

535514: Reinforcement Learning

Lecture 27 — Inverse RL & Model-Based RL

Ping-Chun Hsieh

June 3, 2024

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	MCTS Model-Predictive Control (MPC) PETS	IRL GAIL WAIL
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)		

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

Review: Occupancy Measure Matching

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

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

Property: $V^{\pi}(\mu) = \mathbb{E}_{(s,a) \sim d_{\mu}^{\pi}}[R(s, a)]$

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 $V^{\pi}(\mu) = V^{\pi_e}(\mu)$

(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!

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

OR

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

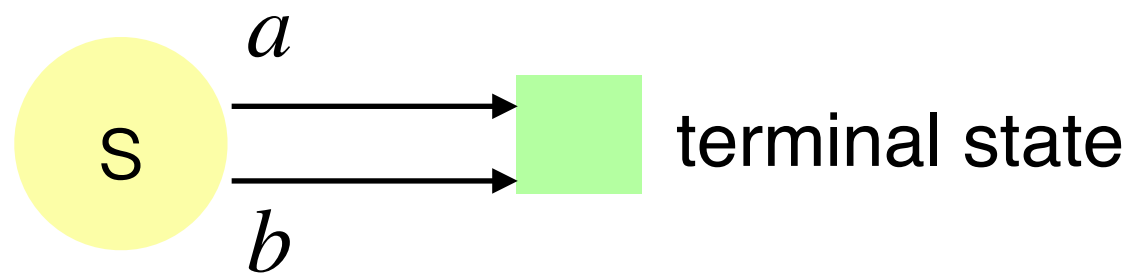
(Easier for training!)

Apprenticeship Learning (APPLE)

A Motivating Example: Connecting OMM & APPLE

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

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

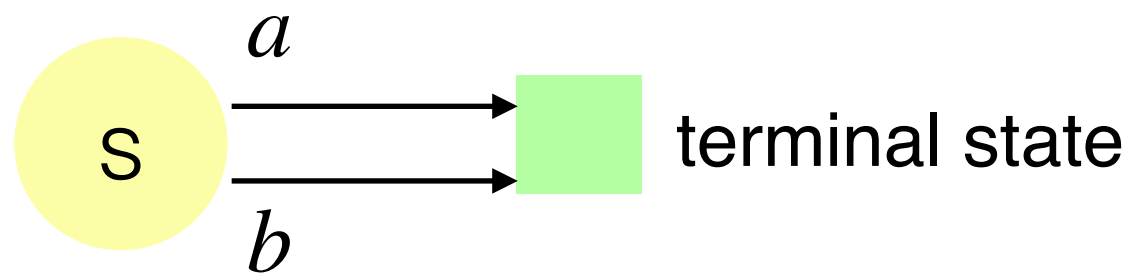
For (s, a) with $d_{\mu}^{\pi}(s, a) < d_{\mu}^{\pi_e}(s, a)$:

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

A Motivating Example: Connecting OMM & APPLE

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

Consider a simple 1-state, 2-action MDP



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

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

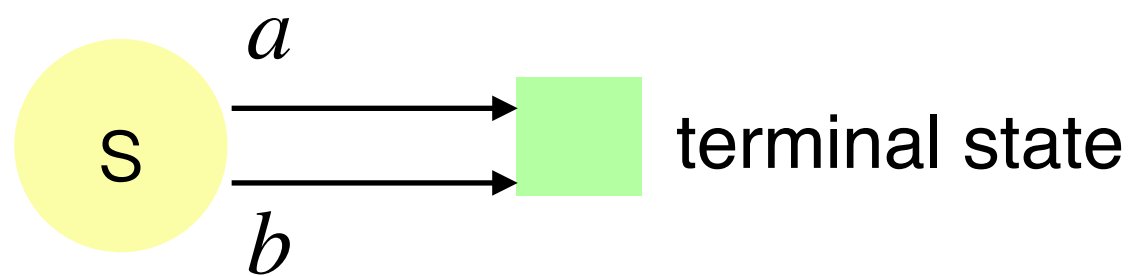
Nice Property: Under $\mathcal{R} = \mathbb{R}^2$, the corresponding metric D is

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

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

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

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

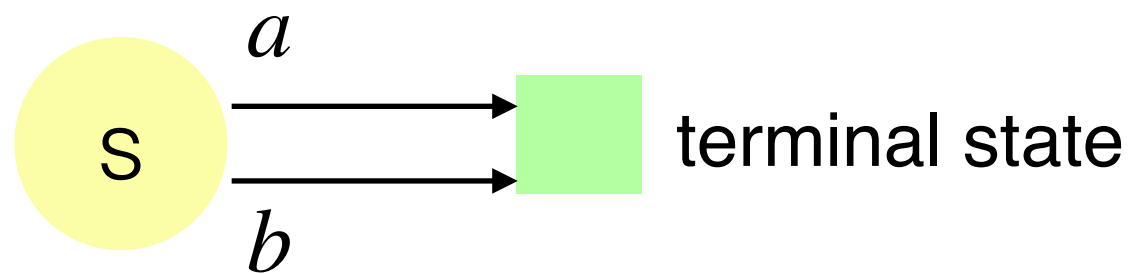
For (s, a) with $d_{\mu}^{\pi}(s, a) < d_{\mu}^{\pi_e}(s, a)$:

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

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

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

Consider a simple 1-state, 2-action MDP



Suppose $\mathcal{R} = \left\{ R \in \mathbb{R}^2 \mid \|R\|_{\infty} \leq 1 \right\}$
 $\pi_e(a|s) = \pi_e(b|s) = 0.5$

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

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

(usually called “*total variation distance*”)

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

Example #1: Wasserstein Metric and APPLE

$$\begin{aligned} \min_{\pi \in \Pi} L(\pi) &:= W_p(d_\mu^\pi, d_\mu^{\pi_e}) && \longleftrightarrow && \min_{\pi \in \Pi} \max_{R \in \mathcal{R}} \left[\underbrace{\left(E_{d_\mu^{\pi_e}}[R(s, a)] - E_{d_\mu^\pi}[R(s, a)] \right)}_{:=L(\pi, R)} \right] \\ & \text{(Wasserstein)} && && \text{where } \mathcal{R} = \left\{ R \in \mathbb{R}^{\mathcal{S} \times \mathcal{A}} \mid \text{Lip}(R) \leq 1 \right\} \end{aligned}$$

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 $1 \leq p < \infty$: $\|U\|_p := \left(\mathbb{E} [\|U(\omega)\|_p^p] \right)^{\frac{1}{p}}$

