

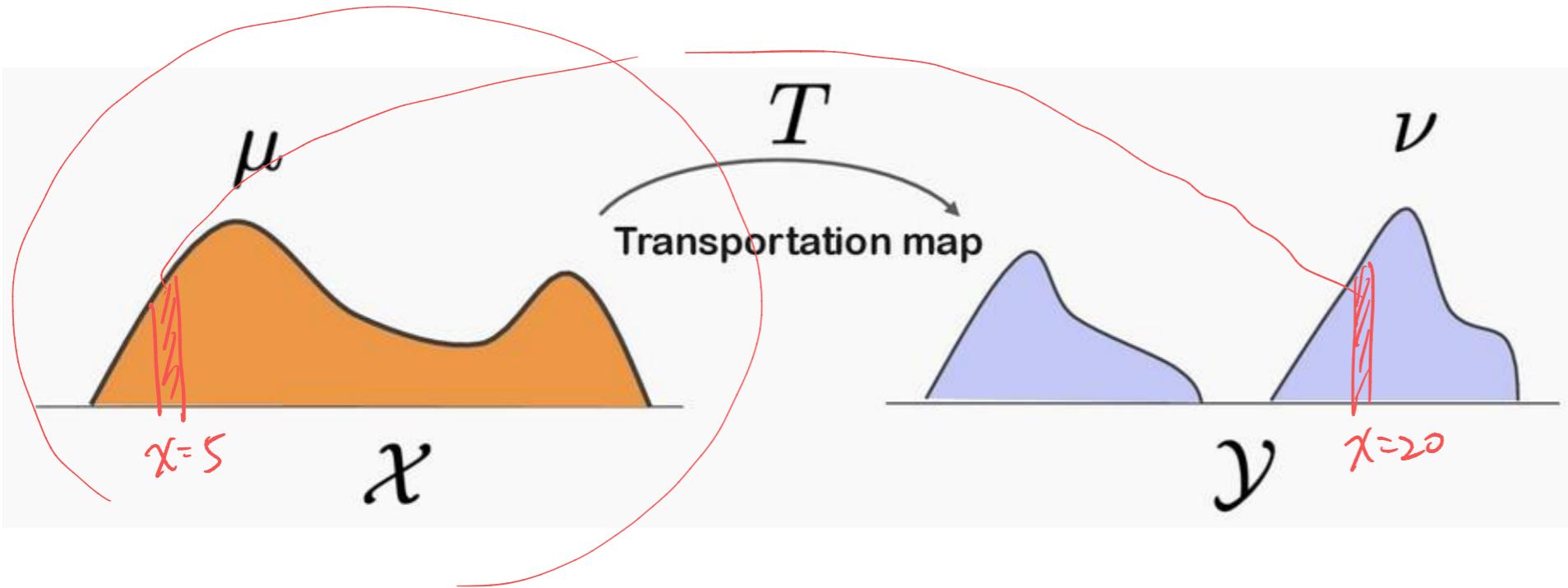
535514: Reinforcement Learning

Lecture 25 – SAC and Imitation Learning

Ping-Chun Hsieh

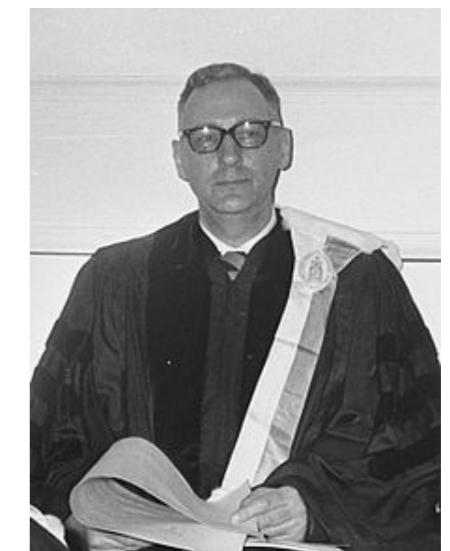
May 27, 2024

Optimal Transport



Leonid Kantorovich
(From Soviet Union)

Kantorovich received Nobel Laureate of Economics in 1975
(with Tjalling Koopmans)



Tjalling Koopmans
(From the US)

On-Policy vs Off-Policy Methods

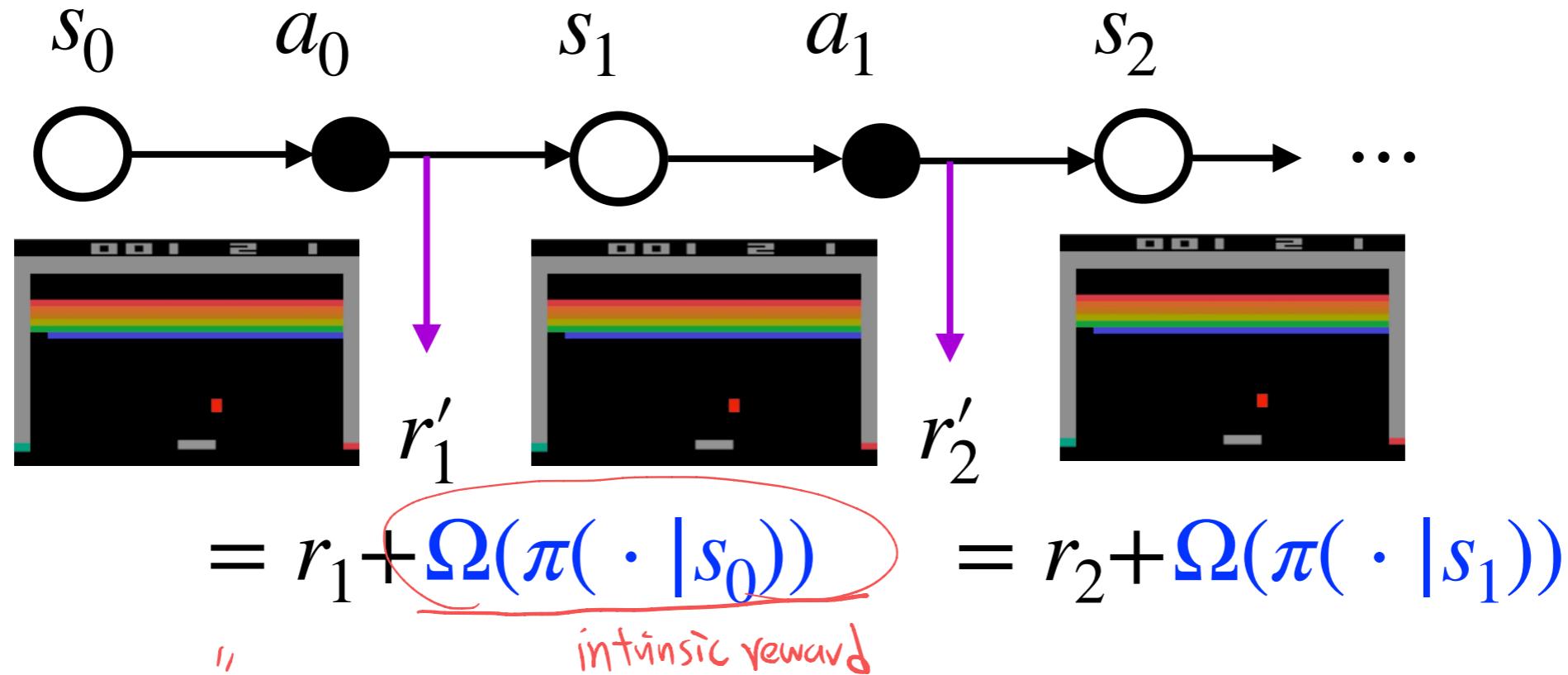
	Policy Optimization	Value-Based	Model-Based	Imitation-Based
On-Policy	Exact PG REINFORCE (w/i baseline) A2C On-policy DAC TRPO Natural PG (NPG) PPO-KL & PPO-Clip RLHF by PPO-KL	Epsilon-Greedy MC Sarsa Expected Sarsa	Model-Predictive Control (MPC) PETS	IRL GAIL IQ-Learn
Off-Policy	Off-policy DPG & DDPG Twin Delayed DDPG (TD3)	Q-learning Double Q-learning DQN & DDQN Rainbow C51 / QR-DQN / IQN Soft Actor-Critic (SAC)		

Soft Policy Iteration

1. Policy evaluation
2. Policy improvement

Review: Regularized MDPs

Regularized MDP = Standard MDP + Regularized rewards!



- ▶ A regularized MDP can be specified by $(\mathcal{S}, \mathcal{A}, P, R, \Omega, \gamma)$
- ▶ $\Omega(\cdot)$: A function that maps an *action distribution* to a *real number*

Question: How to define $Q^\pi(s, a)$ and $V^\pi(s)$ with a regularizer?

Review: Value Functions of Regularized MDPs

O
 s_t

$E \left[E \left[G_t \mid s_t = s, a_t = a; \pi \right] \right] \mid s_t = s; \pi$

	"Unregularized" MDP	"Entropy-Regularized" MDP
Return	$G_t := r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \dots$	$G_t := r_{t+1} + \gamma(r_{t+2} + \Omega(\pi(\cdot s_{t+1}))) + \gamma^2(r_{t+3} + \Omega(\pi(\cdot s_{t+2}))) + \dots$
Value function	$V^\pi(s) := \mathbb{E}[G_t s_t = s; \pi]$	$V_\Omega^\pi(s) := \mathbb{E}[G_t s_t = s; \pi] + \Omega(\pi(\cdot s))$
Q function	$Q^\pi(s, a) := \mathbb{E}[G_t s_t = s, a_t = a; \pi]$	$Q_\Omega^\pi(s, a) := \mathbb{E}[G_t s_t = a, a_t = a; \pi]$
Bellman expectation equations	$V^\pi(s) = \sum_{a \in \mathcal{A}} \pi(a s) Q^\pi(s, a)$ $Q^\pi(s, a) = R_{s,a} + \gamma \sum_{s' \in \mathcal{S}} P_{s,s'}^a V^\pi(s')$	$V_\Omega^\pi(s) = \sum_{a \in \mathcal{A}} \pi(a s) Q_\Omega^\pi(s, a) + \Omega(\pi(\cdot s))$ $Q_\Omega^\pi(s, a) = R_{s,a} + \gamma \sum_{s' \in \mathcal{S}} P_{s,s'}^a V_\Omega^\pi(s')$

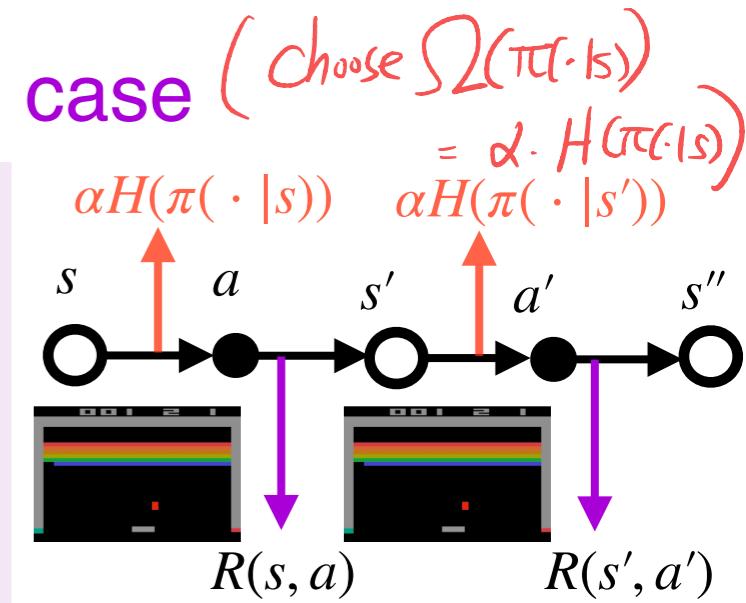
If $\Omega(\pi(\cdot | s)) \equiv \alpha \cdot H(\pi(\cdot | s))$, the value functions are called “soft functions”

Soft Policy Evaluation for Soft Q-Function

Let's extend **policy evaluation** to entropy-regularized case

$$Q_{soft}^{\pi}(s, a) = R_s^a + \gamma E_{s' \sim P(\cdot | s, a)} [V_{soft}^{\pi}(s')] \quad (\text{choose } \Sigma(\pi) \cdot Q)$$

$$(= R_s^a + \gamma E_{s' \sim P(\cdot | s, a), a' \sim \pi(\cdot | s')} [Q_{soft}^{\pi}(s', a') - \alpha \log(\pi(a' | s'))])$$



