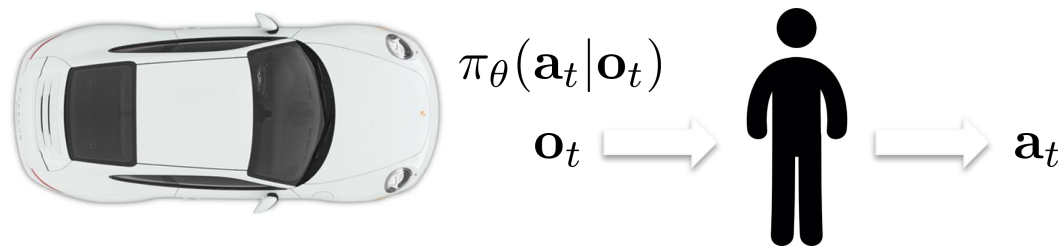
How might we fix this?

- "Generate" corrective labels automatically
1. train $\pi_{\theta}(\mathbf{a}_t|\mathbf{o}_t)$ from human data $\mathcal{D} = \{\mathbf{o}_1, \mathbf{a}_1, \dots, \mathbf{o}_N, \mathbf{a}_N\}$
 2. run $\pi_{\theta}(\mathbf{a}_t|\mathbf{o}_t)$ to get dataset $\mathcal{D}_{\pi} = \{\mathbf{o}_1, \dots, \mathbf{o}_M\}$
 3. Ask human to label \mathcal{D}_{π} with actions \mathbf{a}_t
 4. Aggregate: $\mathcal{D} \leftarrow \mathcal{D} \cup \mathcal{D}_{\pi}$
- Do at data collection time



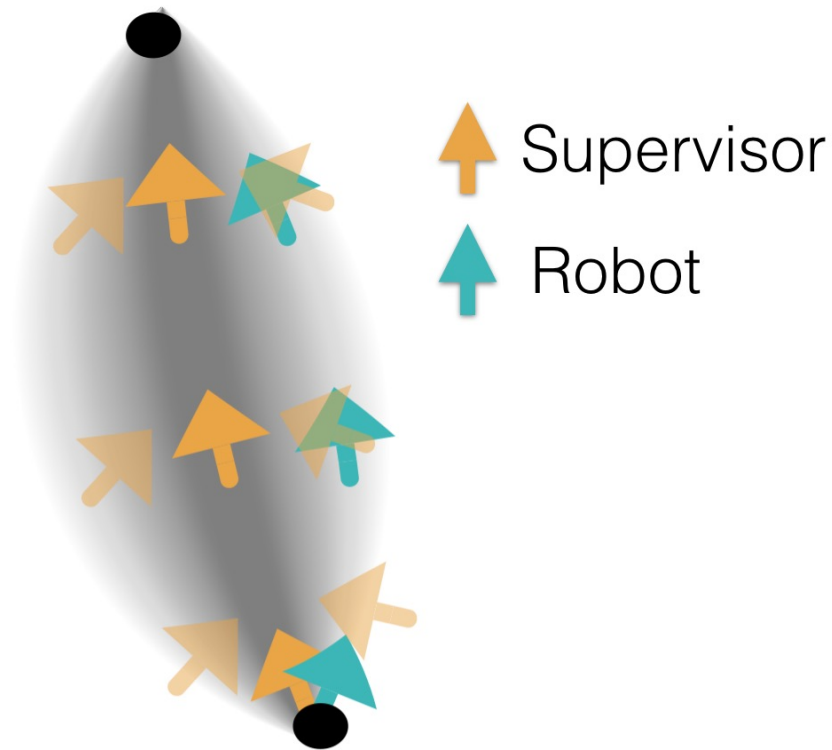
How might we fix this?

1. train $\pi_\theta(\mathbf{a}_t|\mathbf{o}_t)$ from human data $\mathcal{D} = \{\mathbf{o}_1, \mathbf{a}_1, \dots, \mathbf{o}_N, \mathbf{a}_N\}$
 2. run $\pi_\theta(\mathbf{a}_t|\mathbf{o}_t)$ to get dataset $\mathcal{D}_\pi = \{\mathbf{o}_1, \dots, \mathbf{o}_M\}$
 3. Ask human to label \mathcal{D}_π with actions \mathbf{a}_t
 4. Aggregate: $\mathcal{D} \leftarrow \mathcal{D} \cup \mathcal{D}_\pi$
- Do at data collection time



Noising the Data Collection Process

Key idea: force the human to correct for noise during training



Noise Injection

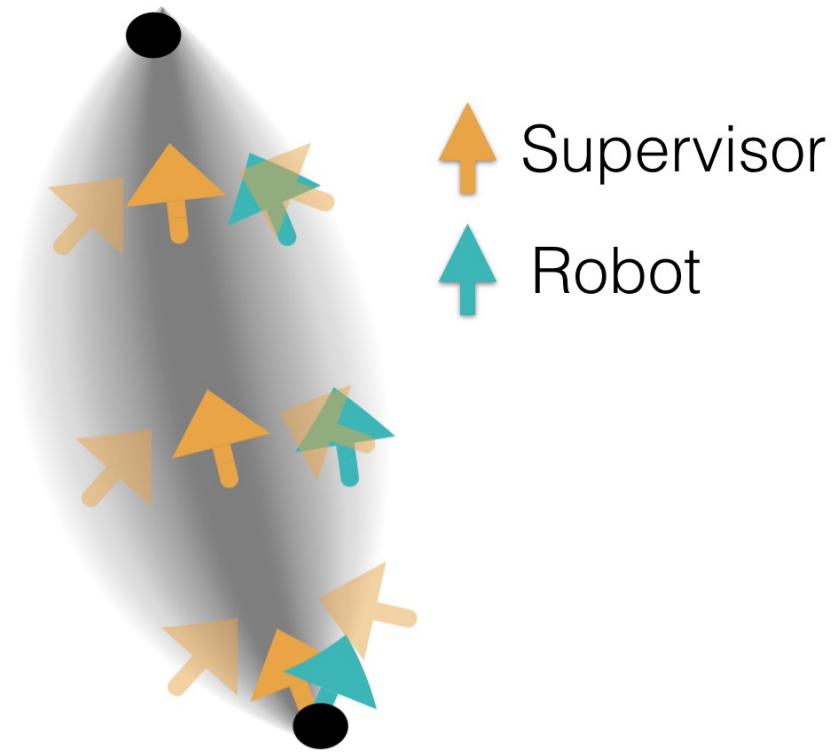
$$\hat{\psi}_{k+1} = \underset{\psi}{\operatorname{argmin}} E_{p(\xi|\pi_{\theta^*}, \psi_k)} - \sum_{t=0}^{T-1} \log [\pi_{\theta^*}(\pi_{\hat{\theta}}(\mathbf{x}_t)|\mathbf{x}_t, \psi)]$$