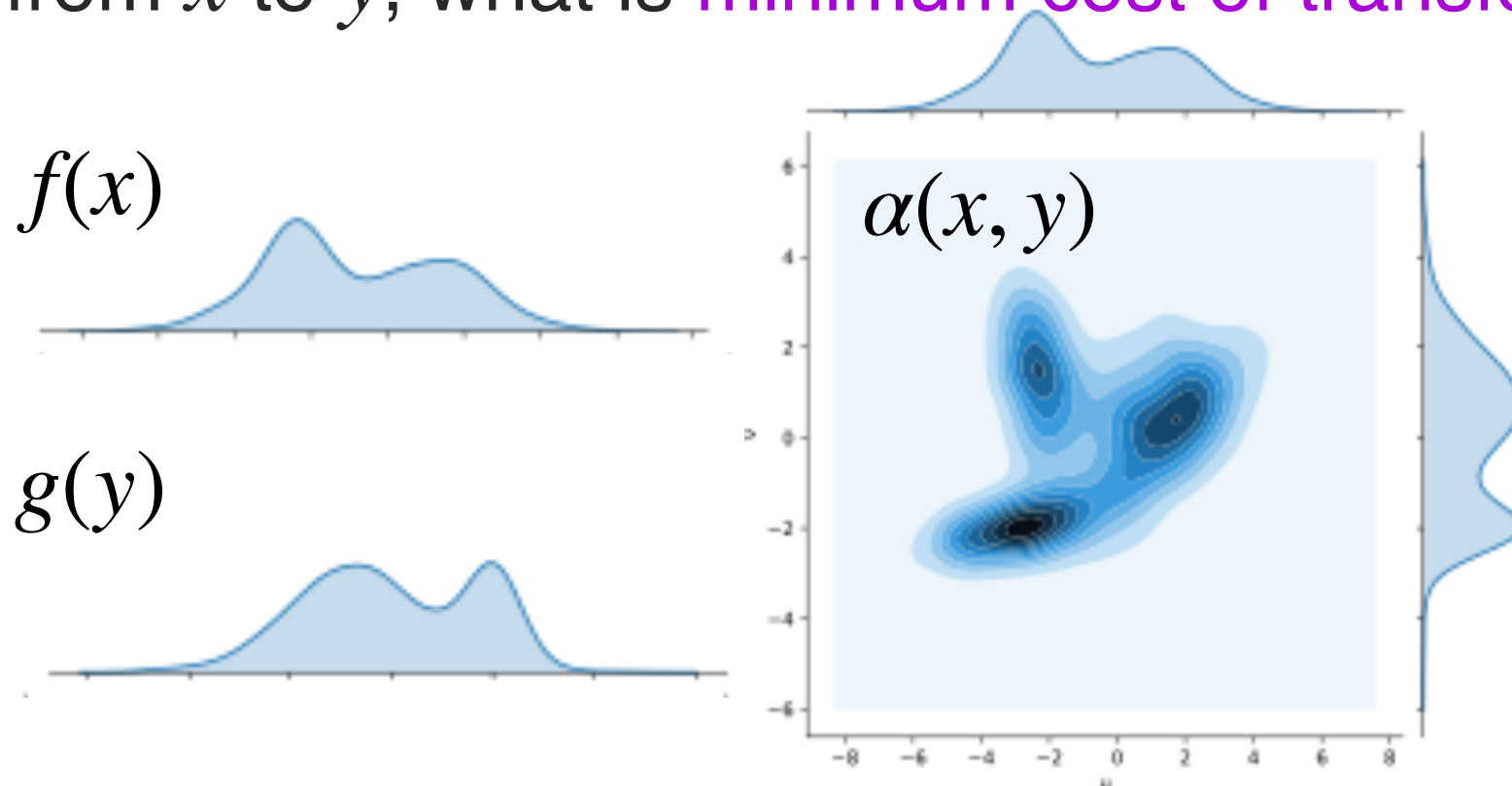
- ▶ **Wasserstein Metric**: For two CDFs F, G over the reals, the Wasserstein metric is defined as

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

- ▶ 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

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

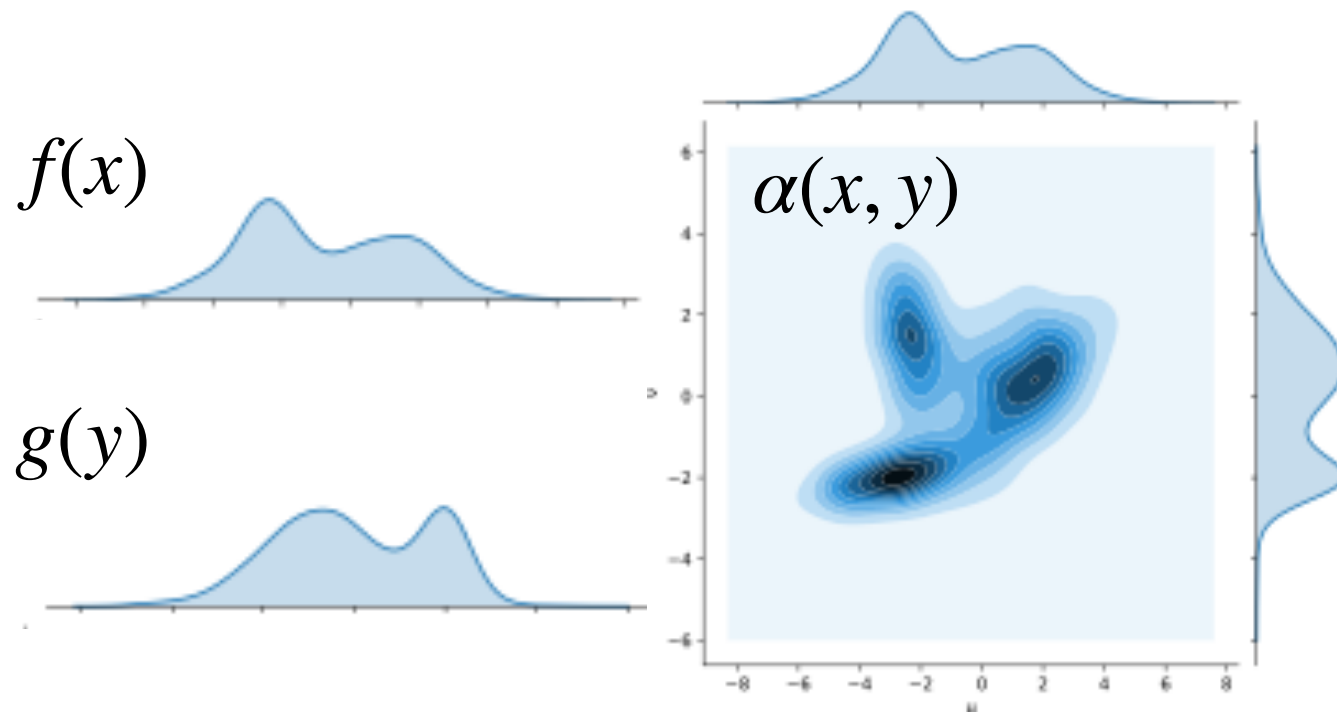
$\alpha(x, y)$: amount of mass to move from x to y
 $\alpha(x, y)$ describes a feasible transport plan if