Soft policy evaluation: Find $Q_{soft}^{\pi}(s, a)$ for a policy π

$$\text{Entropy} = -\left(\sum_a \pi(a | s) \log \pi(a | s)\right)$$

1. (Optimal control) Given R, P , and a policy π :

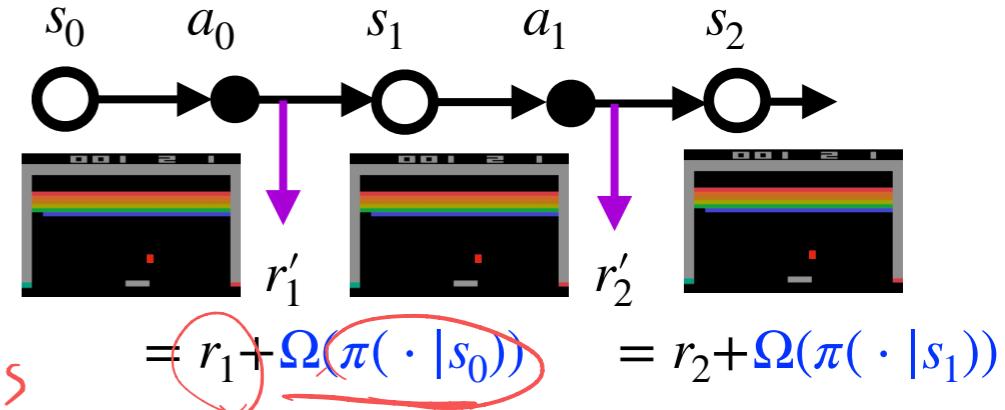
$$\text{Soft IPE operator} = \left[T_{soft}(Q) \right]_{s,a} := R_s^a + \gamma \cdot E_{\substack{s' \sim P(\cdot | s, a) \\ a' \sim \pi(\cdot | s')}} [Q(s', a') - \alpha \cdot \log(\pi(a' | s'))]$$

2. (Learning) Given a policy π (with unknown R, P):

E.g.: TD(0)

$$Q_{soft}(s, a) \leftarrow Q_{soft}(s, a) + \eta \cdot \underbrace{(\gamma + \gamma \cdot (Q(s, a) - \alpha \log \pi(a | s)))}_{\text{TD error}} \underbrace{Q_{soft}(s, a)}_{Q_{soft}(s, a)}$$

Review: Optimal Value Functions and Bellman Optimality Equations of Regularized MDPs



	Unregularized MDP	Regularized MDP
Bellman optimality equations	$V^*(s) := \max_{\pi \in \Pi} V^\pi(s)$ $Q^*(s, a) := \max_{\pi \in \Pi} Q^\pi(s, a)$ <p style="text-align: center;">$\sum_{a \in \mathcal{A}} \pi(a s) \cdot R_s^a$</p>	$V_\Omega^*(s) := \max_{\pi \in \Pi} V_\Omega^\pi(s)$ $Q_\Omega^*(s, a) := \max_{\pi \in \Pi} Q_\Omega^\pi(s, a)$
Bellman optimality equations	$V^*(s) = \max_{a \in \mathcal{A}} R_s^a + \gamma P_s^a V^*$ $= \max_{\pi \in \Pi} R_s^\pi + \gamma P_s^\pi V^*$ $Q^*(s, a) = R_s^a + \gamma E_{s' \sim P(\cdot s, a)}[V^*(s')]$	$V_\Omega^*(s) = \max_{\pi \in \Pi} R_s^\pi + \gamma P_s^\pi V_\Omega^*$ $Q_\Omega^*(s, a) = R_s^a + \gamma E_{s' \sim P(\cdot s, a)}[V_\Omega^*(s')]$ $= R_s^a + \gamma \cdot E \left[\max_{\substack{\pi \in \Pi \\ s' \sim P(\cdot s, a)}} R_{s'}^{\pi} + \gamma P_{s'}^{\pi} V_\Omega^* \right]$

Soft Policy Improvement

Bellman optimality equation

$$Q_{soft}^*(s, a) = R_s^a + \gamma E_{s' \sim P(\cdot | s, a)} [V_{soft}^*(s')] \\ (= R_s^a + \gamma E_{s' \sim P(\cdot | s, a)} [\max_{\pi} \{ \langle \pi(\cdot | s'), Q_{soft}^{\pi}(s', \cdot) \rangle + \alpha H(\pi(\cdot | s')) \}])$$

Soft policy improvement: Given π_k , improve the policy by

$$\pi_{k+1}(\cdot | s) = \arg \max_{\pi} \left\{ \langle \pi(\cdot | s), Q_{soft}^{\pi_k}(s, \cdot) \rangle + \alpha H(\pi(\cdot | s)) \right\}$$

$$\sum_{a \in A} \pi(a | s) \cdot Q_{soft}^{\pi_k}(s, a)$$

If $\alpha=0$:

$$\bar{\pi}_{k+1}(\cdot | s) = \text{argmax}_{\pi} \sum_{a \in A} \pi(a | s) \cdot Q_{soft}^{\pi_k}(s, a)$$

$$\begin{bmatrix} 1 \\ 3 \\ 2 \end{bmatrix} \begin{bmatrix} \pi(a_1) \\ \pi(a_2) \\ \pi(a_3) \end{bmatrix}$$

Solution to Soft Policy Improvement

- **Theorem:** Under soft policy iteration, we have

$$\pi_{k+1}(\cdot | s) = \arg \max_{\pi} \left\{ \langle \pi(\cdot | s), Q_{soft}^{\pi_k}(s, \cdot) \rangle + \alpha H(\pi(\cdot | s)) \right\}$$

$$\Rightarrow \frac{\exp\left(\frac{1}{\alpha}Q_{soft}^{\pi_k}(s, \cdot)\right)}{\sum_{a \in \mathcal{A}} \exp\left(\frac{1}{\alpha}Q_{soft}^{\pi_k}(s, a)\right)}$$

temperature coefficient

$\Rightarrow \begin{cases} 1. \alpha=0: \text{Greedy policy} \\ 2. \alpha \rightarrow \infty: \text{Uniformly random policy} \end{cases}$

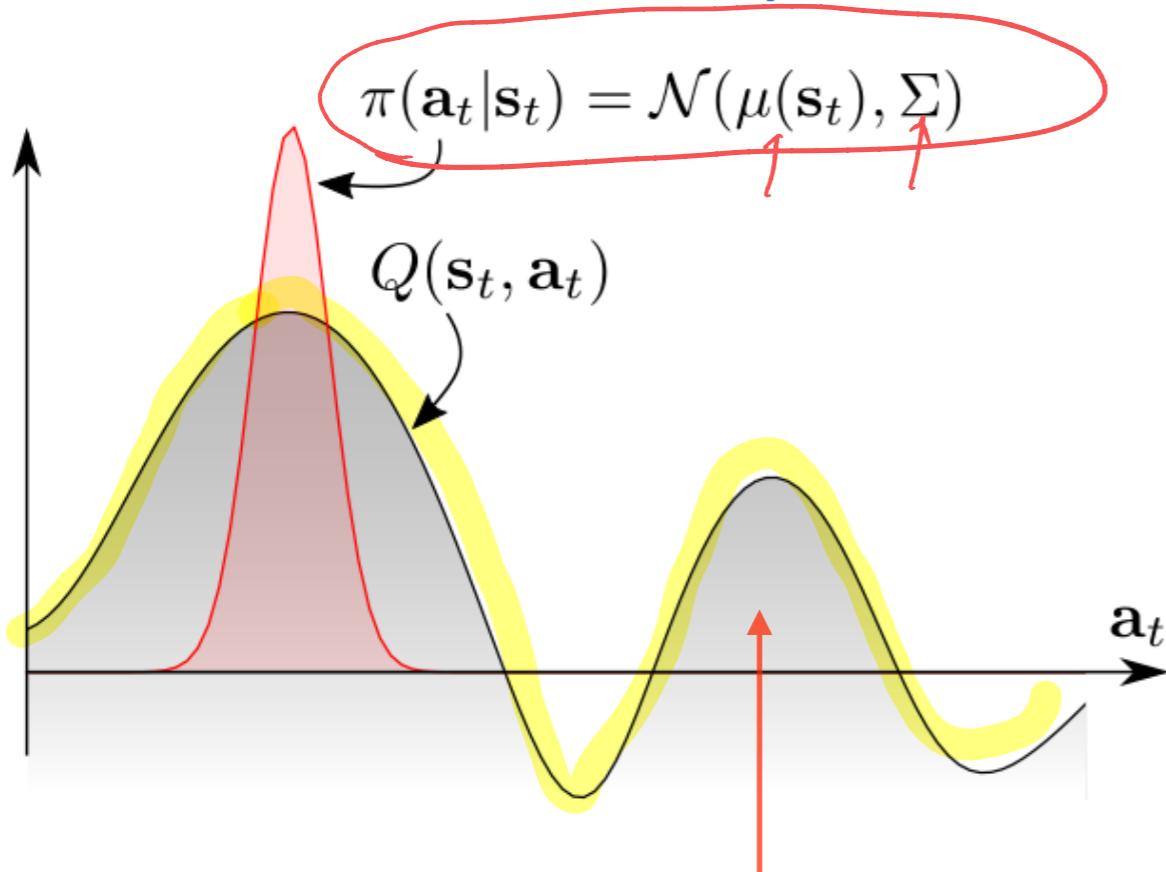
- **Question:** Could you explain why this is called “soft” policy improvement?

| Direct Preference Optimization (DPO)

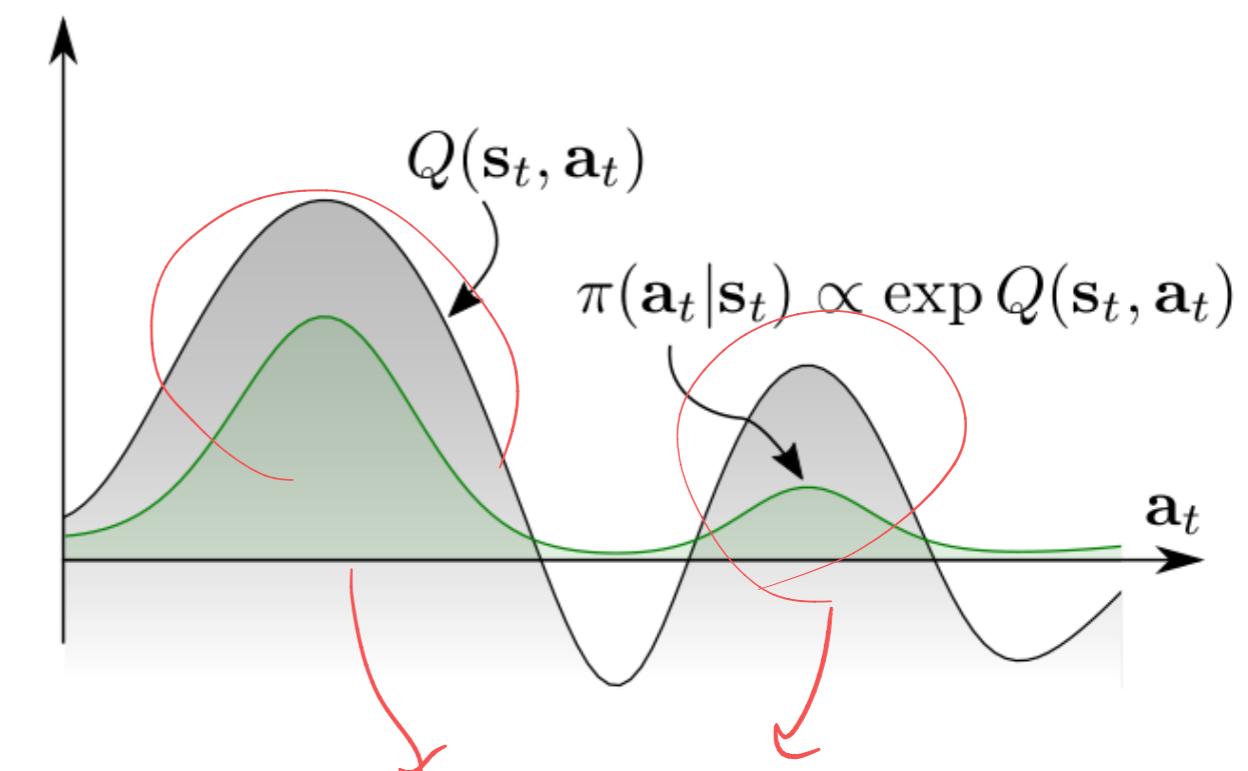
Why is Soft Policy Improvement a Good Idea?

