# Maximum Entropy IRL (continue)

# Recap:

**Constraint Optimization**

$$\min_x f(x) \quad s\,.\,t\,.\,, g_1(x) = 0, g_2(x) = 0, \ldots, g_d(x) = 0$$

# Recap:

## Constraint Optimization

$$\min_{x} f(x) \quad s.t., g_1(x) = 0, g_2(x) = 0, \ldots, g_d(x) = 0$$

Introduce Lagrange multipliers $w \in \mathbb{R}^d$, we have:

$$\min_{x} \max_{w \in \mathbb{R}^d} \underbrace{f(x) + w^\top g(x)}_{:=\ell(x,w)}, \quad (g(x) := [g_1(x), \ldots, g_d(x)]^\top)$$

# Recap:

## Constraint Optimization

$$\min_x f(x) \quad s.t., g_1(x) = 0, g_2(x) = 0, \ldots, g_d(x) = 0$$

Introduce Lagrange multipliers $w \in \mathbb{R}^d$, we have:

$$\min_x \max_{w \in \mathbb{R}^d} \underbrace{f(x) + w^\top g(x)}_{:=\ell(x,w)}, \quad (g(x) := [g_1(x), \ldots, g_d(x)]^\top)$$

We can optimize the dual $\max_{w \in \mathbb{R}^d} \min_x \ell(x, w)$ instead using iterative approach:

# Recap:

## Constraint Optimization

$$\min_x f(x) \quad s.t., g_1(x) = 0, g_2(x) = 0, \ldots, g_d(x) = 0$$

Introduce Lagrange multipliers $w \in \mathbb{R}^d$, we have:

$$\min_x \max_{w \in \mathbb{R}^d} \underbrace{f(x) + w^\top g(x)}_{:= \ell(x,w)}, \quad (g(x) := [g_1(x), \ldots, g_d(x)]^\top)$$

We can optimize the dual $\max_{w \in \mathbb{R}^d} \min_x \ell(x, w)$ instead using iterative approach:

$$x^t = \arg\min_x \ell(x, w^t), \quad \leftarrow \text{Best Response}$$

$$w^{t+1} = w_t + \eta \nabla_w \ell(x^t, w)|_{w=w^t} \quad \leftarrow \text{Gradient Ascent}$$

# Recap on Inverse RL setting:

# Recap on Inverse RL setting:

Q: we want to find a policy $\pi$ such that $\mathbb{E}_{s,a \sim d_\mu^\pi} \phi(s,a) = \mathbb{E}_{s,a \sim d_\mu^{\pi^\star}} \phi(s,a)$

(Note **linear cost assumption** implies $\pi$ is as good as $\pi^\star$)

But there are potentially many such policies...

$$C(sa) = (\theta^*)^T \phi(s,a)$$

$$A$$

$$\mathbb{E}_{sa \sim \pi} \Big[ C(sa) \Big] = \mathbb{E}_{sa \sim \pi^*} \Big[ C(sa) \Big]$$

## Recap on Inverse RL setting:

Q: we want to find a policy $\pi$ such that $\mathbb{E}_{s,a \sim d_\mu^\pi} \phi(s,a) = \mathbb{E}_{s,a \sim d_\mu^{\pi^\star}} \phi(s,a)$

(Note **linear cost assumption** implies $\pi$ is as good as $\pi^\star$)
But there are potentially many such policies…

## The principle of Maximum Entropy:

Find a policy $\pi$ that maximizes some entropy while subject to the constraint:

# Recap on Inverse RL setting:

Q: we want to find a policy $\pi$ such that $\mathbb{E}_{s,a\sim d_\mu^\pi}\phi(s,a) = \mathbb{E}_{s,a\sim d_\mu^{\pi^\star}}\phi(s,a)$