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

Summary: Optimal Transport & Wasserstein Metric

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

**Optimal
Transport
(OT)**



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

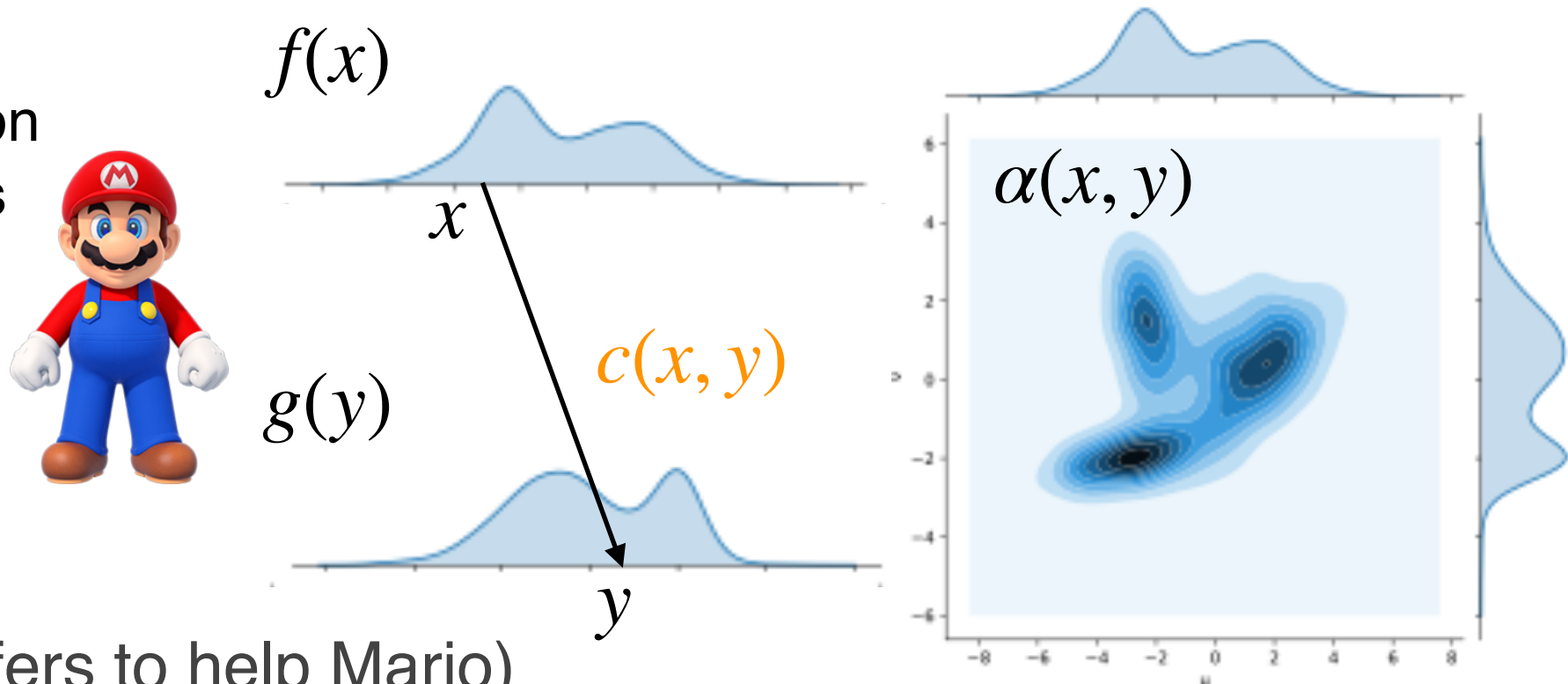
- ▶ OT can be written as an optimization problem:

$$\begin{aligned} & \min_{\alpha} \sum_{x, y} c(x, y) \alpha(x, y) \\ \text{subject to } & (1) \sum_y \alpha(x, y) = f(x), \forall x \quad (2) \sum_x \alpha(x, y) = g(y), \forall y \\ & (3) \alpha(x, y) \geq 0, \forall x, y \end{aligned}$$

Duality of Optimal Transport: Economic Interpretation

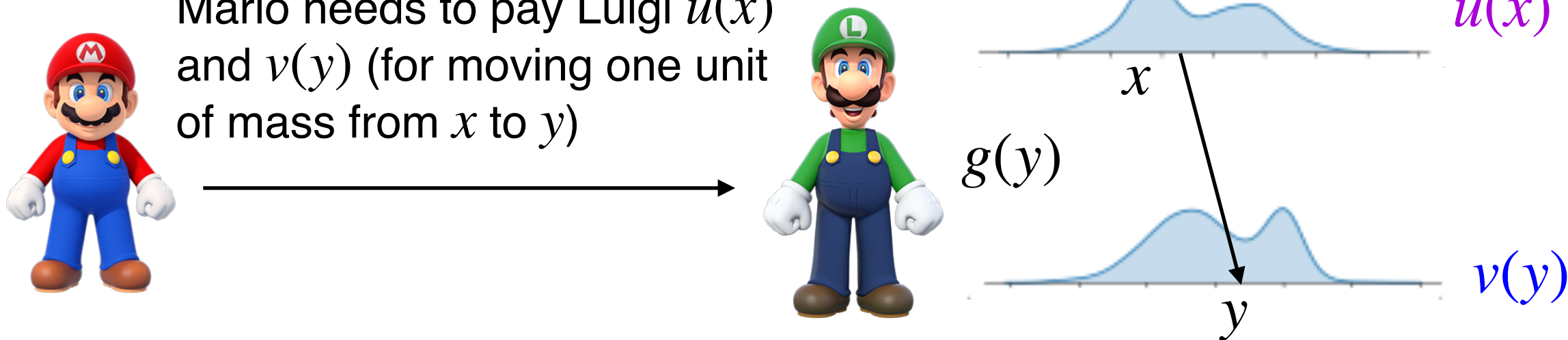
Primal Form of OT (Mario moving the earth by himself)

$c(x, y)$ = 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

$$\begin{aligned} & \min_{\alpha} \sum_{x,y} c(x,y) \alpha(x,y) \\ \text{subject to } & (1) \sum_y \alpha(x,y) = f(x), \forall x \quad (2) \sum_x \alpha(x,y) = g(y), \forall y \\ & (3) \alpha(x,y) \geq 0, \forall x, y \end{aligned}$$