Suppose we'd like to use stochastic policies

Standard Gaussian policies



Energy-based policies



Unimodal policies completely
ignore this part



Soft Policy Iteration (Soft PI)

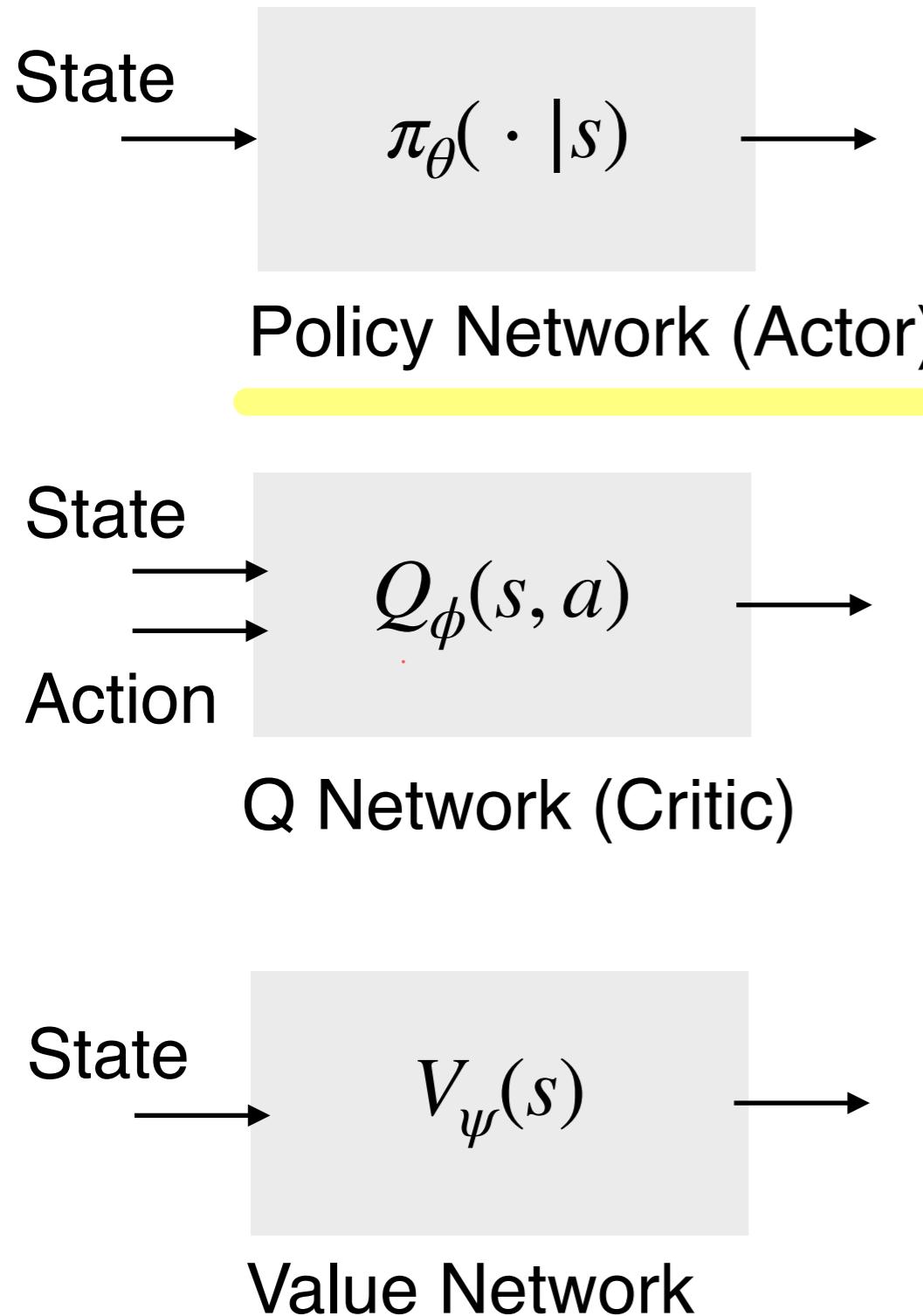
Soft Policy Iteration

1. Initialize $k = 0$ and set $\pi_0(\cdot | s)$ arbitrarily for all states
2. While \underline{k} is zero or $\underline{\pi_k} \neq \underline{\pi_{k-1}}$:
 - ▶ Derive $V_{soft}^{\pi_k}$ and $Q_{soft}^{\pi_k}$ via soft policy evaluation
 - ▶ Derive π_{k+1} by ~~greedy~~ soft policy improvement:

$$\pi_{k+1}(\cdot | s) = \arg \max_{\pi} \left\{ \langle \pi(\cdot | s), Q_{soft}^{\pi_k}(s, \cdot) \rangle + H(\pi(\cdot | s)) \right\}$$

Soft Actor-Critic (SAC)

Soft Actor-Critic: The “Learning” Version of Soft-PI



Actor Loss

$$L_\pi(\theta) = \mathbb{E}_{s \sim D} [D_{KL}(\pi_\theta(\cdot | s) \| \exp(Q_{\bar{\phi}}(s, \cdot)) / Z_{\bar{\phi}}(s))]$$

$Z_{\bar{\phi}}(s) = \int_a \exp(Q_{\bar{\phi}}(s, a)) \pi(a|s) da$

Critic Loss

$$= \mathbb{E}_{s \sim D, a \sim \pi_\theta} [\log(\pi_\theta(a|s)) - Q_{\bar{\phi}}(s, a) + \log Z_{\bar{\phi}}(s)]$$

$L_Q(\phi) = \mathbb{E}_{(s, a) \in D} \left[\frac{1}{2} (Q_\phi(s, a) - (\mathbb{E}_{r, s'} [r + \gamma V_{\bar{\psi}}(s') | s, a]))^2 \right]$

Value Loss

$$\nabla_\phi L(\phi) = \mathbb{E}_{(s, a)} [(Q_\phi - (\dots)). \nabla Q_\phi(s, a)]$$

(by soft Bellman expectation equation)

Value Loss

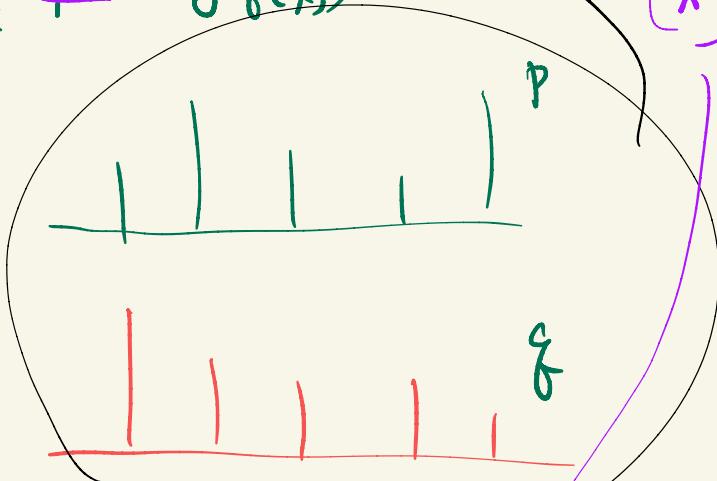
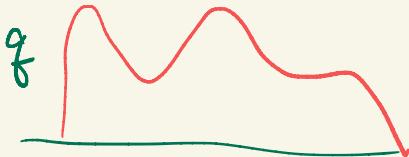
$$L(\psi) = \mathbb{E}_{s \sim D} \left[\frac{1}{2} (V_\psi(s) - \mathbb{E}_{a \sim \pi_{\bar{\theta}}(\cdot | s)} [Q_{\bar{\phi}}(s, a) - \log \pi_{\bar{\theta}}(a | s)])^2 \right]$$

(by chain rule)

Value Loss

$$\nabla_\psi L_V(\psi) =$$

$$D_{KL}(P||q) = \sum_x p(x) \cdot \log\left(\frac{p(x)}{q(x)}\right) \leftarrow = E_{x \sim p} [\log p(x) - \log q(x)]$$



$$D_{KL}(P||q) = \int p(x) \log\left(\frac{p(x)}{q(x)}\right) dx = E_{x \sim p} [\log p(x) - \log q(x)]$$

Inherent Difficulty in Finding $\nabla L_\pi(\theta)$

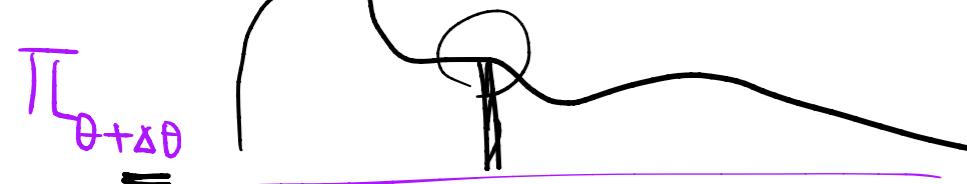
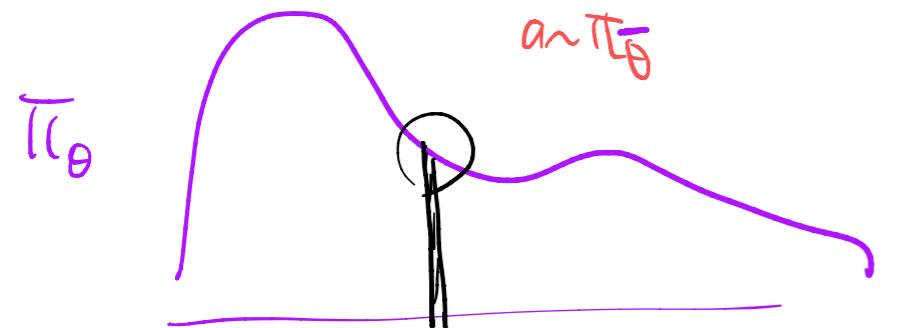
SAC policy loss :

$$L_\pi(\theta) = \mathbb{E}_{s \sim D, a \sim \pi_\theta} \left[\log \pi_\theta(a|s) - Q_{\bar{\phi}}(s, a) + \cancel{\log Z_{\bar{\phi}}(s)} \right]$$

Issue: Is it easy to directly compute the gradient of this KL divergence?