(Note **linear cost assumption** implies $\pi$ is as good as $\pi^\star$)
But there are potentially many such policies…

match
$$d_\mu^\pi \ \& \ d_\mu^{\pi^\star}$$

**The principle of Maximum Entropy:**

Find a policy $\pi$ that maximizes some entropy while subject to the constraint:

$$\min_\pi \ \| d_\mu^\pi(\cdot,\cdot) - d_\mu^{\pi^\star}(\cdot,\cdot) \|_{TV}$$
$$\underbrace{\qquad\qquad\qquad\qquad}_{\Delta}$$

$$\max_\pi \ \mathbb{E}_{s\sim d_\mu^\pi}\left[ \text{entropy}\left( \pi(\,\cdot\,|s) \right) \right]$$

$$s.t, \ \mathbb{E}_{s,a\sim d_\mu^\pi}\phi(s,a) = \mathbb{E}_{s,a\sim d_\mu^{\pi^\star}}\phi(s,a)$$

# Recap on Inverse RL setting:

Q: we want to find a policy $\pi$ such that $\mathbb{E}_{s,a\sim d_\mu^\pi}\phi(s, a) = \mathbb{E}_{s,a\sim d_\mu^{\pi^\star}}\phi(s, a)$

(Note **linear cost assumption** implies $\pi$ is as good as $\pi^\star$)

But there are potentially many such policies…

## The principle of Maximum Entropy:

Find a policy $\pi$ that maximizes some entropy while subject to the constraint:

$$\max_\pi \mathbb{E}_{s\sim d_\mu^\pi}\left[\text{entropy}\left(\pi(\cdot\,|\,s)\right)\right]$$

$$s.t, \mathbb{E}_{s,a\sim d_\mu^\pi}\phi(s, a) = \mathbb{E}_{s,a\sim d_\mu^{\pi^\star}}\phi(s, a)$$

This can be estimated using expert data:

$$\sum_{i=1}^{N}\phi(s_i^\star, a_i^\star)/N$$

**Plan for Today:**

1. The Iterative Algorithm framework

2. How to compute best response: Soft Value Iteration (DP again)

3. The MaxEnt-IRL algorithm

# Setting

Finite horizon MDP $\mathcal{M} = \{S, A, H, c, P, \mu, \pi^\star\}$

## Setting

Finite horizon MDP $\mathcal{M} = \{S, A, H, c, P, \mu, \pi^\star\}$

We have a dataset $\mathcal{D} = (s_i^\star, a_i^\star)_{i=1}^M \sim d_\mu^{\pi^\star}$

# Setting

Finite horizon MDP $\mathcal{M} = \{S, A, H, c, P, \mu, \pi^\star\}$

We have a dataset $\mathcal{D} = (s_i^\star, a_i^\star)_{i=1}^M \sim d_\mu^{\pi^\star}$

**Key Assumption on cost:**

$c(s, a) = \langle \theta^\star, \phi(s, a) \rangle$, **linear w.r.t feature** $\phi(s, a) \in \mathbb{R}^d$

found $\pi$. s.t

$\theta^T \left[ E_\pi \phi(s,a) \right]$

$= (\theta^\star)^T \left( E_{\pi^\star} \phi(s,a) \right)$

(1) $E_{s \sim d_\mu^\pi} \phi(sa) = E_{s \sim d_\mu^{\pi^\star}} \phi(s,a)$

(2) $E_{s \sim d_\mu^\pi} \underbrace{(\theta^\star)^T \phi(s,a)}_{c(s,a)} = E_{s \sim d_\mu^{\pi^\star}} \underbrace{(\theta^\star)^T \phi(s,a)}_{c(s,a)}$

# Notation on Distributions

$\mathbb{P}_h^\pi(s, a; \mu)$: probability of visiting $(s, a)$ at time step $h$ following $\pi$

$$d_\mu^\pi(s, a) = \sum_{h=0}^{H-1} \mathbb{P}_h^\pi(s, a; \mu)/H$$: average state-action distribution

$$d_\mu^\pi(s) = \sum_a d_\mu^\pi(s, a)$$: average state distribution

# Maximum Entropy Inverse RL:

Let's simplify the objective $\max_{\pi} \mathbb{E}_{s \sim d_\mu^\pi} \left[ \text{entropy}(\pi( \cdot \,|\, s)) \right]$:

$$ := - \mathbb{E}_{a \sim \pi(\cdot|s)} \ln \pi(a|s) $$

# Maximum Entropy Inverse RL:

Let's simplify the objective $\max_{\pi} \mathbb{E}_{s \sim d_\mu^\pi} \left[ \text{entropy}(\pi( \cdot \mid s)) \right]$:

$$\mathbb{E}_{s \sim d_\mu^\pi} \left[ \text{entropy}(\pi( \cdot \mid s)) \right] = - \mathbb{E}_{s \sim d_\mu^\pi} \mathbb{E}_{a \sim \pi(\cdot \mid s)} \ln \pi(a \mid s) = - \mathbb{E}_{s,a \sim d_\mu^\pi} \ln \pi(a \mid s)$$