↑
Maximize likelihood

↑
Under noise during data collection

Noising the Data Collection Process

Key idea: force the human to correct for noise during training



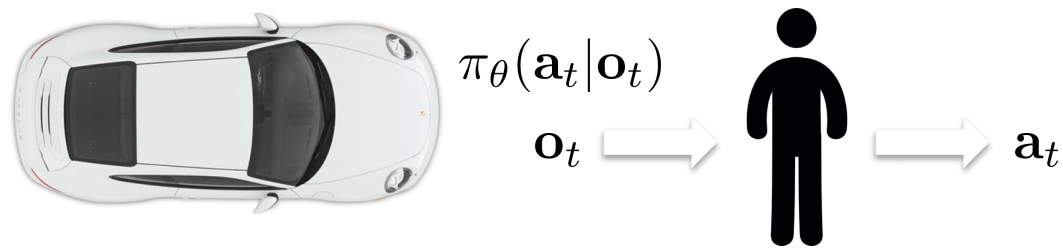
Noise Injection



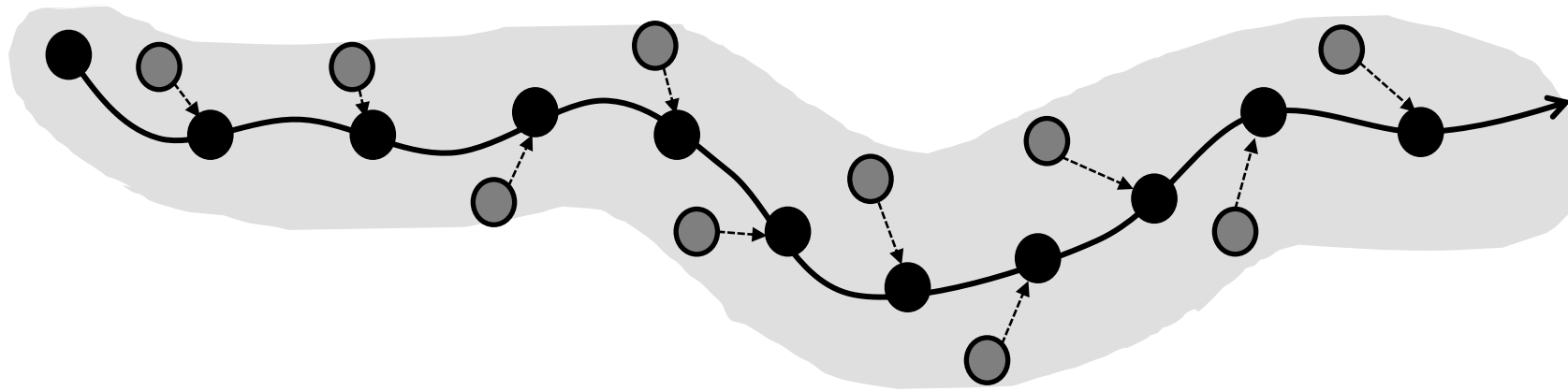
How might we fix this?

"Generate"
corrective labels
automatically

1. train $\pi_{\theta}(\mathbf{a}_t|\mathbf{o}_t)$ from human data $\mathcal{D} = \{\mathbf{o}_1, \mathbf{a}_1, \dots, \mathbf{o}_N, \mathbf{a}_N\}$
2. run $\pi_{\theta}(\mathbf{a}_t|\mathbf{o}_t)$ to get dataset $\mathcal{D}_{\pi} = \{\mathbf{o}_1, \dots, \mathbf{o}_M\}$
3. Ask human to label \mathcal{D}_{π} with actions \mathbf{a}_t
4. Aggregate: $\mathcal{D} \leftarrow \mathcal{D} \cup \mathcal{D}_{\pi}$



Generating Corrective Labels for Imitation Learning with Learned Dynamics



minimizing MSE on expert data
+ spectral norm

When can we trust learned dynamics \hat{f}_ϕ ?



Trust models “around” training data

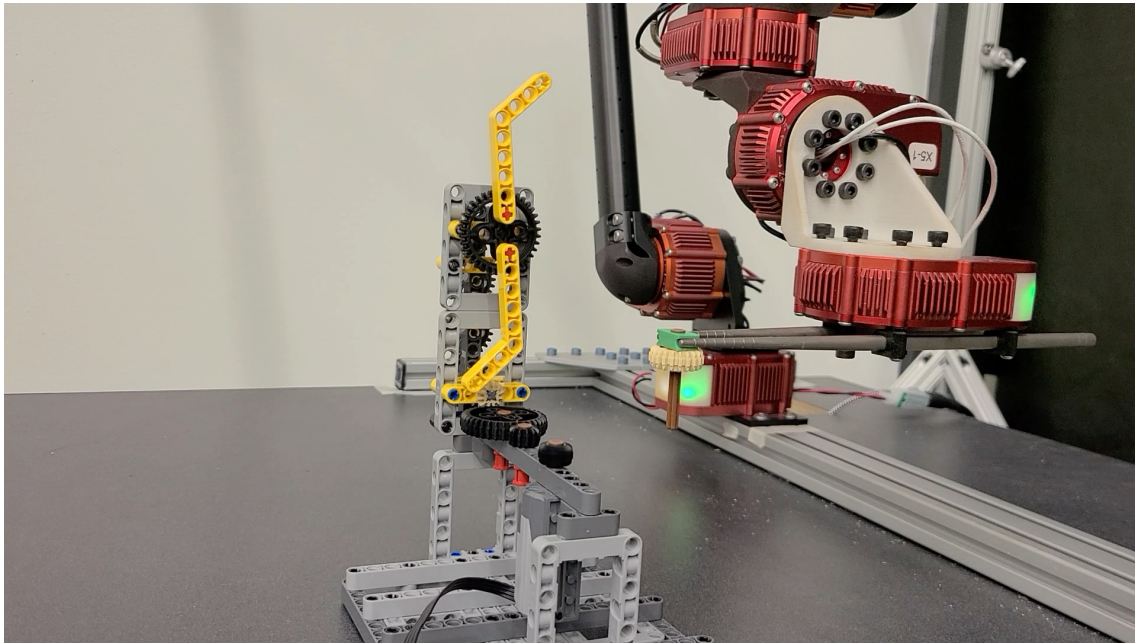
$$\|s_{t+1}^* - \hat{f}_\phi(s_t, a_t)\| \leq \epsilon$$

Find states (s_t), actions (a_t) that lead back to optimal states under ~~true~~ learned dynamics,
where learned dynamics can be trusted

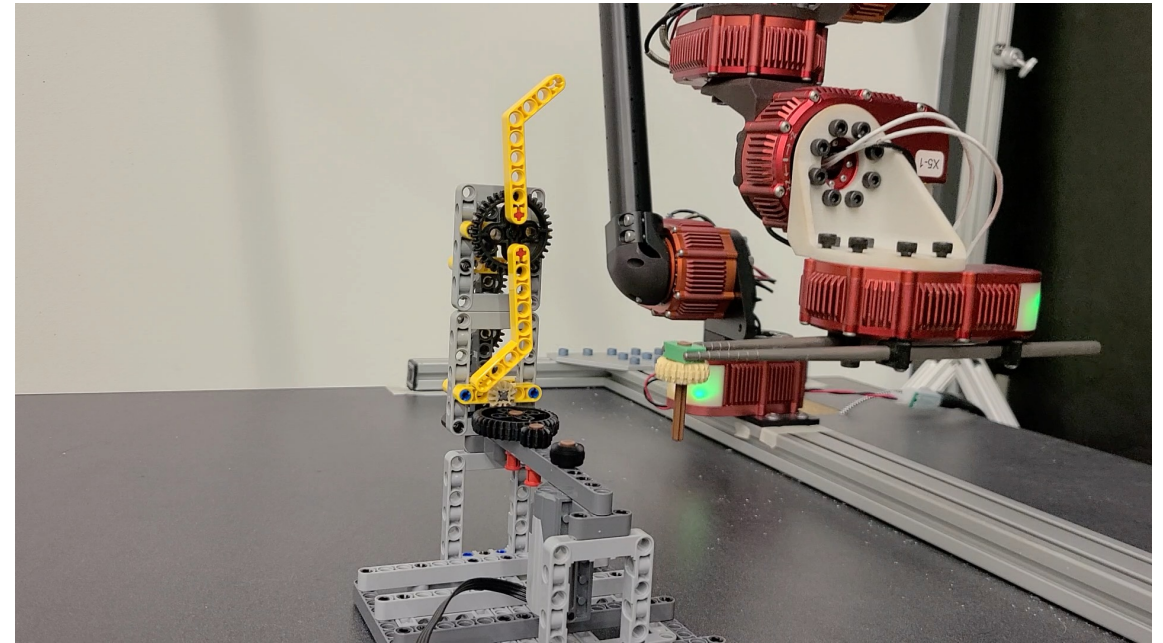
Overall Learning Pipeline with Corrective Labels