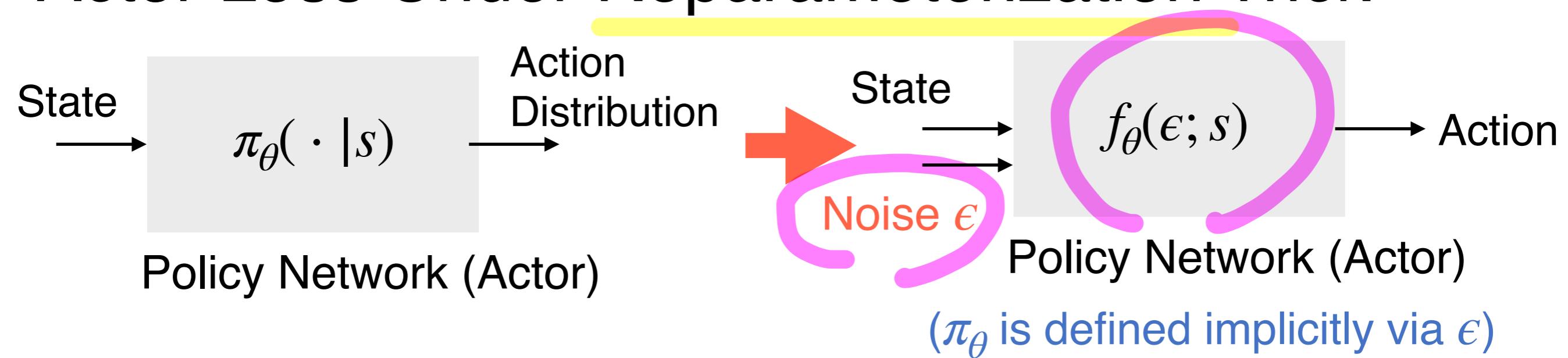
$$\nabla_\theta L_\pi(\theta) = \nabla_\theta \mathbb{E}_{s \sim D, a \sim \pi_\theta} \left[\log \pi_\theta(a|s) - Q_{\bar{\phi}}(s, a) + \cancel{\log Z_{\bar{\phi}}(s)} \right]$$

? $\nabla_\theta \left(\frac{1}{|B|} \sum_{\substack{(s, a) \in \\ s \sim D \\ a \sim \pi_{\bar{\phi}}}} \log \pi_\theta(a|s) - Q_{\bar{\phi}}(s, a) \right)$



$$\nabla \sum \pi_\theta(a|s) (\log \pi_\theta - \dots)$$

Actor Loss Under Reparameterization Trick



Reparamterization Trick:

$$L_\pi(\theta) = \mathbb{E}_{s \sim D, \epsilon \sim G} [\log \pi_\theta(f_\theta(\epsilon; s) | s) - Q_{\bar{\phi}}(s, f_\theta(\epsilon; s))]$$

Takeaway:

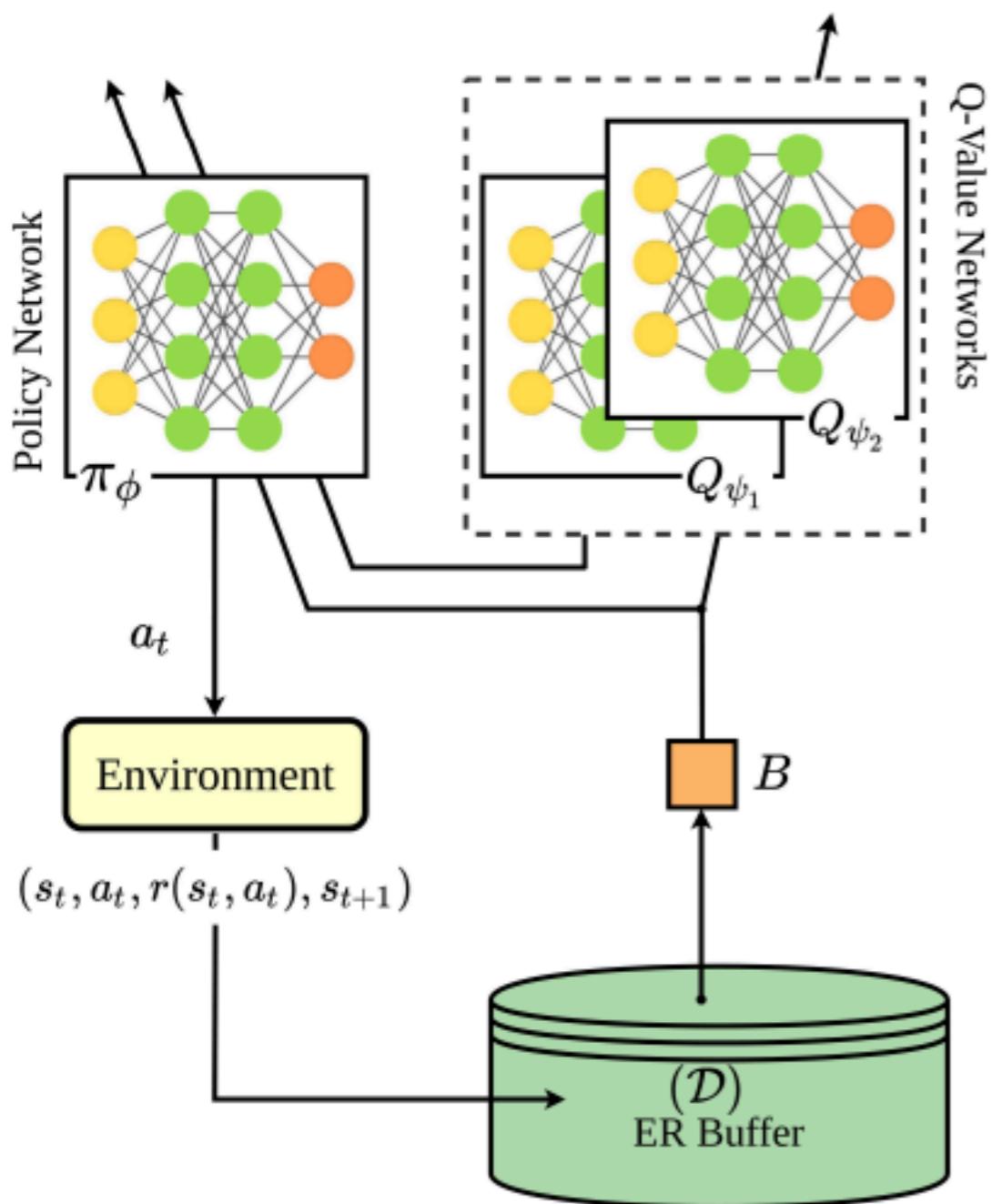
$$\nabla_\theta L_\pi(\theta) = \nabla_\theta \left(\mathbb{E}_{s \sim D, \epsilon \sim G} [\log \pi_\theta(f_\theta(\epsilon; s) | s) - Q_{\bar{\phi}}(s, f_\theta(\epsilon; s))] \right)$$

Reparametrization is
a "generative" approach

that removes the " θ "
in the \mathbb{E}

$$\nabla_\theta \frac{1}{|B|} \sum_{(s, \epsilon) \in B} (\log \pi_\theta(f_\theta(\epsilon; s) | s) - Q_{\bar{\phi}}(s, f_\theta(\epsilon; s)))$$

Architecture of SAC



1. Clipped double Q networks as TD3

2. Gaussian policies

3. Experience replay buffer for off-policy learning

Imitation Learning

Imitation Learning: 2 Major Paradigms

- ▶ Suppose we are given *expert demonstrations*.
How to learn from them?

1. Direct imitation learning

- Copy the **actions** of the expert
- No reasoning about the outcomes of actions

2. Human imitation learning

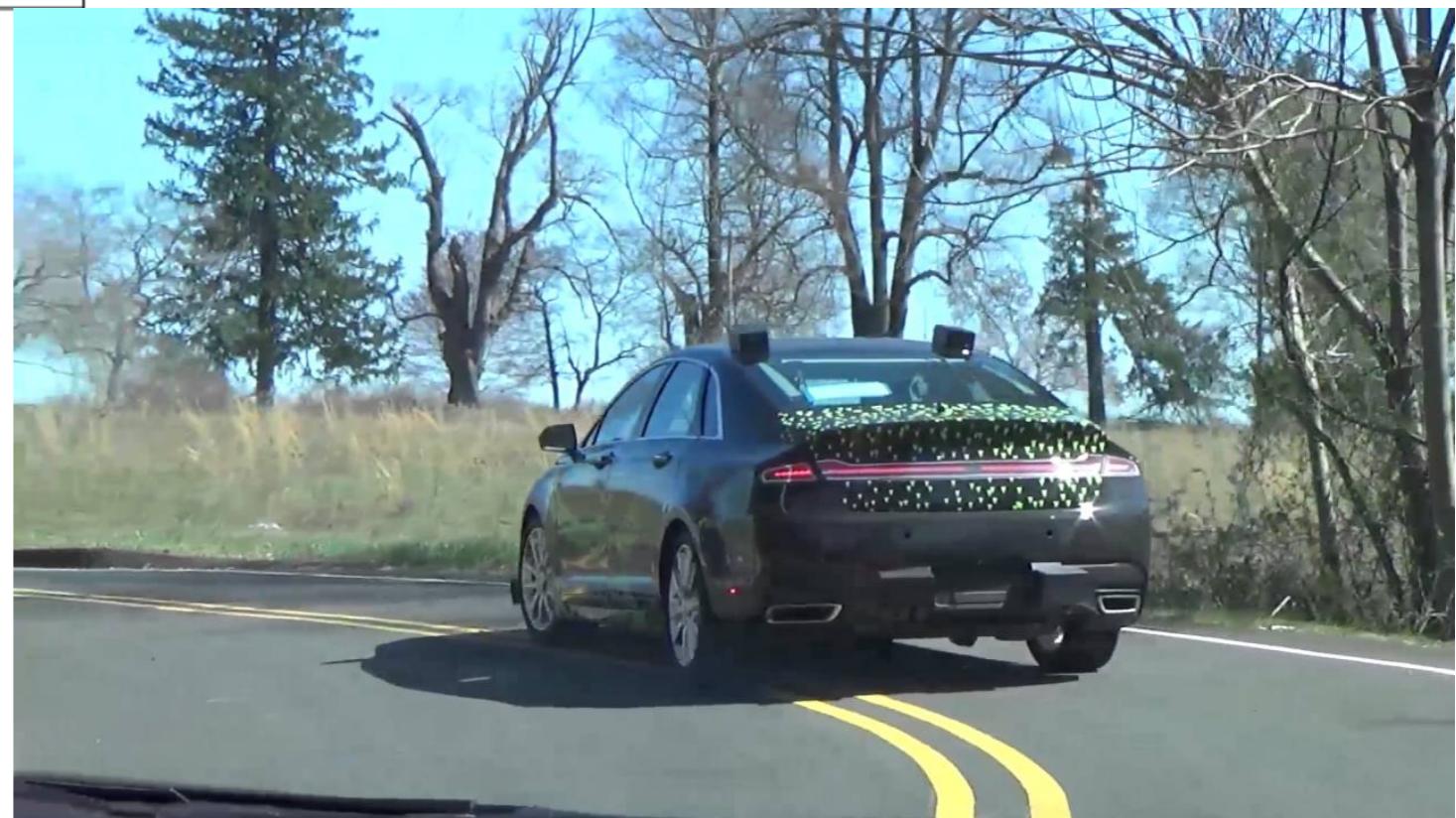
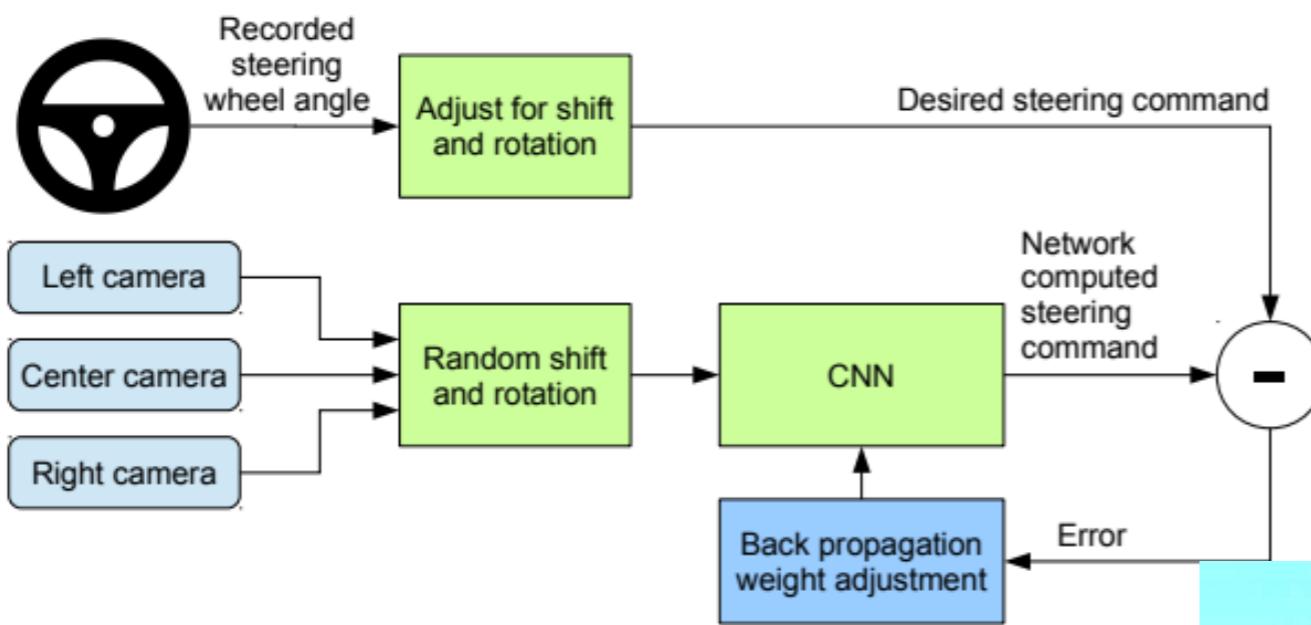
- Copy the **intent** of the expert
- May take very different actions from the expert



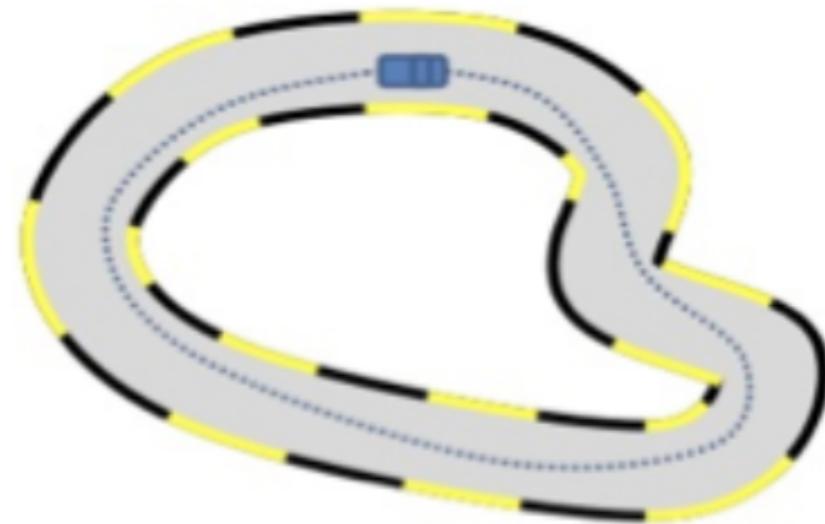
Inverse RL!