Entropy

# Maximum Entropy Inverse RL:

Let's simplify the objective $\max_{\pi} \mathbb{E}_{s \sim d_{\mu}^{\pi}} \left[ \text{entropy}(\pi( \cdot \,|\, s)) \right]$:

$$\mathbb{E}_{s \sim d_{\mu}^{\pi}} \left[ \text{entropy}(\pi( \cdot \,|\, s)) \right] = - \mathbb{E}_{s \sim d_{\mu}^{\pi}} \mathbb{E}_{a \sim \pi(\cdot|s)} \ln \pi(a \,|\, s) = - \mathbb{E}_{s,a \sim d_{\mu}^{\pi}} \ln \pi(a \,|\, s)$$

$$\arg \max_{\pi} \mathbb{E}_{s \sim d_{\mu}^{\pi}} \left[ \text{entropy}(\pi( \cdot \,|\, s)) \right] = \arg \min_{\pi} \mathbb{E}_{s,a \sim d_{\mu}^{\pi}} \ln \pi(a \,|\, s)$$

# Maximum Entropy Inverse RL formulation

We arrive at the following constraint optimization problem:

$$\arg\min_{\pi} \mathbb{E}_{s,a \sim d_{\mu}^{\pi}} \ln \pi(a \mid s)$$

$$s.t, \mathbb{E}_{s,a \sim d_{\mu}^{\pi}} \phi(s,a) = \mathbb{E}_{s,a \sim d_{\mu}^{\pi^\star}} \phi(s,a)$$

# Maximum Entropy Inverse RL formulation

We arrive at the following constraint optimization problem:

$$\arg\min_{\pi} \mathbb{E}_{s,a \sim d_\mu^\pi} \ln \pi(a \mid s)$$

$$s.t, \mathbb{E}_{s,a \sim d_\mu^\pi} \phi(s,a) = \mathbb{E}_{s,a \sim d_\mu^{\pi\star}} \phi(s,a)^{\in \mathbb{R}^d}$$

Introduce the Lagrange multiplier $w \in \mathbb{R}^d$ (we have d many constraints), consider the max-min dual version:

# Maximum Entropy Inverse RL formulation

We arrive at the following constraint optimization problem:

$$\arg \min_{\pi} \mathbb{E}_{s,a \sim d_{\mu}^{\pi}} \ln \pi(a \mid s)$$

$$s.t, \mathbb{E}_{s,a \sim d_{\mu}^{\pi}} \phi(s,a) = \mathbb{E}_{s,a \sim d_{\mu}^{\pi^{\star}}} \phi(s,a) \quad \leftarrow \quad d \text{ many constraints}$$

$$\phi \in \mathbb{R}^{d}$$

Introduce the Lagrange multiplier $w \in \mathbb{R}^{d}$ (we have d many constraints), consider the max-min dual version:

$$\max_{w \in \mathbb{R}^{d}} \min_{\pi} \underbrace{\mathbb{E}_{s,a \sim d_{\mu}^{\pi}} \ln \pi(a \mid s) + w^{\top} \left( \mathbb{E}_{s,a \sim d_{\mu}^{\pi}} \phi(s,a) - \mathbb{E}_{s,a \sim d_{\mu}^{\pi^{\star}}} \phi(s,a) \right)}_{:= \ell(\pi, w)}$$