- ▶ Dual Form of Optimal Transport

$$\begin{aligned} & \max_{u,v} \mathbb{E}_{x \sim f(x)}[u(x)] + \mathbb{E}_{y \sim g(y)}[v(y)] \\ \text{subject to } & u(x) + v(y) \leq c(x,y), \forall x, y \end{aligned}$$

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:

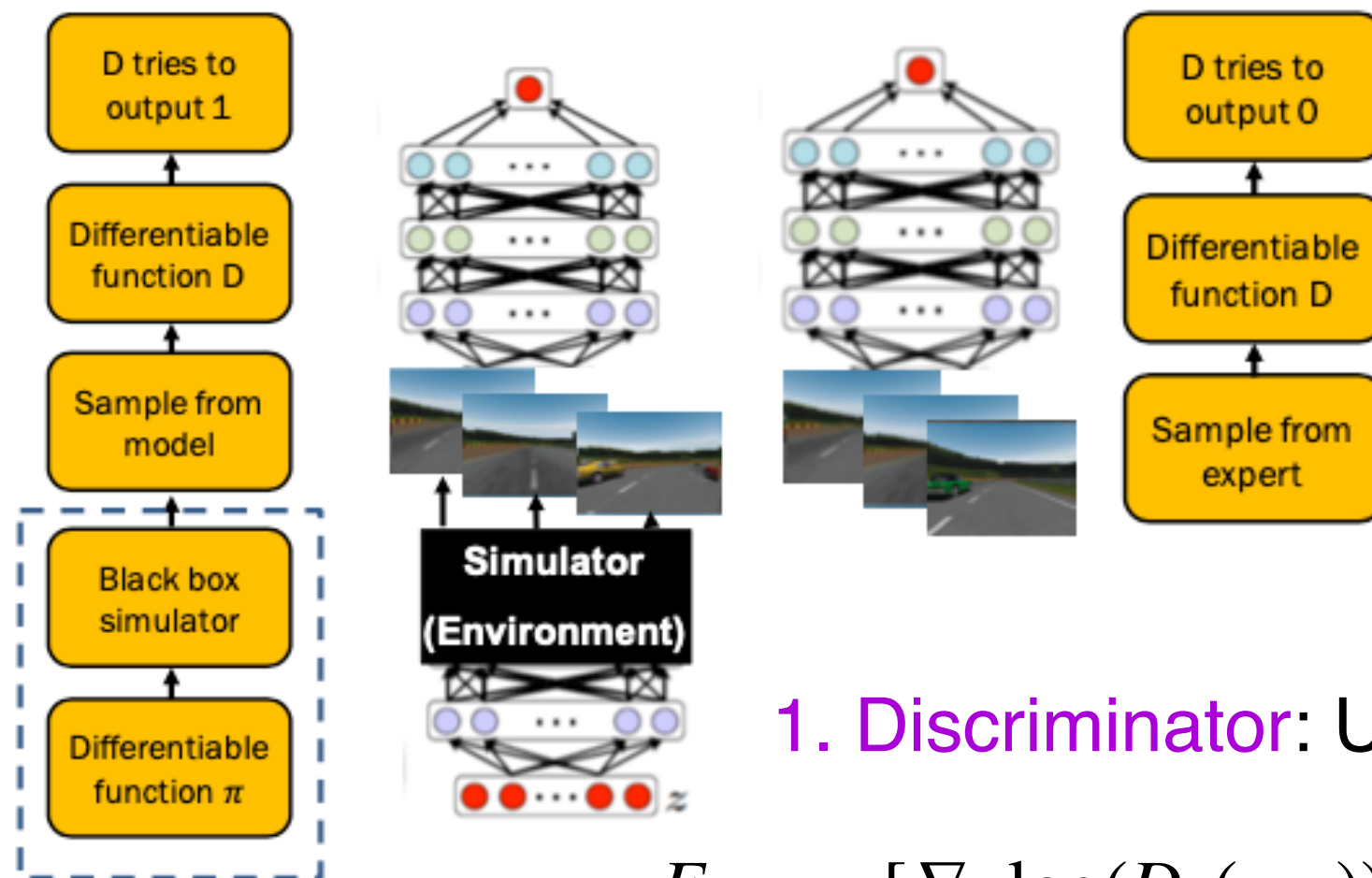
(1) $f(x) \equiv d_\mu^\pi(s, a)$

(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

Expert state-action pairs (projected onto a 2D board via PCA)

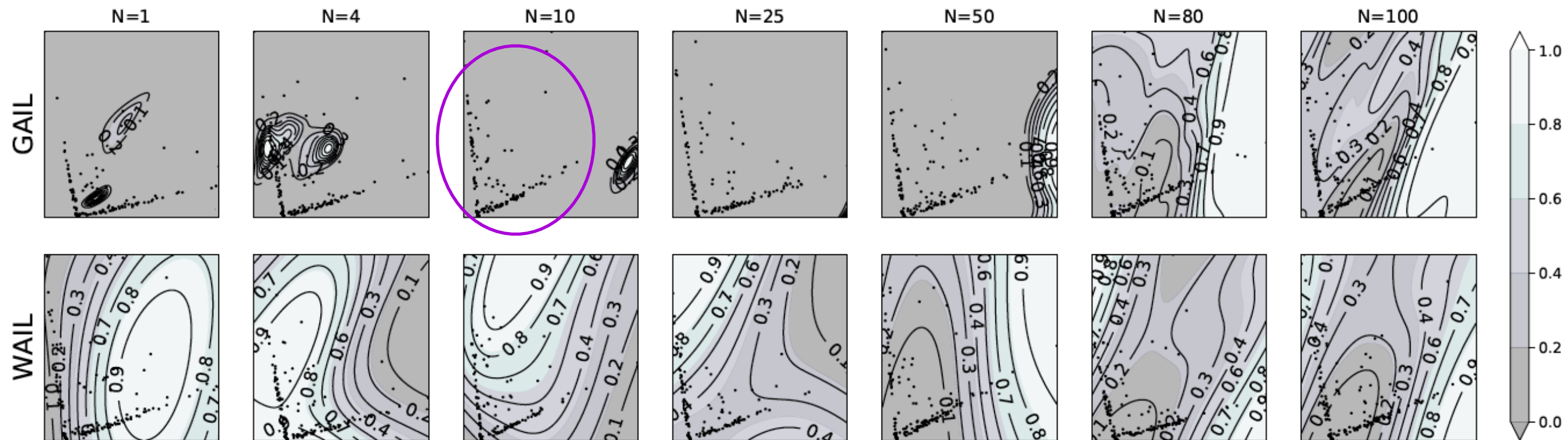


Figure 2: Reward surfaces of WAIL and GAIL on *Humanoid* with respect to different expert data sizes.

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

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

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!