Direct Imitation Learning

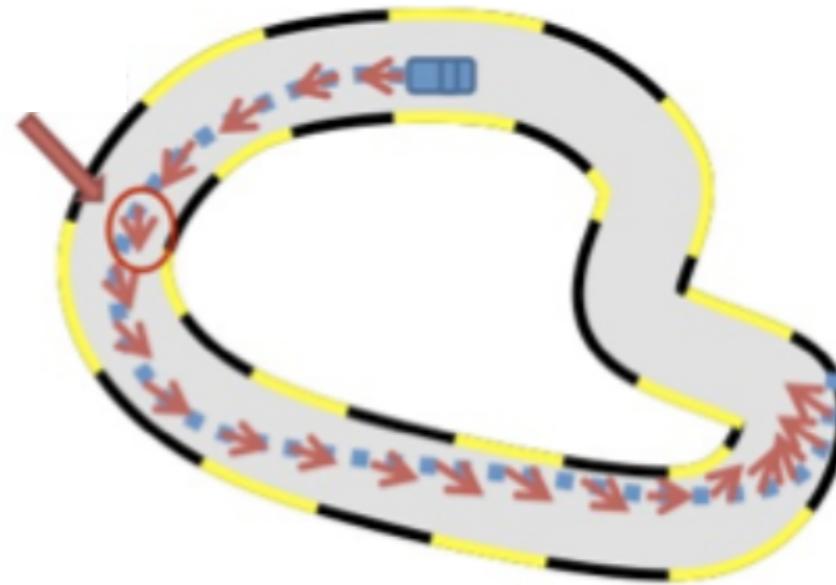
- ▶ Example: Self-driving cars



Direct Imitation Learning: Out-of-Distribution Issue



Expert Trajectories
(Training Distributions)

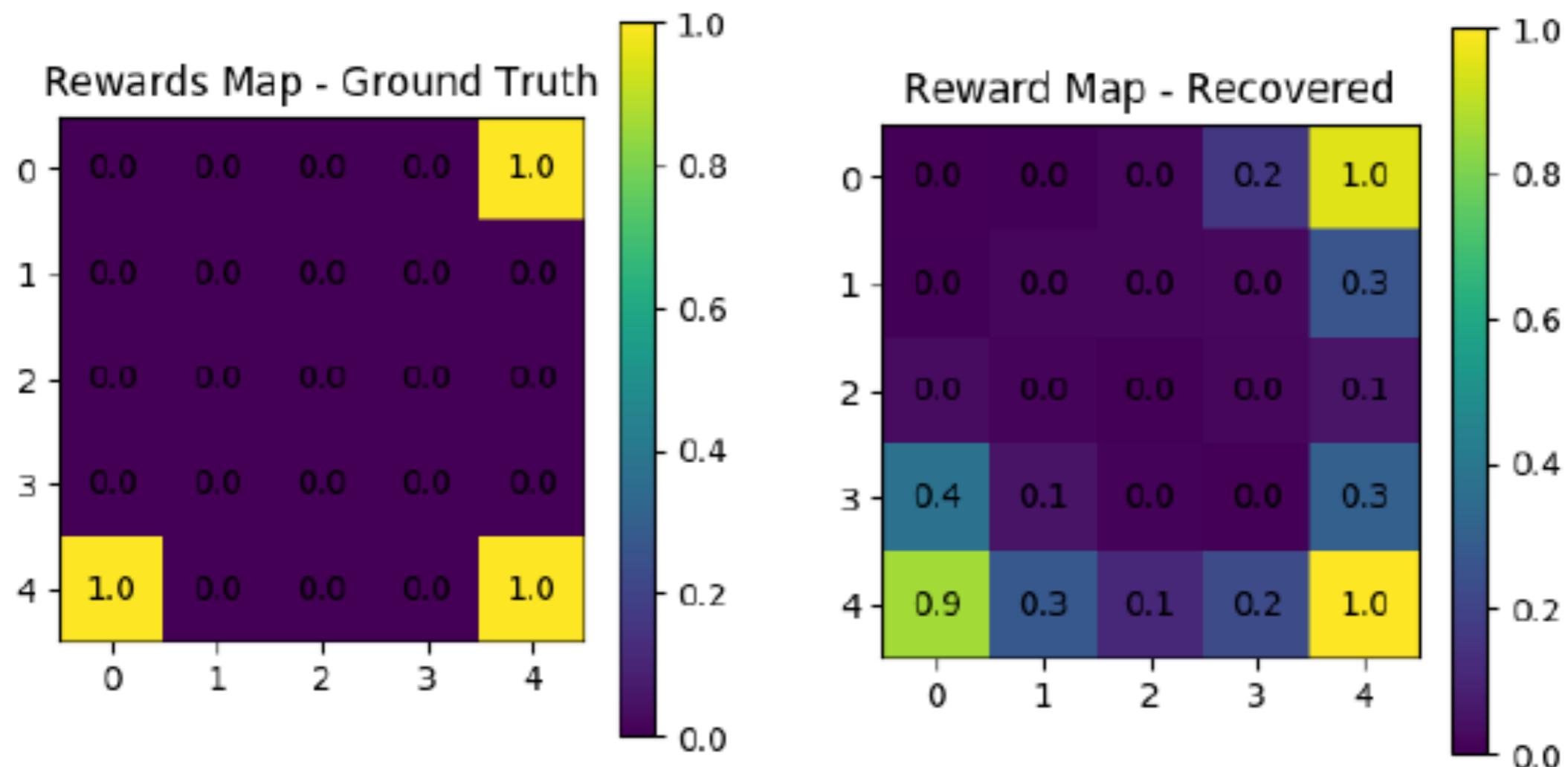


Direct Imitation Learning
(Testing Distributions)

Makes mistakes, enter new states
Cannot recover from new states

Inverse RL

- ▶ **Example:** Reward recovery in Gridworld

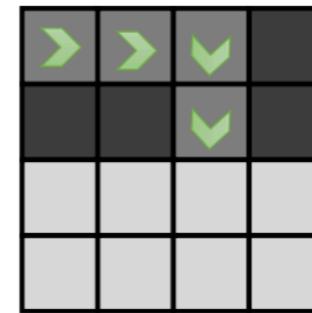
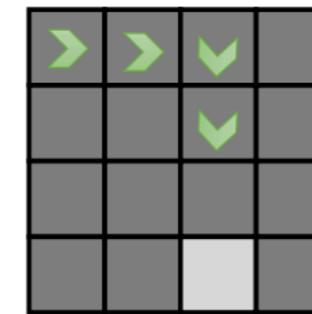
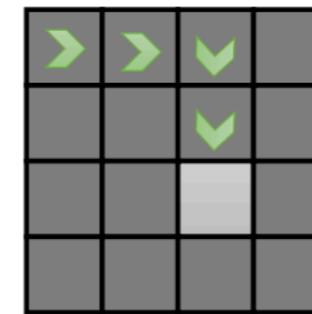
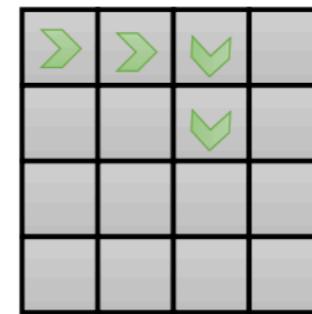
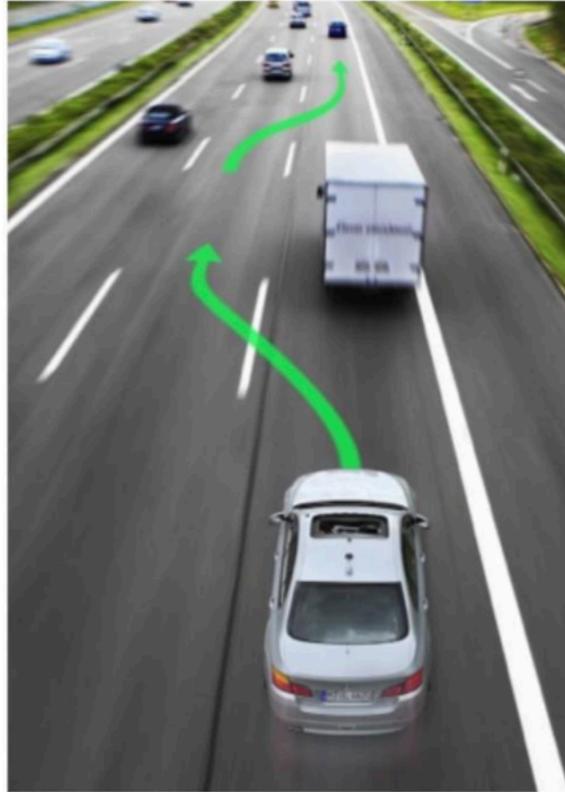


Inverse RL (Informal)

- ▶ Suppose the agent is in an MDP (S, A, P, γ)
- ▶ Suppose we are given **expert demonstrations** (under some unknown policy π_e)
- ▶ **Goal:** Infer the reward function R behind the expert actions solely from expert demonstrations (and thereafter learn a good policy)

First Attempt: Infer Rewards from Demonstrations

- ▶ **Example:** Human driving



Typically, “reward inference” is an underspecified problem

(Multiple reward functions can explain the same behavior)

What's reward $R(s, a)$?

Reward identifiability issue!

Rethinking Imitation?

Recall: *Occupancy measure* (or *discounted state visitation*)

$$d_\mu^\pi(s) := (1 - \gamma) \mathbb{E}_{s_0 \sim \mu} \left[\sum_{t=0}^{\infty} \gamma^t P(s_t = s \mid s_0, \pi) \right]$$

$$d_\mu^\pi(s, a) := (1 - \gamma) \mathbb{E}_{s_0 \sim \mu} \left[\sum_{t=0}^{\infty} \gamma^t P(s_t = s, a_t = a \mid s_0, \pi) \right]$$

(Q1) If $\pi_\theta = \pi_e$, then do we have $d_\mu^{\pi_\theta}(s, a) = d_\mu^{\pi_e}(s, a)$?

(Q2) If $d_\mu^{\pi_\theta}(s, a) = d_\mu^{\pi_e}(s, a)$, then do we have $V^{\pi_\theta}(\mu) = V^{\pi_e}(\mu)$?

(Q3) If $d_\mu^{\pi_\theta}(s, a) = d_\mu^{\pi_e}(s, a)$, then do we have $\pi_\theta = \pi_e$?