Learn models of dynamics and find states (s_t), actions (a_t) that lead back to optimal states under this learned dynamics

Standard behavior cloning



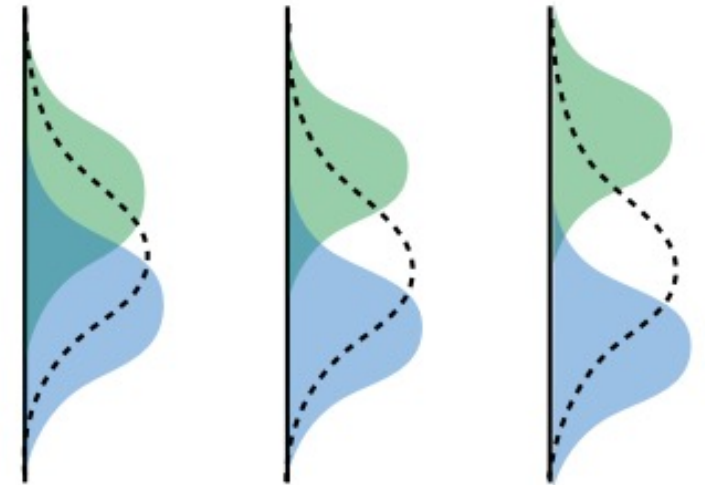
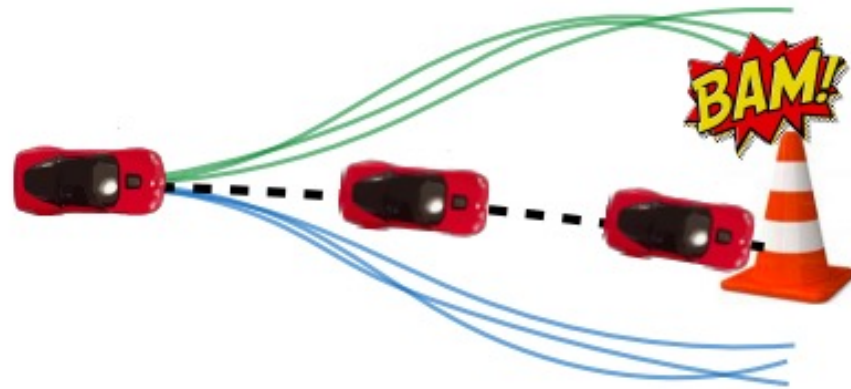
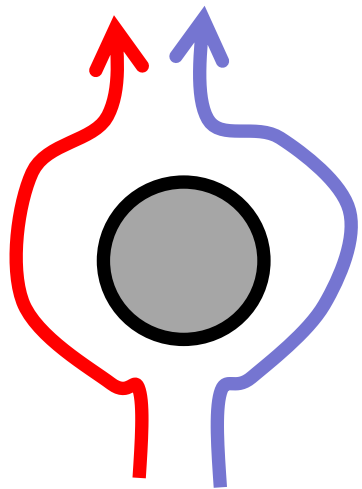
Corrective labels



So does this solve all the issues in imitation?

Why might we fail to fit the expert?

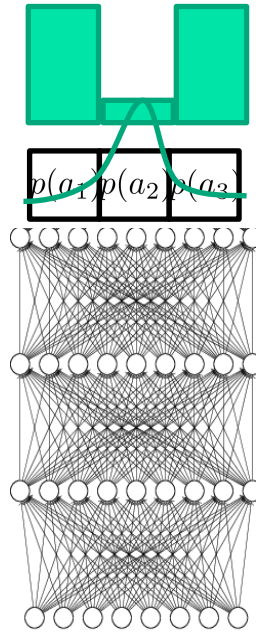
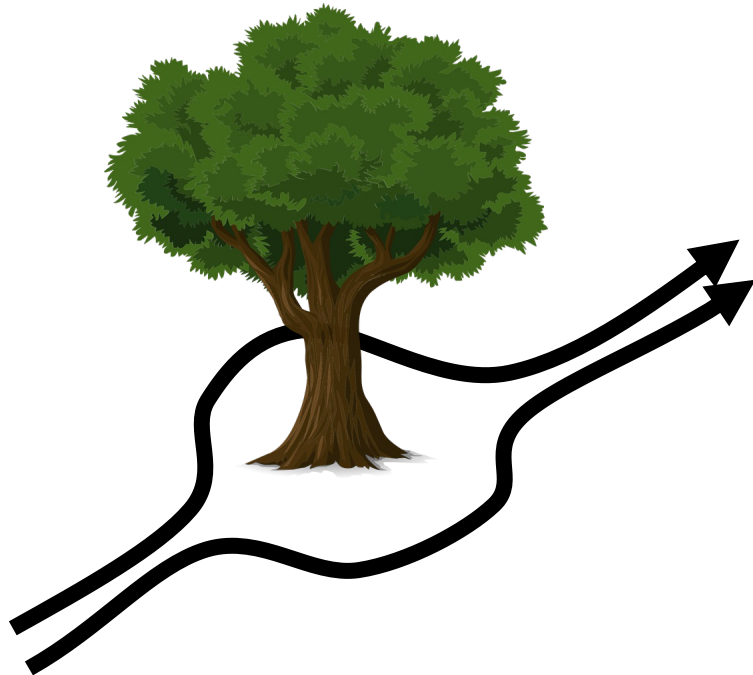
Multimodal behavior.. amongst other reasons



Not a matter of network size! It's about distributional expressivity

Why might we fail to fit the expert?