Model-Based RL

Moerland, Thomas M., et al. "Model-based reinforcement learning: A survey." *Foundations and Trends® in Machine Learning*

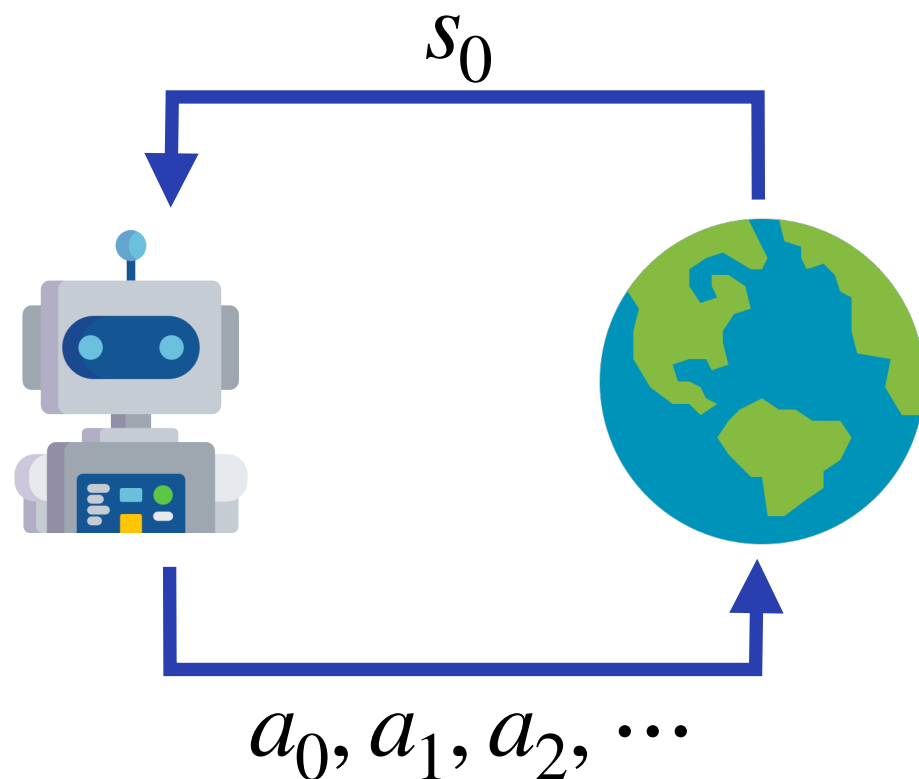
Part of the material is based upon the course material of CS285 by Sergey Levine

Model-Based Reinforcement Learning (MBRL)

- **Idea**: Learn the dynamics model, and then determine the action sequence or a policy
- **Today**: How to determine action sequence if the dynamics model (P, R) is known
- **Next Lecture**: How to learn the dynamics model and apply MBRL in offline settings

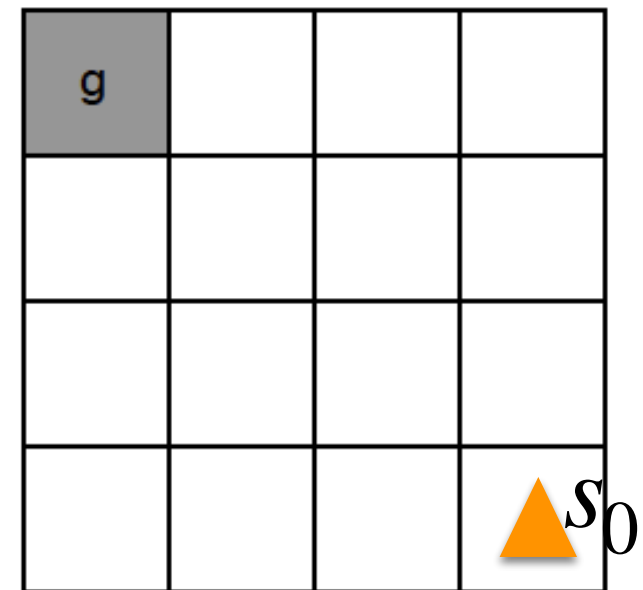
The Simplest Case: Open-Loop Planning

- **Open-loop planning**: Given dynamics model and initial state s_0 , determine the sequence of actions $a_0, a_1, a_2, \dots, a_T$



Example: Deterministic Gridworld

- Reward = -1 , for each step
- Episode ends when reaching “g”



$$\arg \max_{a_0, a_1, \dots, a_T} \mathbb{E} \left[\sum_{t=0}^T R(s_t, a_t) \right]$$

(For simplicity, let's assume $\gamma = 1$)

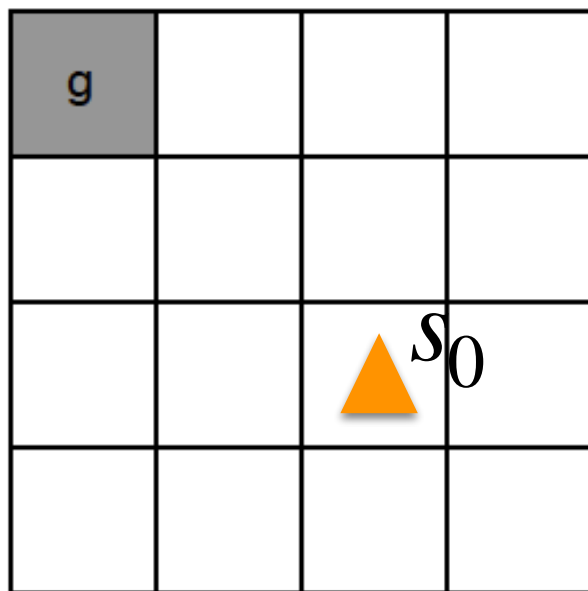
Open-Loop Planning in Stochastic Environments

Open-loop planning:

$$\arg \max_{a_0, a_1, \dots, a_T} \mathbb{E} \left[\sum_{t=0}^T R(s_t, a_t) \right]$$

- ▶ However, open-loop planning can be sub-optimal in stochastic environments

Example: Stochastic Gridworld



Suppose 4 possible actions (Up, Down, Left, Right)

Suppose at each state:

(1) The next state follows the action, w.p. $1 - \epsilon$

(2) The next state can be at any neighboring state, w.p. ϵ

(Suppose ϵ is small)

What's the selected action sequence under open-loop planning?