# Maximum Entropy Inverse RL formulation

We arrive at the following constraint optimization problem:

$$\arg \min_{\pi} \mathbb{E}_{s,a \sim d_{\mu}^{\pi}} \ln \pi(a \mid s)$$

$$s.t, \mathbb{E}_{s,a \sim d_{\mu}^{\pi}} \phi(s,a) = \mathbb{E}_{s,a \sim d_{\mu}^{\pi^{\star}}} \phi(s,a)$$

Introduce the Lagrange multiplier $w \in \mathbb{R}^d$ (we have d many constraints), consider the max-min dual version:

$$\max_{w \in \mathbb{R}^d} \min_{\pi} \underbrace{\mathbb{E}_{s,a \sim d_{\mu}^{\pi}} \ln \pi(a \mid s) + w^{\top} \left( \mathbb{E}_{s,a \sim d_{\mu}^{\pi}} \phi(s,a) - \mathbb{E}_{s,a \sim d_{\mu}^{\pi^{\star}}} \phi(s,a) \right)}_{:=\ell(\pi,w)}$$

Let's solve it by the iterative procedure!

# Maximum Entropy Inverse RL Algorithm framework

$$\max_{w \in \mathbb{R}^d} \min_{\pi} \underbrace{\mathbb{E}_{s,a \sim d_\mu^\pi} \ln \pi(a \mid s) + w^\top \left( \mathbb{E}_{s,a \sim d_\mu^\pi} \phi(s, a) - \mathbb{E}_{s,a \sim d_\mu^{\pi^\star}} \phi(s, a) \right)}_{:= \ell(\pi, w)}$$

# Maximum Entropy Inverse RL Algorithm framework

$$\max_{w \in \mathbb{R}^d} \min_{\pi} \underbrace{\mathbb{E}_{s,a \sim d_\mu^\pi} \ln \pi(a \mid s) + w^\top \left( \mathbb{E}_{s,a \sim d_\mu^\pi} \phi(s,a) - \mathbb{E}_{s,a \sim d_\mu^{\pi^\star}} \phi(s,a) \right)}_{:=\ell(\pi,w)}$$

Initialize $w^0 \in \mathbb{R}^d$

# Maximum Entropy Inverse RL Algorithm framework

$$\max_{w \in \mathbb{R}^d} \min_{\pi} \underbrace{\mathbb{E}_{s,a \sim d_\mu^\pi} \ln \pi(a \mid s) + w^\top \left( \mathbb{E}_{s,a \sim d_\mu^\pi} \phi(s, a) - \mathbb{E}_{s,a \sim d_\mu^{\pi^\star}} \phi(s, a) \right)}_{:= \ell(\pi, w)}$$

Initialize $w^0 \in \mathbb{R}^d$

For $t = 0 \to T - 1$

# Maximum Entropy Inverse RL Algorithm framework

$$\max_{w \in \mathbb{R}^d} \min_{\pi} \underbrace{\mathbb{E}_{s,a \sim d_\mu^\pi} \ln \pi(a \mid s) + w^\top \left( \mathbb{E}_{s,a \sim d_\mu^\pi} \phi(s,a) - \mathbb{E}_{s,a \sim d_\mu^{\pi^\star}} \phi(s,a) \right)}_{:=\ell(\pi,w)}$$

Initialize $w^0 \in \mathbb{R}^d$

For $t = 0 \rightarrow T - 1$

$$\pi^t = \arg\min_{\pi} \mathbb{E}_{s,a \sim d_\mu^\pi} \left[ (w^t)^\top \phi(x,a) + \ln \pi(a \mid s) \right]$$

# Maximum Entropy Inverse RL Algorithm framework

$$\max_{w \in \mathbb{R}^d} \min_{\pi} \underbrace{\mathbb{E}_{s,a \sim d_\mu^\pi} \ln \pi(a \mid s) + w^\top \left( \mathbb{E}_{s,a \sim d_\mu^\pi} \phi(s,a) - \mathbb{E}_{s,a \sim d_\mu^{\pi^\star}} \phi(s,a) \right)}_{:= \ell(\pi, w)}$$

Initialize $w^0 \in \mathbb{R}^d$

$$\min_{\pi} \ell(\pi, w^t) \iff \min_{\pi} \left[ \underset{s \sim d_\mu^\pi}{\mathbb{E}} \left[ (w^t)^\top \phi(s,a) \right] + \ln \pi(a \mid s) \right]$$