About (Q3): A Bijection Theorem

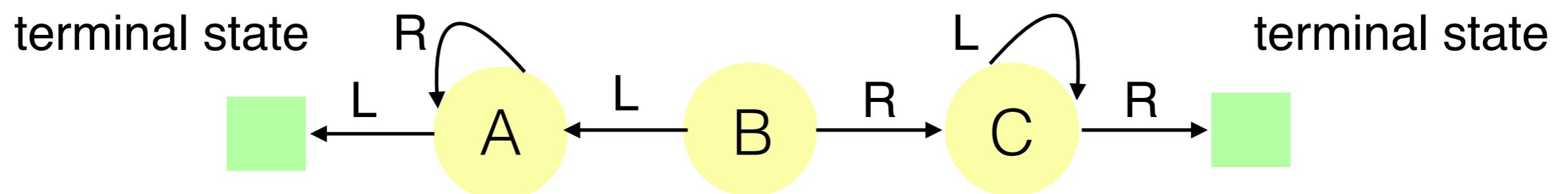
Theorem [Syed et al., 2008]: For any valid discounted state visitation distribution $d^\pi(s, a)$, define a policy $\pi'(a | s) := d^\pi(s, a) / \sum_{a' \in A} d^\pi(s, a')$.

Then, we have $d^{\pi'}(s, a) = d^\pi(s, a)$, for all (s, a) .

(In other words, the mapping from $d^\pi \rightarrow \pi$ is a bijection)

- However, the above Bijection Theorem does NOT implies that (Q3).
- Regarding (Q3), Bijection Theorem only implies that $\pi_\theta(\cdot | s) = \pi_e(\cdot | s)$ at those states with $d^{\pi_e}(s) > 0$

Example:



Inverse RL: Occupancy Measure Matching

Brian Ziebart et al., Maximum entropy inverse reinforcement learning, AAAI 2008

Jonathan Ho and S. Ermon, Generative adversarial imitation learning, NIPS 2016

Xiao et al., Wasserstein Adversarial Imitation Learning, NeurIPS 2019

Garg et al., IQ-Learn: Inverse soft-Q Learning for Imitation, NeurIPS 2021

Occupancy Measure Matching: Formulation

Recall: *Occupancy measure* (or *discounted state visitation*)

$$d_\mu^\pi(s) := (1 - \gamma) \mathbb{E}_{s_0 \sim \mu} \left[\sum_{t=0}^{\infty} \gamma^t P(s_t = s \mid s_0, \pi) \right]$$

$$d_\mu^\pi(s, a) := (1 - \gamma) \mathbb{E}_{s_0 \sim \mu} \left[\sum_{t=0}^{\infty} \gamma^t P(s_t = s, a_t = a \mid s_0, \pi) \right]$$

Claim: $V^\pi(\mu) = \sum_{(s,a)} d_\mu^\pi(s, a) R(s, a)$ (Why?)

Occupancy measure matching:

Find a policy π such that $d_\mu^\pi(s, a) = d^{\pi_e}(s, a), \quad \forall (s, a)$

Occupancy measure matching implies $V^\pi(\mu) = V^{\pi_e}(\mu)$

- ▶ **Question:** Is $d_\mu^\pi(s, a)$ easy to parameterize?

(Direct) Occupancy Measure Matching (OMM)

$$\min_{\pi \in \Pi} L(\pi) := D(d_\mu^\pi, d_\mu^{\pi_e})$$

($D(\cdot, \cdot)$ is some distance)
(d_μ^π could be hard to express!)



Dual of each other!

$$\max_{R \in \mathcal{R}} \min_{\pi \in \Pi} \underbrace{\left[\left(E_{d_\mu^{\pi_e}}[R(s, a)] - E_{d_\mu^\pi}[R(s, a)] \right) \right]}_{:=L(\pi, R)}$$

OR

$$\min_{\pi \in \Pi} \max_{R \in \mathcal{R}} \underbrace{\left[\left(E_{d_\mu^{\pi_e}}[R(s, a)] - E_{d_\mu^\pi}[R(s, a)] \right) \right]}_{:=L(\pi, R)}$$

(Easier for training!)

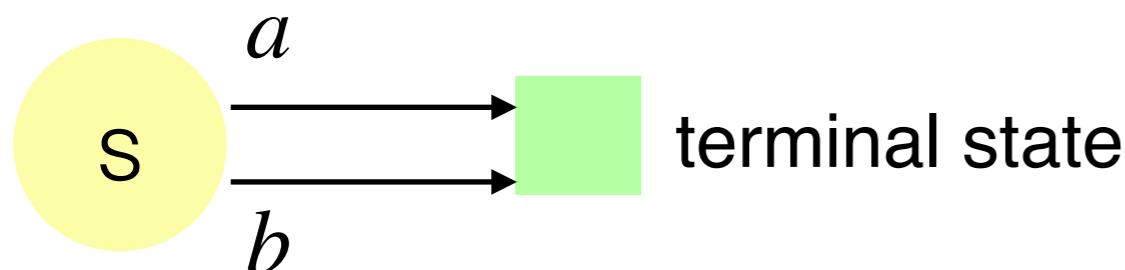
Apprenticeship Learning (APPLE)

A Motivating Example: Connecting OMM & APPLE

$$\min_{\pi \in \Pi} L(\pi) := D(d_\mu^\pi, d_\mu^{\pi_e})$$

$$\min_{\pi \in \Pi} \max_{R \in \mathcal{R}} \underbrace{\left[(E_{d_\mu^{\pi_e}}[R(s, a)] - E_{d_\mu^\pi}[R(s, a)]) \right]}_{:=L(\pi, R)}$$

Consider a simple 1-state, 2-action MDP



Suppose $\mathcal{R} = \mathbb{R}^2$

$$\pi_e(a|s) = \pi_e(b|s) = 0.5$$

Let's write down $R \in \mathcal{R}$ that maximizes $L(\pi, R)$ under a fixed π

For (s, a) with $d_\mu^\pi(s, a) > d_\mu^{\pi_e}(s, a)$:

For (s, a) with $d_\mu^\pi(s, a) < d_\mu^{\pi_e}(s, a)$:

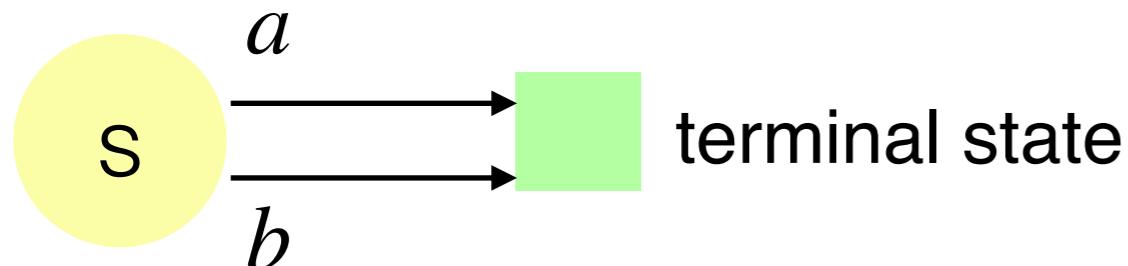
For (s, a) with $d_\mu^\pi(s, a) = d_\mu^{\pi_e}(s, a)$:

A Motivating Example: Connecting OMM & APPLE

$$\min_{\pi \in \Pi} L(\pi) := D(d_\mu^\pi, d_\mu^{\pi_e})$$

$$\min_{\pi \in \Pi} \max_{R \in \mathcal{R}} \underbrace{\left[(E_{d_\mu^{\pi_e}}[R(s, a)] - E_{d_\mu^\pi}[R(s, a)]) \right]}_{:=L(\pi, R)}$$

Consider a simple 1-state, 2-action MDP



Suppose $\mathcal{R} = \mathbb{R}^2$

$$\pi_e(a|s) = \pi_e(b|s) = 0.5$$

Nice Property: Under $\mathcal{R} = \mathbb{R}^2$, the corresponding metric D is

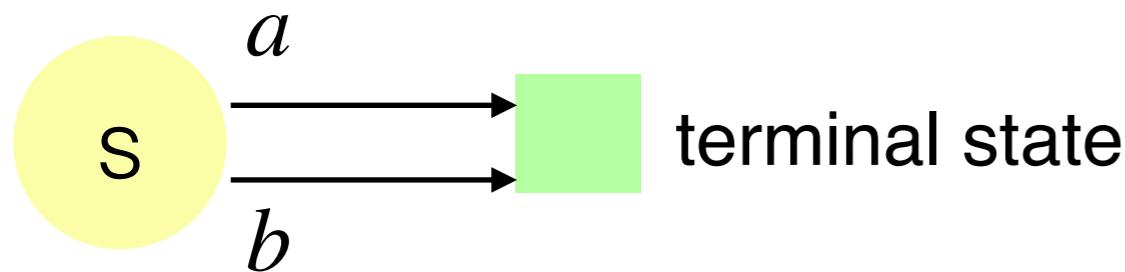
$$D(d_\mu^\pi, d_\mu^{\pi_e}) = \begin{cases} 0, & \text{if } d_\mu^\pi(s, a) = d_\mu^{\pi_e}(s, a), \forall (s, a) \\ \infty, & \text{otherwise} \end{cases}$$

A Motivating Example: Connecting OMM & APPLE (Cont.)

$$\min_{\pi \in \Pi} L(\pi) := D(d_\mu^\pi, d_\mu^{\pi_e})$$

$$\min_{\pi \in \Pi} \max_{R \in \mathcal{R}} \underbrace{\left[(E_{d_\mu^{\pi_e}}[R(s, a)] - E_{d_\mu^\pi}[R(s, a)]) \right]}_{:=L(\pi, R)}$$

Consider a simple 1-state, 2-action MDP



Suppose $\mathcal{R} = \left\{ R \in \mathbb{R}^2 \mid \|R\|_\infty \leq 1 \right\}$

$$\pi_e(a|s) = \pi_e(b|s) = 0.5$$

Let's write down $R \in \mathcal{R}$ that maximizes $L(\pi, R)$ under a fixed π

For (s, a) with $d_\mu^\pi(s, a) > d_\mu^{\pi_e}(s, a)$:

For (s, a) with $d_\mu^\pi(s, a) < d_\mu^{\pi_e}(s, a)$:

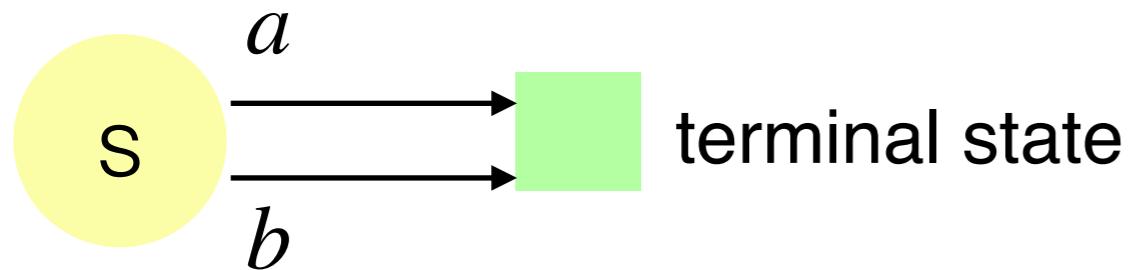
For (s, a) with $d_\mu^\pi(s, a) = d_\mu^{\pi_e}(s, a)$:

A Motivating Example: Connecting OMM & APPLE (Cont.)

$$\min_{\pi \in \Pi} L(\pi) := D(d_\mu^\pi, d_\mu^{\pi_e})$$

$$\min_{\pi \in \Pi} \max_{R \in \mathcal{R}} \underbrace{\left[(E_{d_\mu^{\pi_e}}[R(s, a)] - E_{d_\mu^\pi}[R(s, a)]) \right]}_{:=L(\pi, R)}$$

Consider a simple 1-state, 2-action MDP



Suppose $\mathcal{R} = \{R \in \mathbb{R}^2 \mid \|R\|_\infty \leq 1\}$

$$\pi_e(a|s) = \pi_e(b|s) = 0.5$$

Nice Property: Under $\mathcal{R} = \{R \in \mathbb{R}^2 \mid \|R\|_\infty \leq 1\}$, the metric D is

$$D(d_\mu^\pi, d_\mu^{\pi_e}) = \sum_{(s,a)} |d_\mu^\pi(s, a) - d_\mu^{\pi_e}(s, a)|$$

(usually called “*total variation distance*”)

How to choose \mathcal{R} to get some widely-used D ?