Multimodal behavior → use more expressive probability distributions

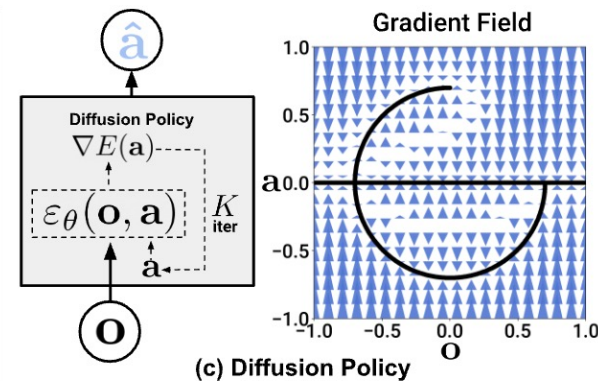
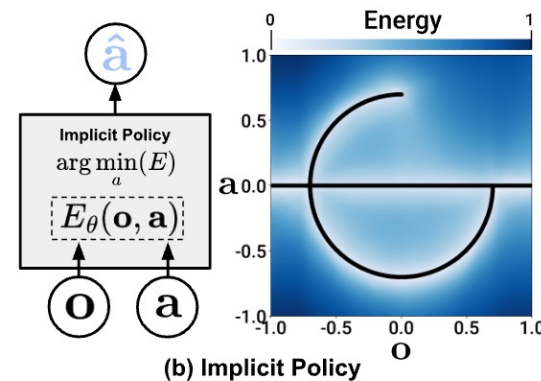
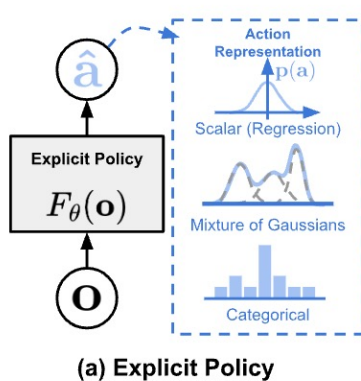
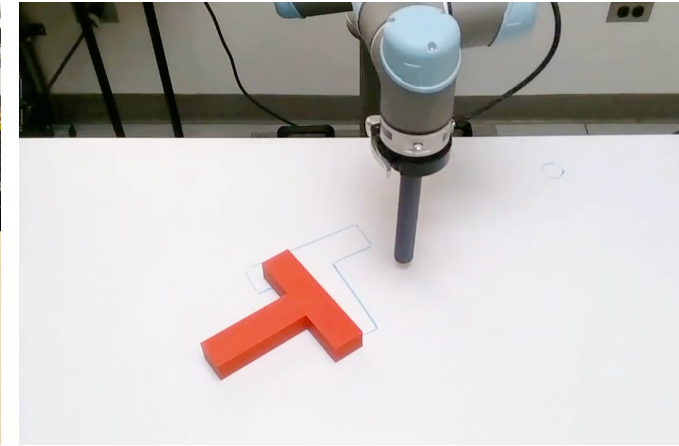
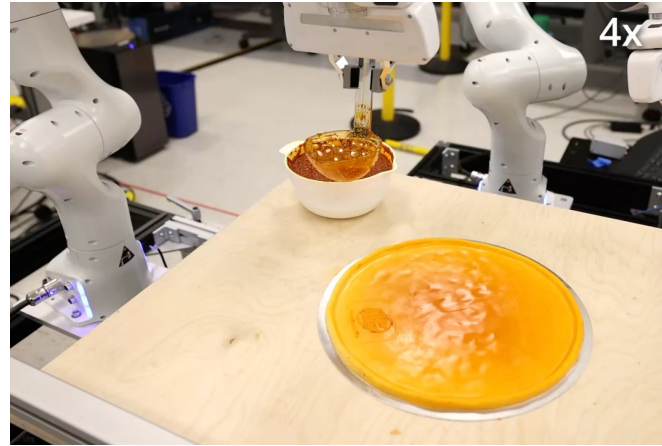


1. Output mixture of Gaussians
2. Latent variable models
3. Autoregressive discretization
4. Diffusion models
5. ...



Why might we fail to fit the expert?

1. Output mixture of Gaussians
2. Latent variable models
3. Autoregressive discretization
- ➔ 4. Diffusion models
5. ...



Some cool imitation videos

1x and tesla humanoid robots

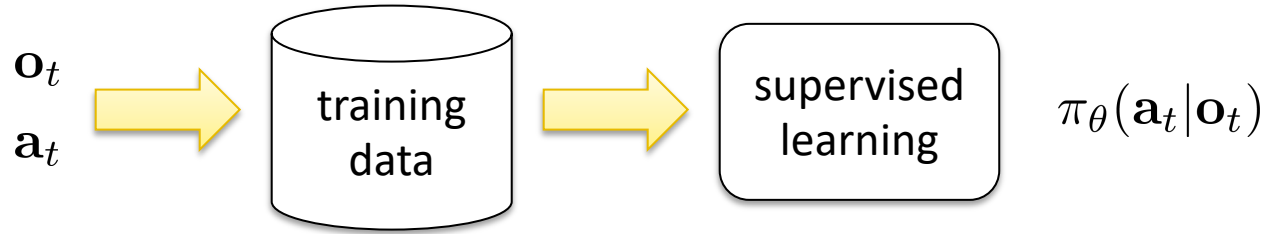


● 1X END-TO-END AUTONOMY
UPDATE, JAN 2024

ALOHA and CherryBot Fine Manipulation

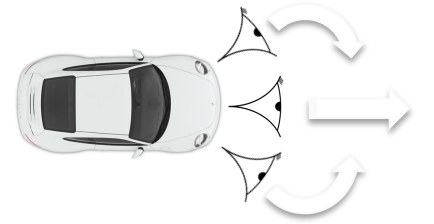


Perspectives on Imitation – don't believe everything you see online



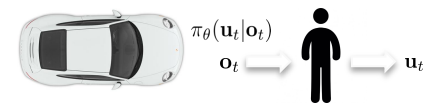
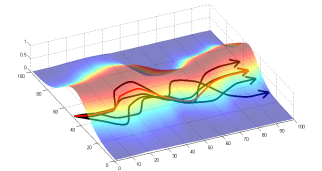
■ Pros:

- Easy to use, no additional infra
- Can sometimes be unreasonably effective



■ Cons:

- Challenges of compounding error, multimodality
- Doesn't really generalize
- Very expensive in terms of data collection!



Lecture Outline

A Formalism for Sequential Decision Making



Imitation Learning: Behavior Cloning



Imitation Learning: Improvements – Compounding Error



Imitation Learning: Improvements – Multimodality

Class Outline

State Estimation

Robotic System Design

Filtering

Localization

SLAM

Control

Feedback Control

PID Control

MPC

LQR

Planning

Search

Heuristic Search

Motion Planning

Lazy Search

Learning

Imitation Learning

Policy Gradient

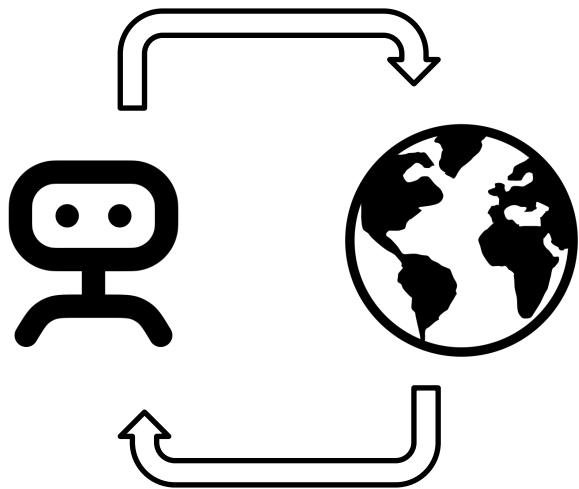
Actor-Critic

Model-Based RL

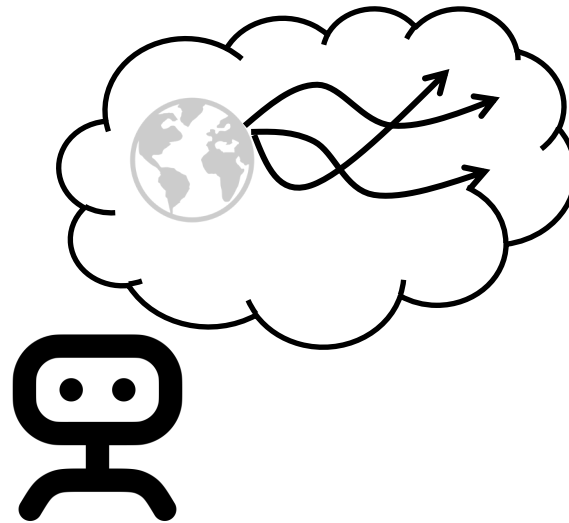
Ok so how can we learn policies?

$$\max_{\theta} \mathbb{E}_{\tau \sim \pi_{\theta}} \left[\sum_{t=0}^T r(s_t, a_t) \right]$$