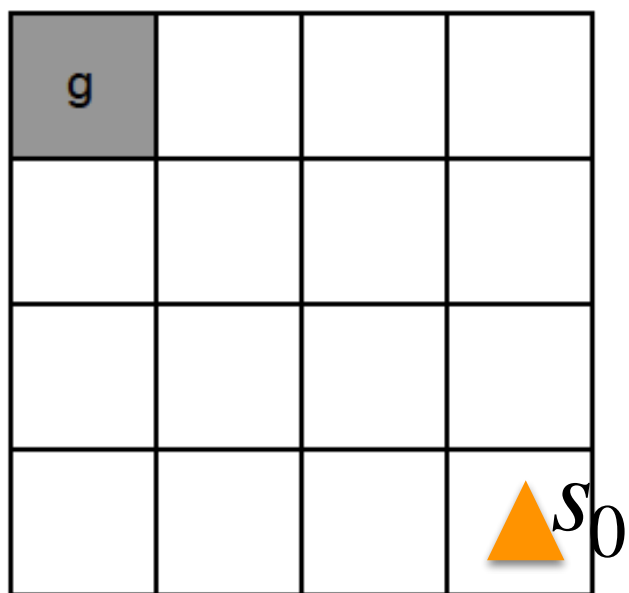
Open-Loop Planning by Stochastic Optimization

- Rethink open-loop planning as an optimization problem (this is often called “**trajectory optimization**”)

$$\arg \max_{a_0, a_1, \dots, a_T} \mathbb{E} \left[\sum_{t=0}^T R(s_t, a_t) \right] \xrightarrow{\theta := (a_0, \dots, a_T)} \arg \max_{\theta} J(\theta)$$

(Here $J(\theta)$ is the total expected reward obtained under θ)

Example: Deterministic Gridworld



Suppose $T = 10$ and set θ as (L, U, D, L, D, L, U, U, U, D, R)
What's the corresponding $J(\theta)$?

Question: In general, given θ , how to obtain $J(\theta)$?

Solution 1: Random Shooting

- ▶ Rethink open-loop planning as an optimization problem (this is often called “**trajectory optimization**”)

$$\arg \max_{a_0, a_1, \dots, a_T} \mathbb{E} \left[\sum_{t=0}^T R(s_t, a_t) \right] \xrightarrow{\theta := (a_0, \dots, a_T)} \arg \max_{\theta} J(\theta)$$

(Here $J(\theta)$ is the total expected reward obtained under θ)

- ▶ **Random shooting** (the simplest open-loop planning approach)

Step 1: Randomly draw action sequences $\theta^{(1)}, \theta^{(2)}, \dots, \theta^{(K)}$ from some distribution (e.g., uniformly random)

Step 2: Choose $\theta^* = \arg \max_{k=1, \dots, K} J(\theta^{(k)})$

Question: Is Random Shooting an efficient method? And why?

Solution 2: Cross-Entropy Method (CEM)

► **Cross-Entropy Method** (a popular open-loop planning approach)

Step 1: Randomly sample action sequences $\theta^{(1)}, \theta^{(2)}, \dots, \theta^{(K)}$ from some distribution $\rho(\theta)$

Step 2: Evaluate $J(\theta^{(1)}), \dots, J(\theta^{(K)})$ and choose $\theta^* = \arg \max_{k=1, \dots, K} J(\theta^{(k)})$

Step 3: Select the top- N candidates $J(\theta^{(i_1)}), \dots, J(\theta^{(i_N)})$

Step 4: Re-fit $\rho(\theta)$ to $\theta^{(i_1)}, \dots, \theta^{(i_N)}$ by maximum likelihood estimation (MLE)

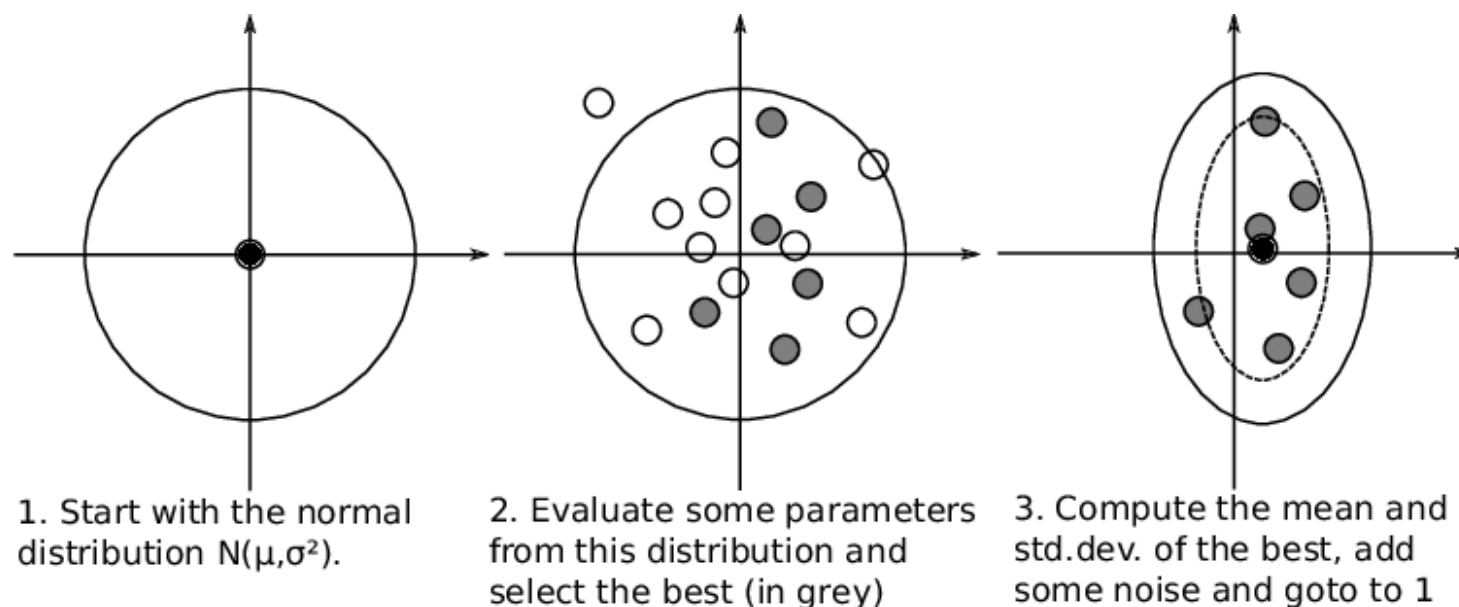


Figure Source: Towards fast and adaptive optimal control policies for robots: A direct policy search approach

Discussions on Open-Loop Planning

- ▶ Main advantages of open-loop planning:

(A1) Simple, zeroth-order optimization (no gradient at all!)