Example #1: Wasserstein Metric and APPLE

$$\min_{\pi \in \Pi} L(\pi) := W(d_\mu^\pi, d_\mu^{\pi_e})$$

(Wasserstein)

$$\leftrightarrow \min_{\pi \in \Pi} \max_{R \in \mathcal{R}} \underbrace{\left[\left(E_{d_\mu^{\pi_e}}[R(s, a)] - E_{d_\mu^\pi}[R(s, a)] \right) \right]}_{:=L(\pi, R)}$$

$$\text{where } \mathcal{R} = \left\{ R \in \mathbb{R}^{\mathcal{S} \times \mathcal{A}} \mid \text{Lip}(R) \leq 1 \right\}$$

This is also known as the *Kantorovich-Rubenstein duality*

Wasserstein Metric

Metric for random vectors

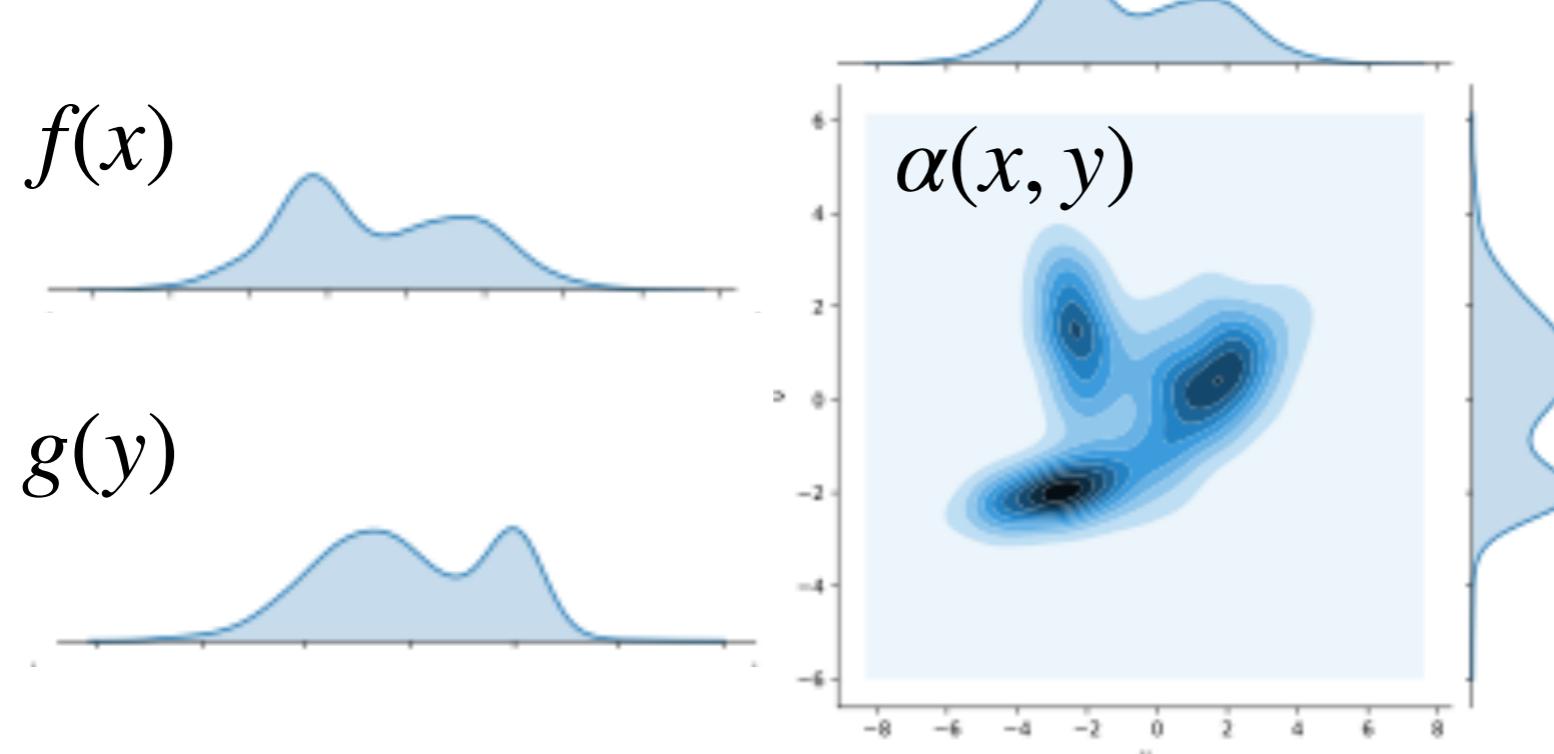
- ▶ $U : \Omega \rightarrow \mathbb{R}^d$: a random vector from the sample space Ω to \mathbb{R}^d
- ▶ For $1 \leq p < \infty$: $\|U\|_p := \left(\mathbb{E}[\|U(\omega)\|_p^p] \right)^{\frac{1}{p}}$
- ▶ **Wasserstein Metric**: For two CDFs F, G over the reals, the Wasserstein metric is defined as

$$d_p(F, G) := \inf_{(U,V): U \sim F, V \sim G} \|U - V\|_p$$

- ▶ Infimum is taken over all joint distributions of random variables (U, V) , whose marginal distributions are F, G

Intuition Behind Wasserstein Metric

- Also known as: optimal transport problem or earth mover's distance
- Given two density $f(x)$, $g(x)$ and a cost function $c(x, y)$ of moving mass from x to y , what is **minimum cost of transforming from $f(x)$ to $g(y)$?**



Minimum cost

$$C^* := \inf_{\alpha} \int c(x, y) \alpha(x, y) dx dy$$

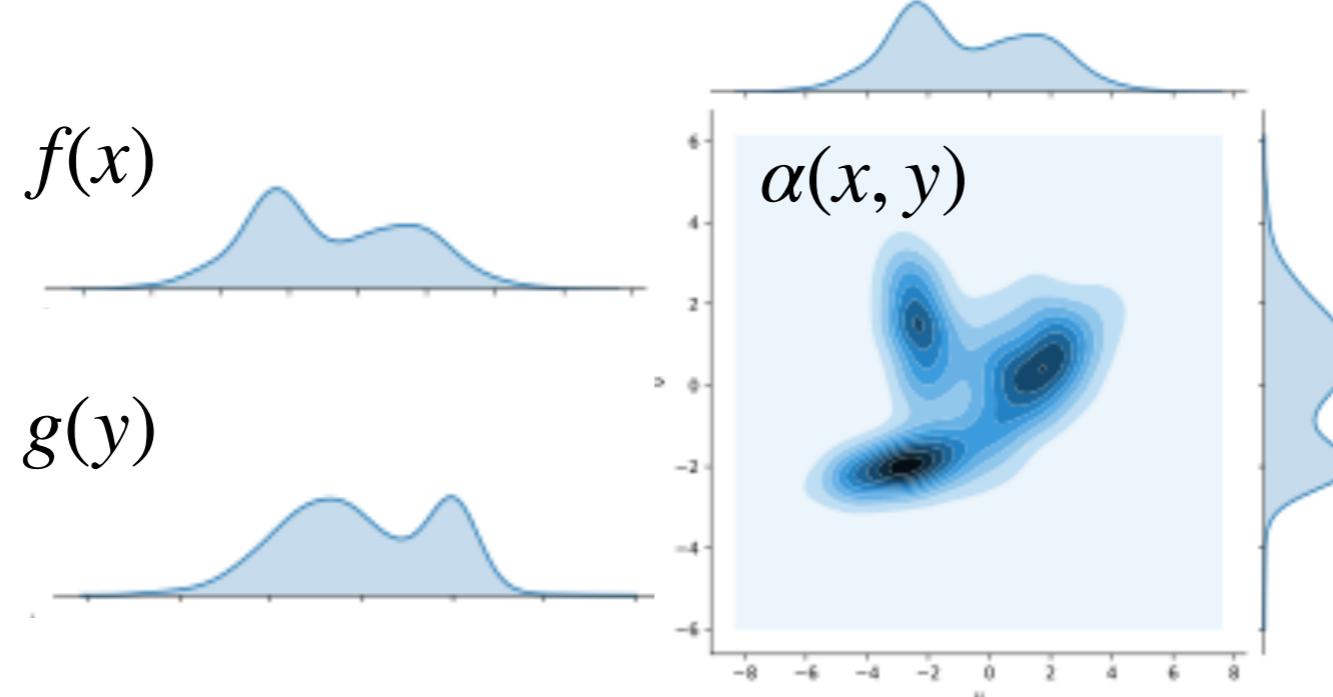
$\alpha(x, y)$: amount of mass to move from x to y
 $\alpha(x, y)$ describes a feasible transport plan if

$$\int \alpha(x, y) dy = f(x), \quad \int \alpha(x, y) dx = g(y)$$

Summary: Optimal Transport & Wasserstein Metric

Wasserstein $d_p(F, G) := \inf_{(U,V): U \sim F, V \sim G} ||U - V||_p$

Optimal
Transport
(OT)



$c(x, y)$ = cost function of moving one unit of mass from x to y

- ▶ OT can be written as an optimization problem:

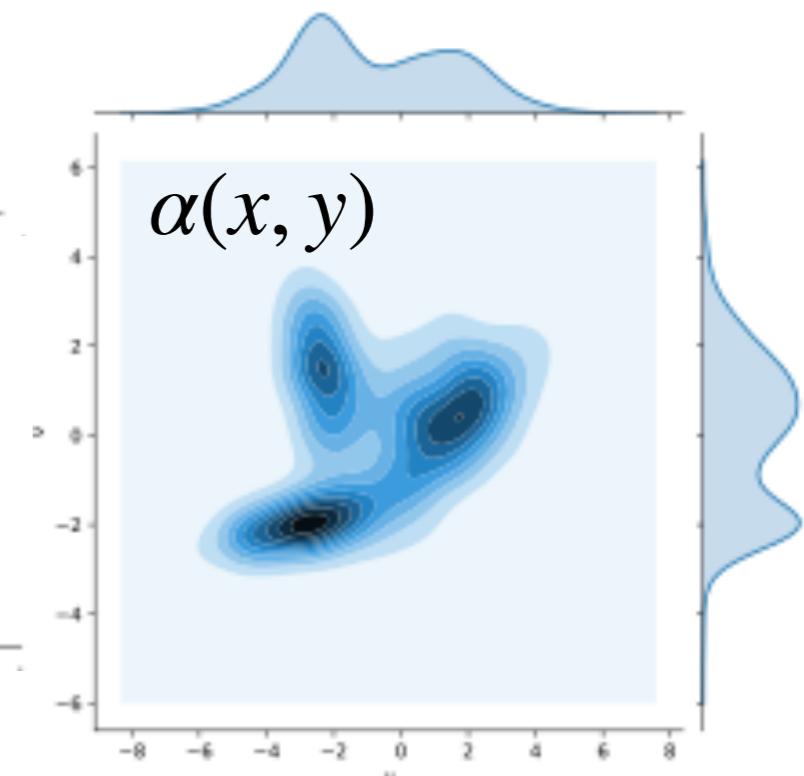
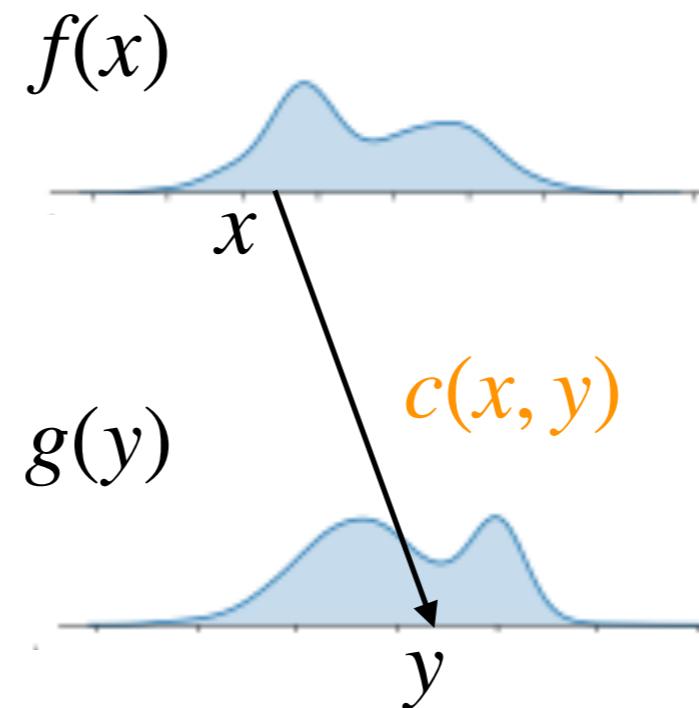
$$\min_{\alpha} \sum_{x,y} c(x, y) \alpha(x, y)$$

subject to (1) $\sum_y \alpha(x, y) = f(x), \forall x$ (2) $\sum_x \alpha(x, y) = g(y), \forall y$
(3) $\alpha(x, y) \geq 0, \forall x, y$