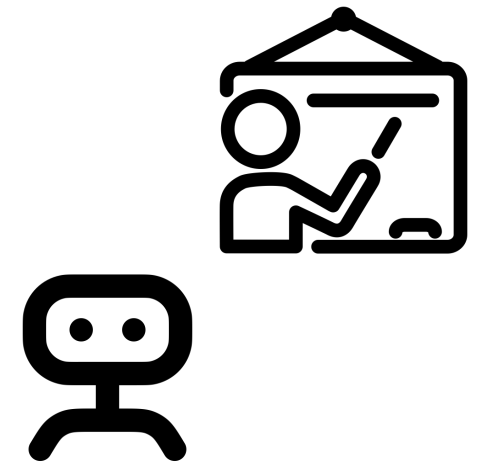
Model-free RL



Model-based RL



Imitation Learning



What if we just performed gradient ascent?

$$\begin{aligned} \max_{\theta} \mathbb{E}_{\tau \sim \pi_{\theta}} \left[\sum_{t=0}^T r(s_t, a_t) \right] \\ = \int p_{\theta}(\tau) R(\tau) d\tau \end{aligned}$$

Standard gradient descent (supervised learning)

$$\nabla_{\theta} \mathbb{E}_{x \sim g(x)} [f_{\theta}(x)]$$

REINFORCE gradient descent (RL)

$$\nabla_{\theta} \mathbb{E}_{x \sim p_{\theta}(x)} [f(x)]$$

Gradient wrt expectation variable, not of integrand!

Taking the gradient of sum of rewards

$$J(\theta) = \int p_{\theta}(\tau) R(\tau) d(\tau)$$

$$\nabla_{\theta} J(\theta) = \nabla_{\theta} \int p_{\theta}(\tau) R(\tau) d(\tau)$$

$$= \int \nabla_{\theta} p_{\theta}(\tau) R(\tau) d(\tau) = \int \frac{p_{\theta}(\tau)}{p_{\theta}(\tau)} \nabla_{\theta} p_{\theta}(\tau) R(\tau) d(\tau)$$

$$= \int p_{\theta}(\tau) \nabla_{\theta} \log p_{\theta}(\tau) R(\tau) d(\tau) = \mathbb{E}_{p_{\theta}(\tau)} [\nabla_{\theta} \log p_{\theta}(\tau) R(\tau)]$$

REINFORCE trick

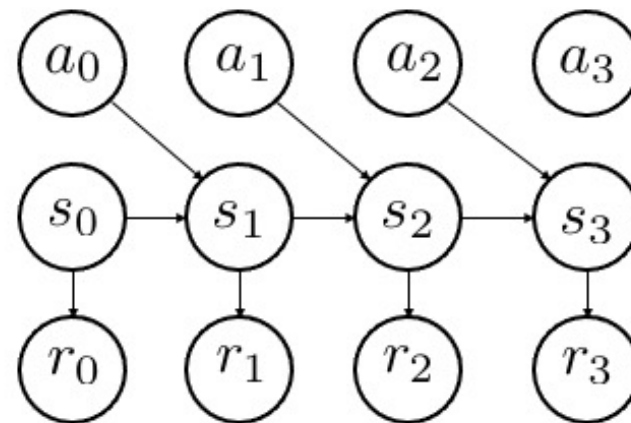
Taking the gradient of return

Initial State

Dynamics

Policy

$$p_{\theta}(\tau) = p(s_0) \prod_{t=0}^{T-1} p(s_{t+1} | s_t, a_t) \pi(a_t | s_t)$$



$$\log p_{\theta}(\tau) = \log p(s_0) + \sum_{t=0}^{T-1} \log p(s_{t+1} | s_t, a_t) + \log \pi(a_t | s_t)$$

$$\nabla_{\theta} \log p_{\theta}(\tau) = \cancel{\nabla_{\theta} \log p(s_0)} + \sum_{t=0}^{T-1} \cancel{\nabla_{\theta} \log p(s_{t+1} | s_t, a_t)} + \nabla_{\theta} \log \pi(a_t | s_t)$$

$$\nabla_{\theta} \log p_{\theta}(\tau) = \sum_{t=0}^{T-1} \nabla_{\theta} \log \pi(a_t | s_t)$$

Model Free!!

Taking the gradient of return

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\tau \sim p_{\theta}(\tau)} \left[\nabla_{\theta} \log p_{\theta}(\tau) \sum_{t=0}^T r(s_t, a_t) \right]$$

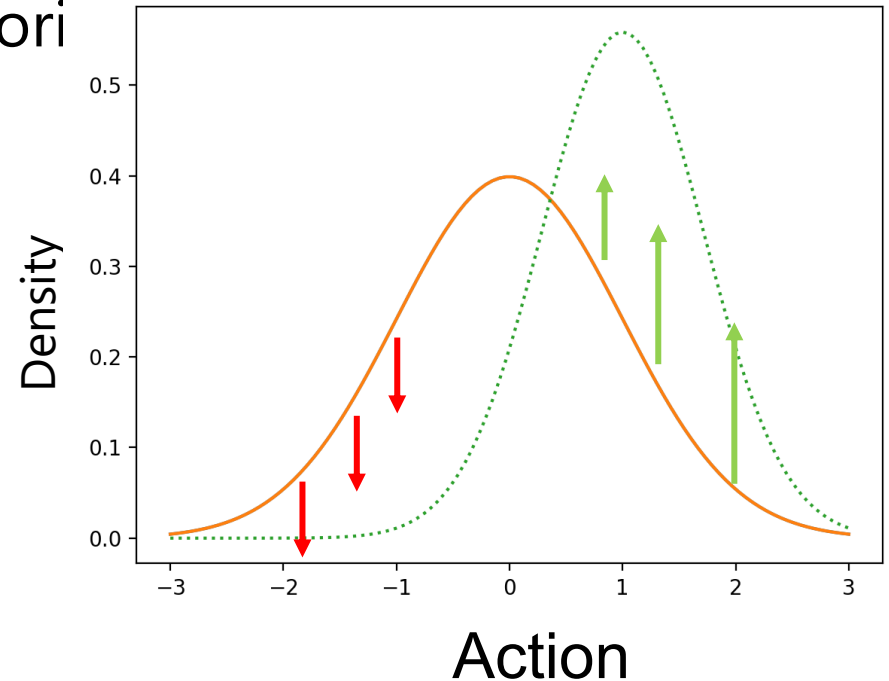
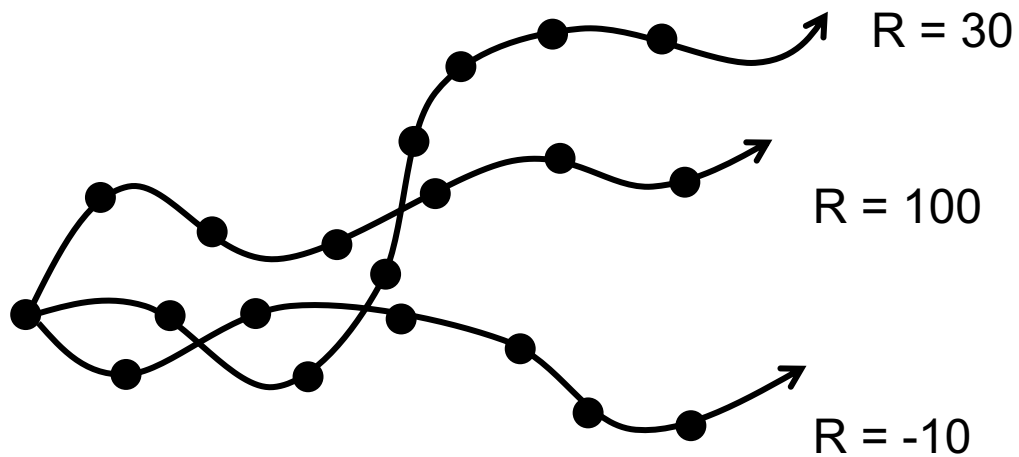
$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\substack{s_0 \sim p(s_0) \\ s_{t+1} \sim p(s_{t+1} | s_t, a_t) \\ a_t \sim \pi(a_t | s_t)}} \left[\sum_{t=0}^T \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \sum_{t'=0}^T r(s_{t'}, a_{t'}) \right]$$

$$\approx \frac{1}{N} \sum_{i=0}^N \sum_{t=0}^T \nabla_{\theta} \log \pi_{\theta}(a_t^i | s_t^i) \sum_{t'=0}^T r(s_{t'}^i, a_{t'}^i) \quad (\text{approximating using samples})$$

What does this mean?