(A2) Very computationally efficient under parallelization

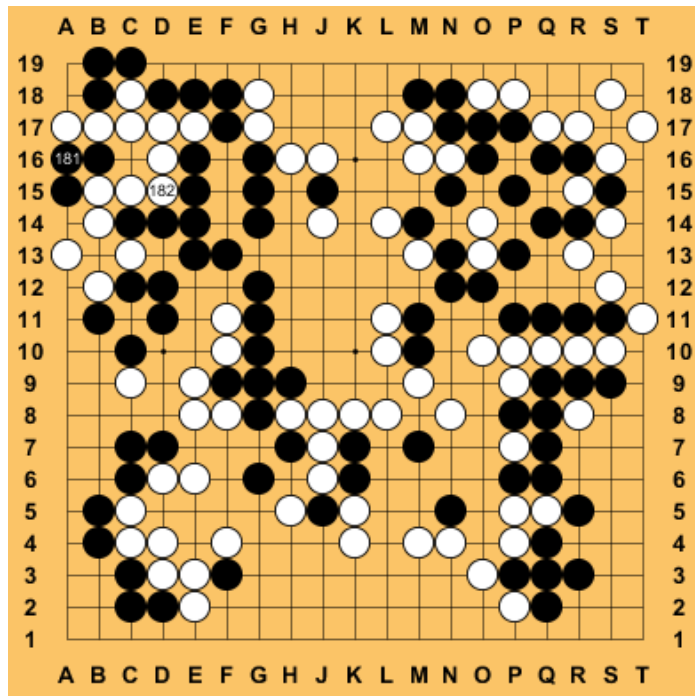
- ▶ Main drawbacks of open-loop planning:

(D1) Not scalable (dimensionally grows with T)

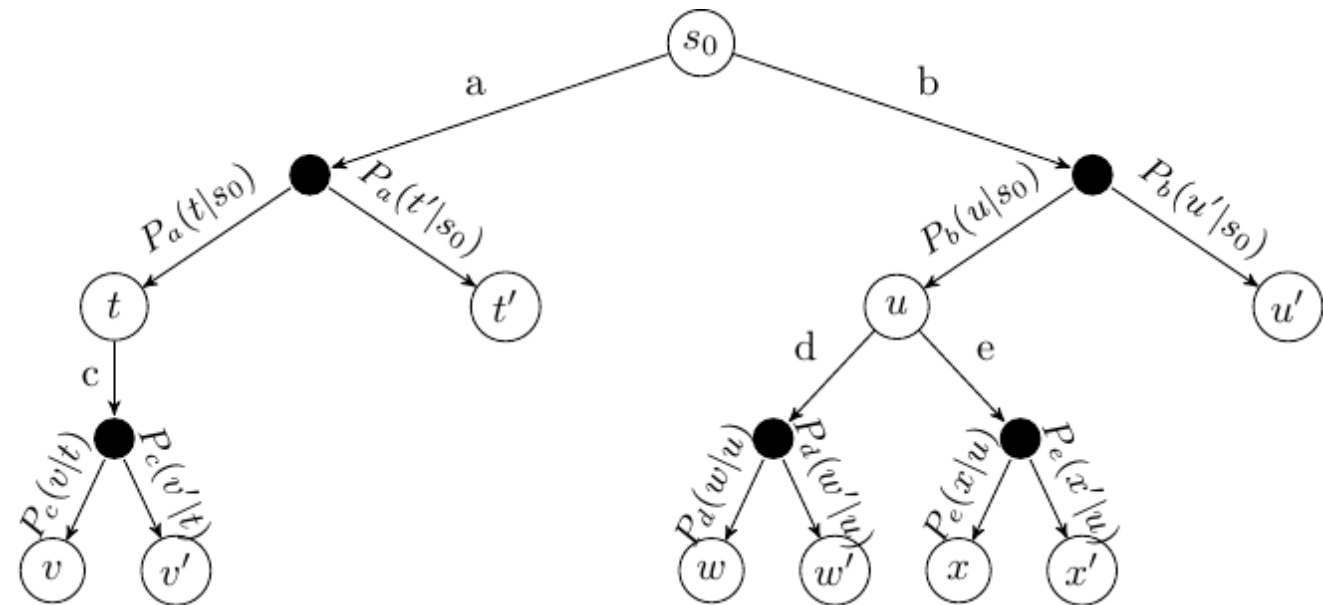
(D2) Can be sub-optimal (e.g., in stochastic environments)

Question: How to do more efficient “planning”?

For (D1): Monte-Carlo Tree Search (MCTS)



In discrete cases (state and action spaces are discrete), “ θ ” can be expressed as a tree



However, the number of nodes can grow exponentially with planning horizon T

Question: How to plan without building a full tree?

Intuition: Select the nodes with highest reward (**exploit**), but also spend some time visiting some unfamiliar nodes (**explore**)

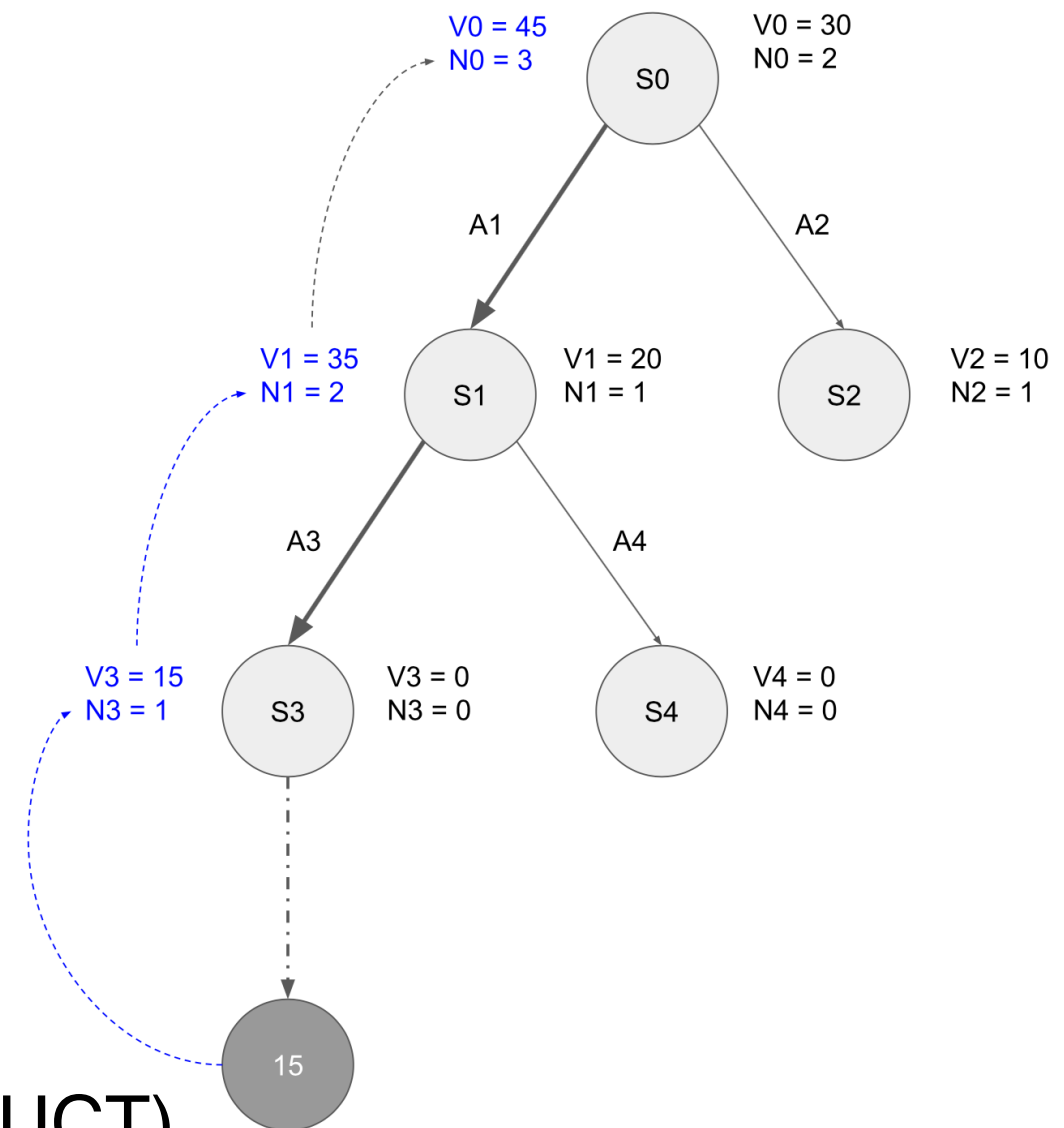
Case Study: Monte-Carlo Tree Search (MCTS)