Duality of Optimal Transport: Economic Interpretation

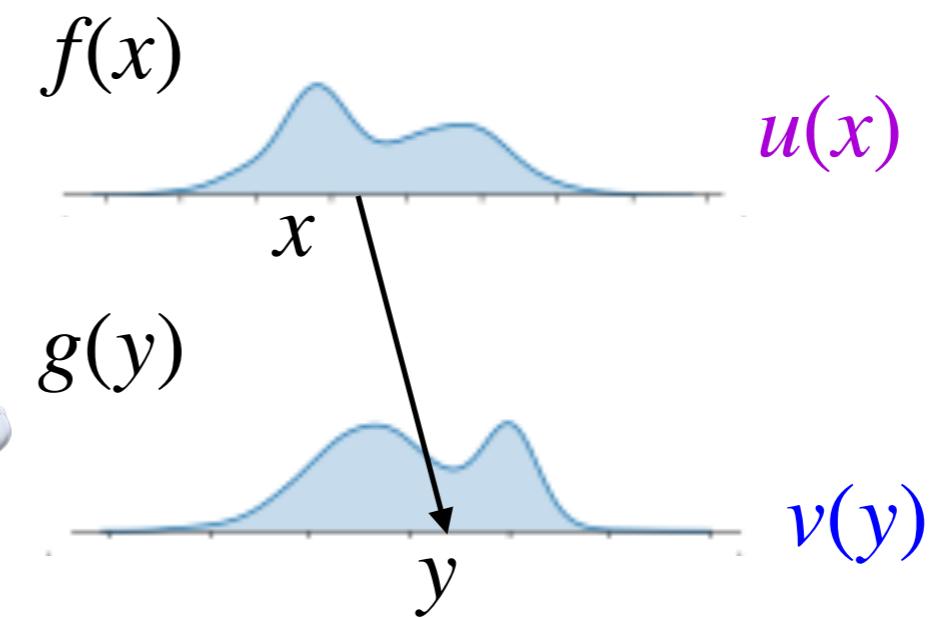
Primal Form of OT (Mario moving the earth by himself)

$c(x, y) = \text{Mario's cost function}$
(for moving one unit of mass
from x to y)



Dual Form of OT (Luigi offers to help Mario)

Mario needs to pay Luigi $u(x)$
and $v(y)$ (for moving one unit
of mass from x to y)



Question: Under what condition would Mario ask for Luigi's help?

Duality of Optimal Transport (Formally)

- ▶ Primal Form of Optimal Transport

$$\min_{\alpha} \sum_{x,y} c(x, y) \alpha(x, y)$$

subject to (1) $\sum_y \alpha(x, y) = f(x), \forall x$ (2) $\sum_x \alpha(x, y) = g(y), \forall y$
(3) $\alpha(x, y) \geq 0, \forall x, y$

- ▶ Dual Form of Optimal Transport

$$\max_{u,v} \mathbb{E}_{x \sim f(x)}[u(x)] + \mathbb{E}_{y \sim g(y)}[v(y)]$$

subject to $u(x) + v(y) \leq c(x, y), \forall x, y$

The dual form looks
exactly like APPLE!

- ▶ Both forms lead to the same optimal values (called “strong duality”)

Example #2: Generative Adversarial Imitation Learning (GAIL)

- ▶ **Recall**: Dual Form of Optimal Transport

$$\max_{u,v} \mathbb{E}_{x \sim f(x)}[u(x)] + \mathbb{E}_{y \sim g(y)}[v(y)]$$

subject to $u(x) + v(y) \leq c(x, y), \forall x, y$

$D_\phi(s, a)$: A **binary classifier** that predicts the probability of the event that “the observed (s, a) is drawn from π ”

Let’s choose the following:

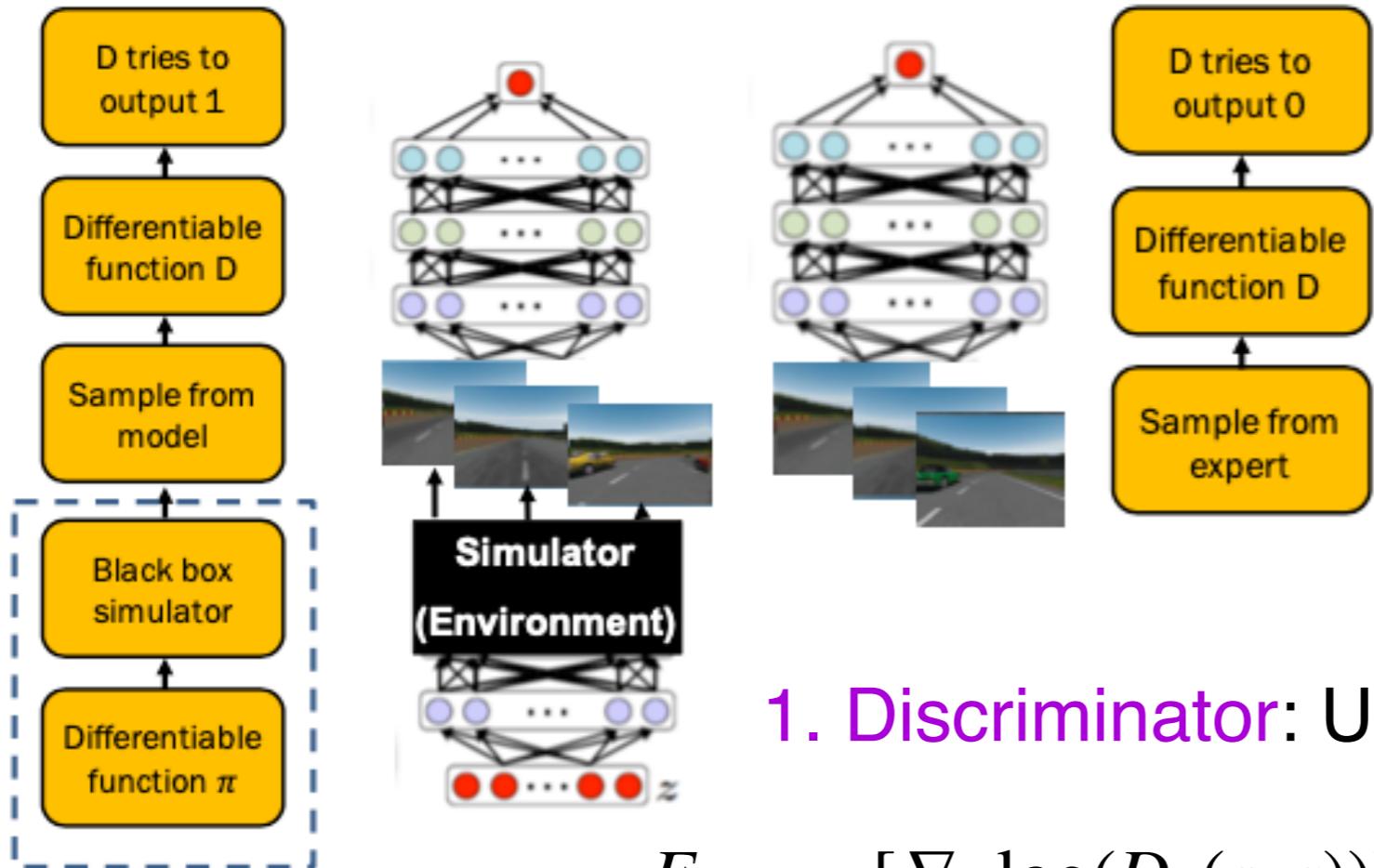
$$(1) f(x) \equiv d_\mu^\pi(s, a)$$

$$(2) g(y) \equiv d_\mu^{\pi_e}(s, a)$$

$$(3) u(x) \equiv \log(D_\phi(s, a))$$

$$(4) v(y) \equiv \log(1 - D_\phi(s, a))$$

GAIL: Discriminator and Generator



1. **Discriminator:** Update ϕ by

$$E_{(s,a) \sim d_\mu^\pi} [\nabla_\phi \log(D_\phi(s, a))] + E_{(s,a) \sim d_\mu^{\pi_e}} [\nabla_\phi \log(1 - D_\phi(s, a))]$$

2. **Generator:** Use any RL algorithm with reward function $\log(D_\phi(s, a))$

A Comparison Between Wasserstein AIL and GAIL

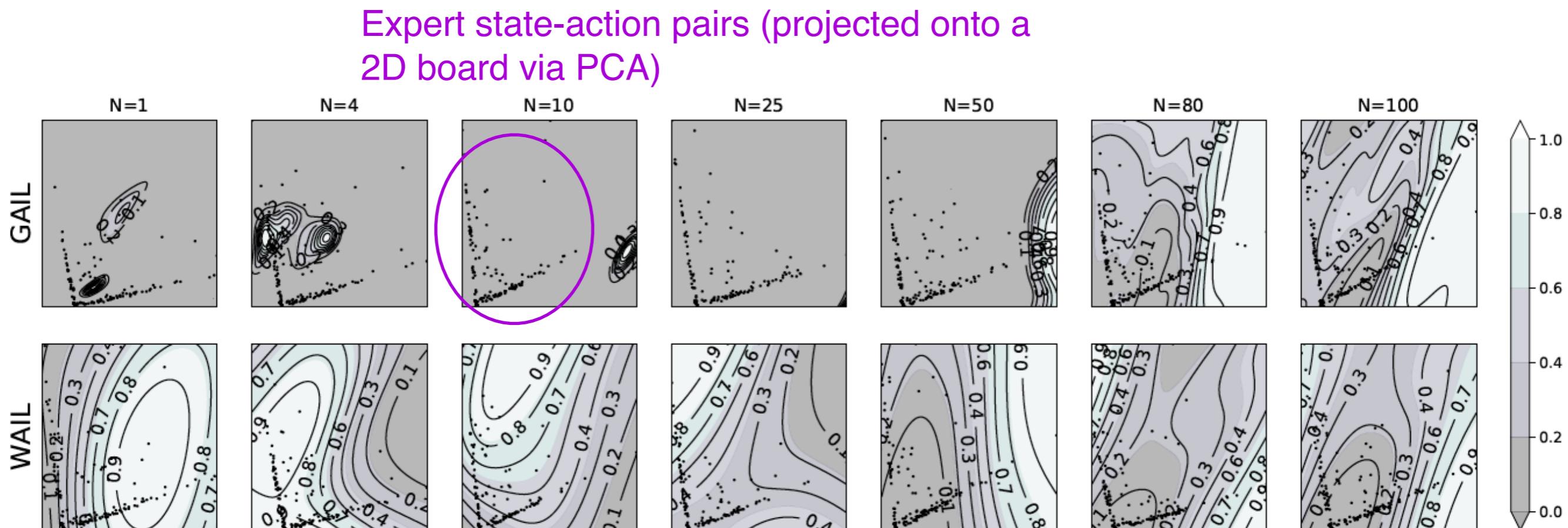


Figure 2: Reward surfaces of WAIL and GAIL on *Humanoid* with respect to different expert data sizes.

Summary: Occupancy Measure Matching via Apprenticeship Learning (With Regularization)

$$\min_{\pi \in \Pi} \max_{R \in \mathcal{R}} \underbrace{\left[\left(E_{(s,a) \sim d_\mu^{\pi_e}}[R(s,a)] - E_{(s,a) \sim d_\mu^\pi}[R(s,a)] \right) - H(\pi) + \psi(R) \right]}_{:= L(\pi, R)}$$

where $H(\pi) := E \left[\sum_t -\gamma^t \log \pi_t(a_t|s_t) \right]$ is the discounted causal entropy

$\psi(R)$ is a regularizer for the reward function

Key Idea: By choosing different “reward function classes \mathcal{R} ”, we obtain various OMM approaches!