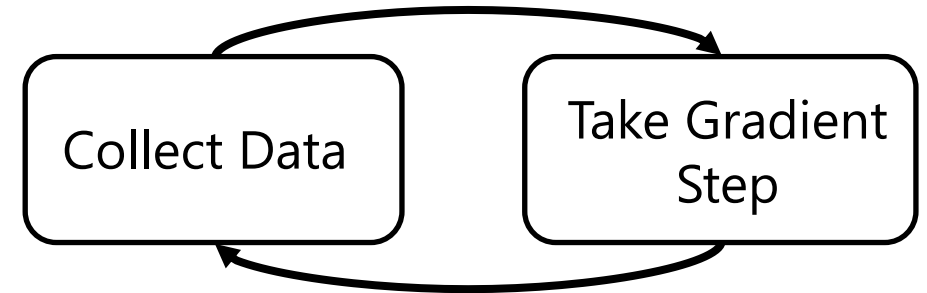
$$\nabla_{\theta} J(\theta) = \int p_{\theta}(\tau) \nabla_{\theta} \log p_{\theta}(\tau) d\tau \approx \frac{1}{N} \sum_{i=0}^N \sum_{t=0}^T \nabla_{\theta} \log \pi_{\theta}(a_t^i | s_t^i) \sum_{t'=0}^T r(s_{t'}^i, a_{t'}^i)$$

Increase the likelihood of actions in high return trajectories



Resulting Algorithm (REINFORCE)

$$\nabla_{\theta} J(\theta) = \int p_{\theta}(\tau) \nabla_{\theta} \log p_{\theta}(\tau) d\tau$$



REINFORCE algorithm:

On-policy



1. sample $\{\tau^i\}$ from $\pi_{\theta}(\mathbf{a}_t|\mathbf{s}_t)$ (run it on the robot)
2. $\nabla_{\theta} J(\theta) \approx \sum_i (\sum_t \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t^i|\mathbf{s}_t^i)) (\sum_t r(\mathbf{s}_t^i, \mathbf{a}_t^i))$
3. $\theta \leftarrow \theta + \alpha \nabla_{\theta} J(\theta)$

How is this related to supervised learning?

Reinforcement Learning

$$\nabla_{\theta} J(\theta) = \int p_{\theta}(\tau) \nabla_{\theta} \log p_{\theta}(\tau) d\tau$$

$$\approx \frac{1}{N} \sum_{i=0}^N \sum_{t=0}^T \nabla_{\theta} \log \pi_{\theta}(a_t^i | s_t^i) \sum_{t'=0}^T r(s_{t'}^i, a_{t'}^i)$$

Supervised Learning

$$\max_{\theta} \mathbb{E}_{(x,y) \sim \mathcal{D}} [\log p_{\theta}(y|x)]$$

$$\approx \frac{1}{N} \sum_i \nabla_{\theta} \log p_{\theta}(y^i | x^i)$$

PG = select good data + increase likelihood of selected data

What makes policy gradient challenging?

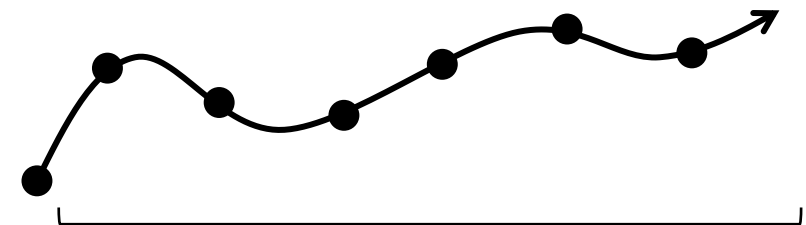
$$\nabla_{\theta} J(\theta) = \int p_{\theta}(\tau) \nabla_{\theta} \log p_{\theta}(\tau) d\tau$$

$$\approx \frac{1}{N} \sum_{i=0}^N \sum_{t=0}^T \nabla_{\theta} \log \pi_{\theta}(a_t^i | s_t^i) \sum_{t'=0}^T r(s_{t'}^i, a_{t'}^i)$$

High variance estimator!!

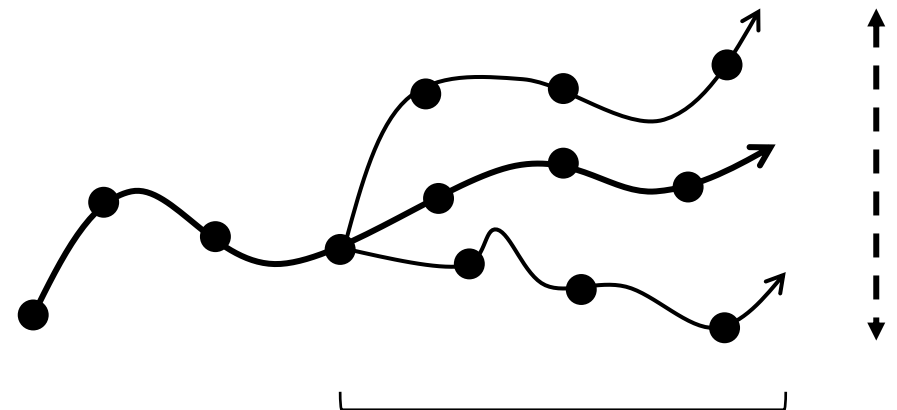
Hard to tell what matters without many samples

What we do



Single sample estimate

What we actually want



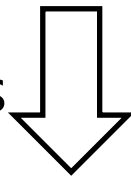
Averaged return estimate

Variance Reduction with Causality

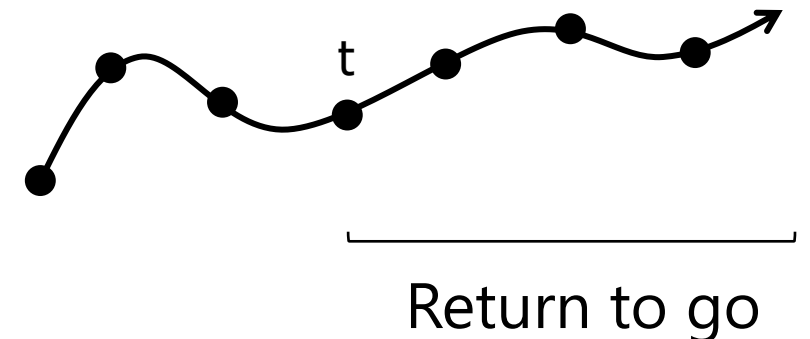
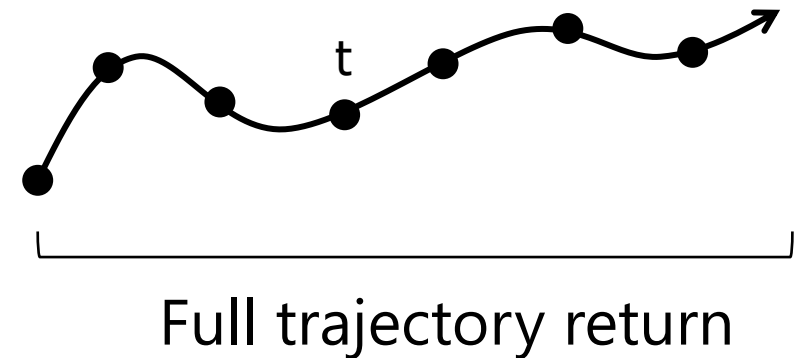
Idea: Trajectory returns depend on past and future, but we only care about the future, since actions cannot affect the past. Instead, consider “return-to-go”