► General Recipe of MCTS

Step 1: Traverse the tree from s_0 to a leaf node s_n by using TreePolicy

Step 2: Evaluate the leaf node s_n by some EvaluatePolicy

Step 3: Update the values of all the nodes between s_0 and s_n



A popular TreePolicy: Upper-Confidence Tree (UCT)

$$\text{Score of } (s, a) = \frac{V(s, a)}{N(s, a)} + C \sqrt{\frac{\log N(s)}{N(s, a)}}$$

The Classic UCT Paper

Bandit based Monte-Carlo Planning

Levente Kocsis and Csaba Szepesvári

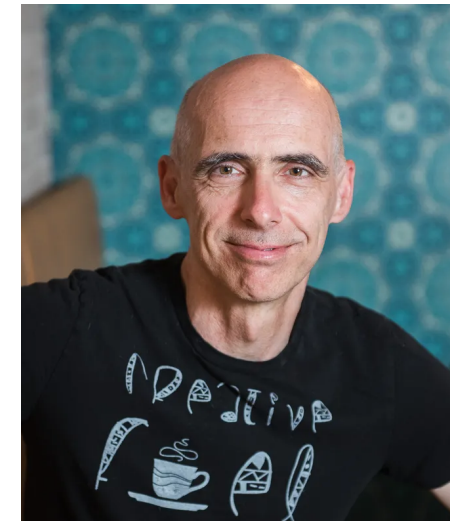
Computer and Automation Research Institute of the
Hungarian Academy of Sciences, Kende u. 13-17, 1111 Budapest, Hungary
kocsis@szttaki.hu

[ECML 2006]

Abstract. For large state-space Markovian Decision Problems Monte-Carlo planning is one of the few viable approaches to find near-optimal solutions. In this paper we introduce a new algorithm, UCT, that applies bandit ideas to guide Monte-Carlo planning. In finite-horizon or discounted MDPs the algorithm is shown to be consistent and finite sample bounds are derived on the estimation error due to sampling. Experimental results show that in several domains, UCT is significantly more efficient than its alternatives.



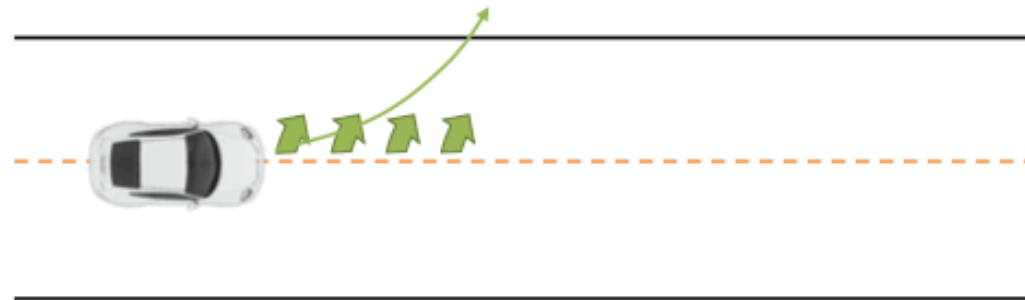
Levente Kocsis



Csaba Szepesvari

For (D2): Closed-Loop Planning

- Why is open-loop planning sub-optimal?



1. Stochastic transitions can put us in unexpected or even dangerous states (even if the dynamics model is fully known)
 2. If the dynamics model is NOT known, then we can make more mistakes
- **Close-Loop Planning:** Replan as you go (replan to fix the mistakes)

Closed-Loop Planning: Model-Predictive Control (MPC)

► Model-Predictive Control

At each time step $t = 0, 1, 2, \dots, T$

Step 1: Plan from current state s_t for a future horizon $H < T$
(e.g., by random shooting or CEM)

Step 2: Execute the first planned action and observe the next state s_{t+1}

Nice features:

1. Short-horizon planning would work
2. Replanning helps with model errors
3. Even random shooting would work sufficiently well!



A Practical Example of MPC: PDDM for Robot Arm

PDDM = Online **P**lanning with **D**eep **D**ynamics **M**odels

Deep Dynamics Models for Learning Dexterous Manipulation

Anusha Nagabandi, Kurt Konoglie, Sergey Levine, Vikash Kumar
Google Brain

Abstract: Dexterous multi-fingered hands can provide robots with the ability to flexibly perform a wide range of manipulation skills. However, many of the more complex behaviors are also notoriously difficult to control: Performing in-hand object manipulation, executing finger gaits to move objects, and exhibiting precise fine motor skills such as writing, all require finely balancing contact forces, breaking and reestablishing contacts repeatedly, and maintaining control of un-actuated objects. Learning-based techniques provide the appealing possibility of acquiring these skills directly from data, but current learning approaches either require large amounts of data and produce task-specific policies, or they have not yet been shown to scale up to more complex and realistic tasks requiring fine motor skills. In this work, we demonstrate that our method of online planning with deep dynamics models (PDDM) addresses both of these limitations; we show that improvements in learned dynamics models, together with improvements in on-line model-predictive control, can indeed enable efficient and effective learning of flexible contact-rich dexterous manipulation skills – and that too, on a 24-DoF anthropomorphic hand in the real world, using just 4 hours of purely real-world data to learn to simultaneously coordinate multiple free-floating objects. Videos can be found at <https://sites.google.com/view/pddm/>