$$\approx \frac{1}{N} \sum_{i=0}^N \sum_{t=0}^T \nabla_{\theta} \log \pi_{\theta}(a_t^i | s_t^i) \underbrace{\sum_{t'=0}^T r(s_{t'}^i, a_{t'}^i)}_{\text{Includes } t' < t}$$

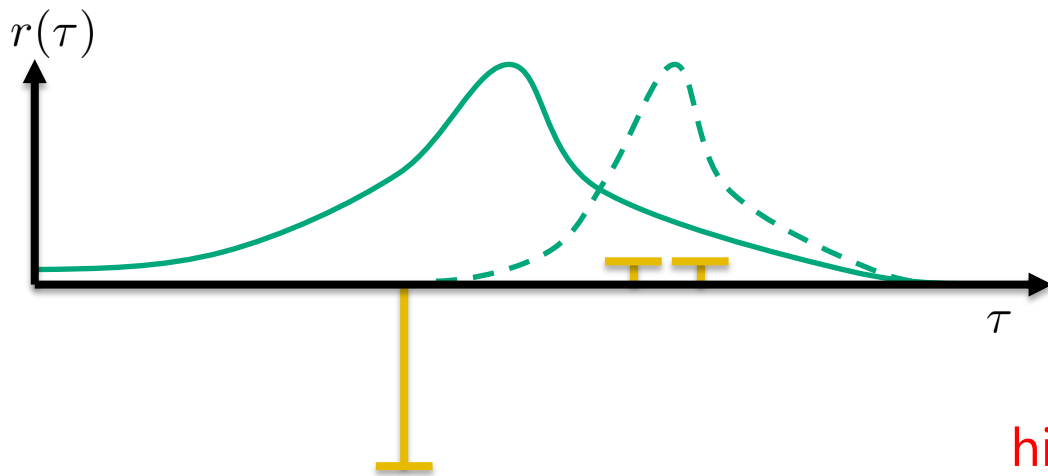
Ignore past terms



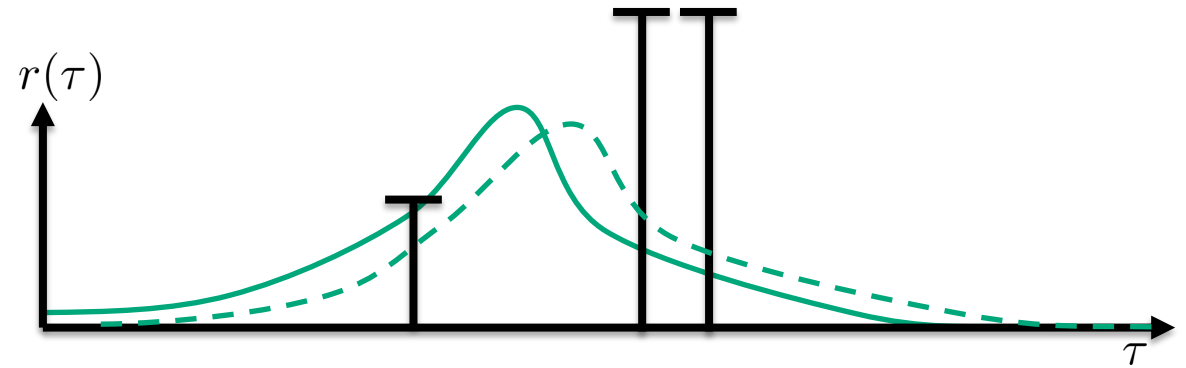
$$\frac{1}{N} \sum_{i=0}^N \sum_{t=0}^T \nabla_{\theta} \log \pi_{\theta}(a_t^i | s_t^i) \sum_{t'=t}^T r(s_{t'}^i, a_{t'}^i)$$



Can we reduce variance further?



high variance




Arbitrarily uncentered, scaled returns can lead to huge variance:

- Imagine all rewards were positive, every action would be pushed up, some more than others
- What if instead, we pushed down some actions and pushed up some others (even if rewards are positive)

Variance Reduction with a Baseline

Idea: We can reduce variance by subtracting a current state dependent function from the policy gradient return

$$\frac{1}{N} \sum_{i=1}^N \sum_{t=0}^T \nabla_{\theta} \log \pi_{\theta}(a_t^i | s_t^i) \left[\sum_{t'=t}^T r(s_{t'}^i, a_{t'}^i) - b(s_t) \right]$$


Baseline: Centers the returns, reduces variance

But does this increase bias??

Variance Reduction with a Baseline

$$\int_{\mathcal{S}} \int_{\mathcal{A}} p(s_t, a_t) \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \left[\sum_{t'=t}^T r(s_{t'}, a_{t'}) - b(s_t) \right] ds_t da_t$$

$$\int_{\mathcal{S}} \int_{\mathcal{A}} p(s_t, a_t) \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \left[\sum_{t'=t}^T r(s_{t'}, a_{t'}) \right] ds_t da_t - \int_{\mathcal{S}} \int_{\mathcal{A}} p(s_t, a_t) \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) b(s_t) ds_t da_t$$

Let us show this is 0!

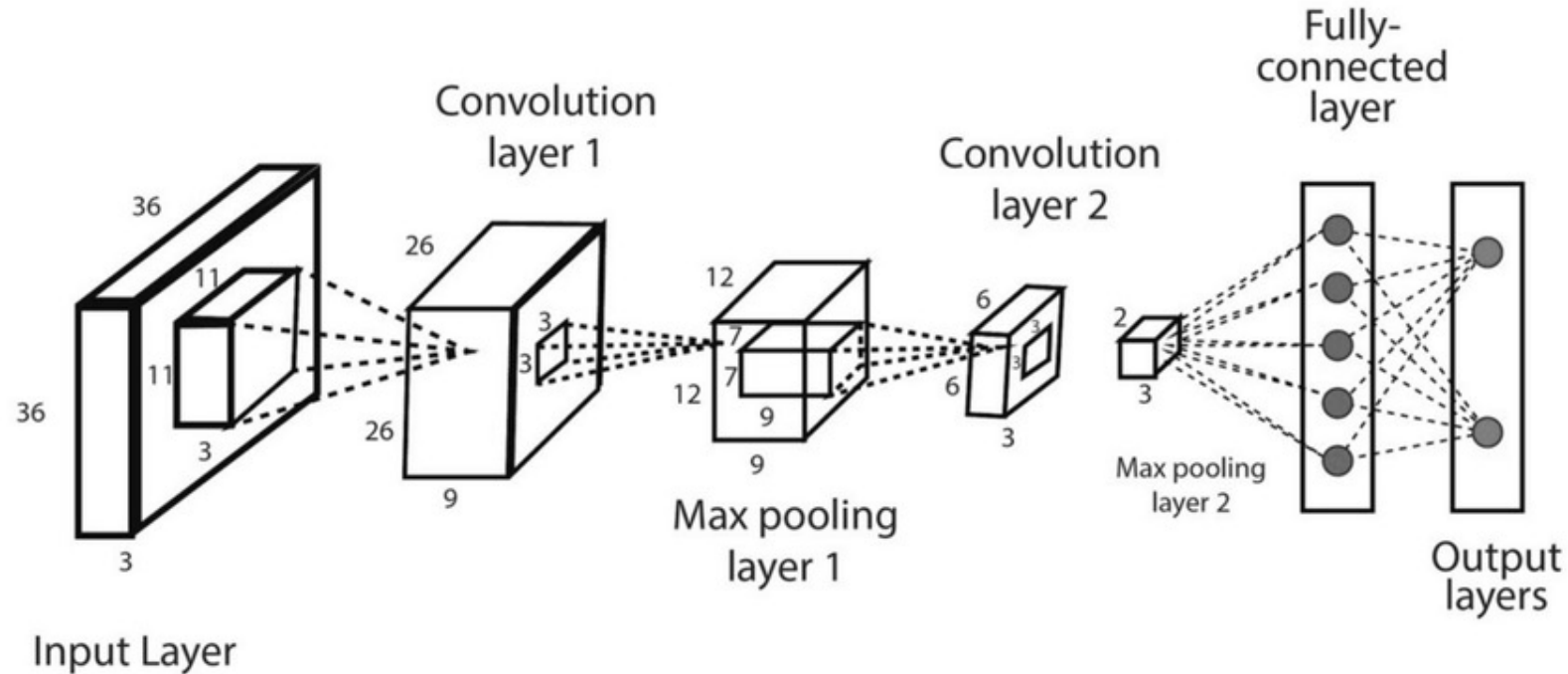
Variance Reduction with a Baseline

$$\begin{aligned}\int \int p(s_t, a_t) \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) [b(s_t)] ds_t da_t &= \int \int p(s_t) \pi_{\theta}(a_t | s_t) \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) [b(s_t)] ds_t da_t \\ &= \int p(s_t) b(s_t) \int \pi_{\theta}(a_t | s_t) \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) da_t ds_t \\ &= \int p(s_t) b(s_t) \int \nabla_{\theta} \pi_{\theta}(a_t | s_t) da_t ds_t \\ &= \int p(s_t) b(s_t) \nabla_{\theta} \int \pi_{\theta}(a_t | s_t) da_t ds_t = \int p(s_t) b(s_t) \nabla_{\theta} (1) ds_t = 0\end{aligned}$$

Unbiased!

Learning Baselines

Baselines are typically learned as deep neural nets from $\mathbb{R}^s \rightarrow \mathbb{R}^1$



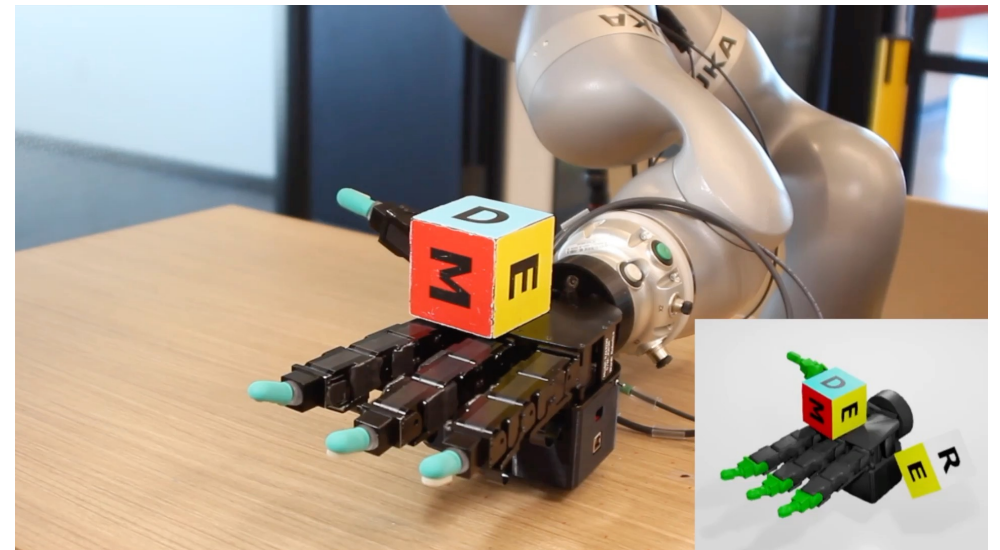
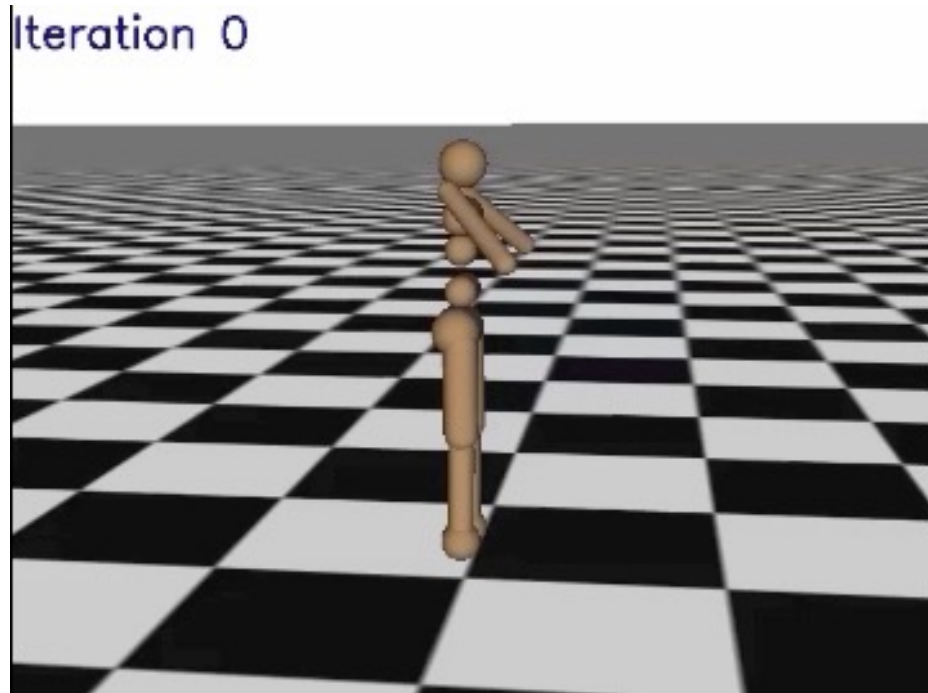
$$\frac{1}{N} \sum_{j=1}^N \left\| \hat{V}(s_t^j, a_t^j) - \sum_{t=1}^H r(s_t^j, a_t^j) \right\|$$

Minimize with Monte-carlo regression at every iteration, club with policy loss

$$A(s_t, a_t) = \sum_{t'=t}^T r(s_{t'}, a_{t'}) - V(s_t)$$

Allows us to define advantages

Policy Gradient in Action



Class Outline

State Estimation

Robotic System Design

Filtering

Localization

SLAM

Control

Feedback Control

PID Control

MPC

LQR

Planning

Search

Heuristic Search

Motion Planning

Lazy Search

Learning

Imitation Learning

Policy Gradient

Actor-Critic

Model-Based RL