For $t = 0 \to T - 1$

$$\pi^t = \arg \min_{\pi} \mathbb{E}_{s,a \sim d_\mu^\pi} \left[ (w^t)^\top \phi(s,a) + \ln \pi(a \mid s) \right]$$

(# best response: $\pi^t = \arg \min_{\pi} \ell(\pi, w^t)$)

# Maximum Entropy Inverse RL Algorithm framework

$$\max_{w \in \mathbb{R}^d} \min_{\pi} \underbrace{\mathbb{E}_{s,a \sim d_\mu^\pi} \ln \pi(a \mid s) + w^\top \left( \mathbb{E}_{s,a \sim d_\mu^\pi} \phi(s,a) - \mathbb{E}_{s,a \sim d_\mu^{\pi^\star}} \phi(s,a) \right)}_{:= \ell(\pi, w)}$$

Initialize $w^0 \in \mathbb{R}^d$

For $t = 0 \to T - 1$

$$\pi^t = \arg\min_{\pi} \mathbb{E}_{s,a \sim d_\mu^\pi} \left[ (w^t)^\top \phi(x,a) + \ln \pi(a \mid s) \right]$$

(# best response: $\pi^t = \arg\min_{\pi} \ell(\pi, w^t)$)

$$w^{t+1} = w^t + \eta \left( \underbrace{\mathbb{E}_{s,a \sim d_\mu^{\pi^t}} \phi(s,a) - \mathbb{E}_{s,a \sim d_\mu^{\pi^\star}} \phi(s,a)}_{:= \nabla_w \ell(\pi^t, w)\big|_{w=w^t}} \right)$$

(# gradient update: $w^{t+1} = w^t + \eta \underline{\nabla_w \ell(\pi^t, w^t)}$)

# Maximum Entropy Inverse RL Algorithm framework

$$\max_{w \in \mathbb{R}^d} \min_{\pi} \underbrace{\mathbb{E}_{s,a \sim d_\mu^\pi} \ln \pi(a \,|\, s) + w^\top \left( \mathbb{E}_{s,a \sim d_\mu^\pi} \phi(s, a) - \mathbb{E}_{s,a \sim d_\mu^{\pi^\star}} \phi(s, a) \right)}_{:= \ell(\pi, w)}$$

Initialize $w^0 \in \mathbb{R}^d$

For $t = 0 \to T - 1$

$$\pi^t = \arg\min_{\pi} \mathbb{E}_{s,a \sim d_\mu^\pi} \left[ (w^t)^\top \phi(x, a) + \ln \pi(a \,|\, s) \right] \quad \text{(\# best response: } \pi^t = \arg\min_{\pi} \ell(\pi, w^t))$$

$$w^{t+1} = w^t + \eta \left( \mathbb{E}_{s,a \sim d_\mu^\pi} \phi(s, a) - \mathbb{E}_{s,a \sim d_\mu^{\pi^\star}} \phi(s, a) \right)$$

$$\text{(\# gradient update: } w^{t+1} = w^t + \eta \nabla_w \ell(\pi^t, w^t))$$

Return $\bar{\pi} = \text{Uniform}(\pi^0, \ldots, \pi^{T-1})$

# Maximum Entropy Inverse RL Algorithm framework

$$\max_{w \in \mathbb{R}^d} \min_{\pi} \underbrace{\mathbb{E}_{s,a \sim d_\mu^\pi} \ln \pi(a \mid s) + w^\top \left( \mathbb{E}_{s,a \sim d_\mu^\pi} \phi(s,a) - \mathbb{E}_{s,a \sim d_\mu^{\pi^\star}} \phi(s,a) \right)}_{:= \ell(\pi, w)}$$

Initialize $w^0 \in \mathbb{R}^d$

For $t = 0 \to T-1$

This is like an RL problem w/ cost
$c(s,a) := (w^t)^\top \phi(s,a)$, but w/ an additional $\ln \pi(a \mid s)$

$\pi^t = \arg\min_{\pi} \boxed{\mathbb{E}_{s,a \sim d_\mu^\pi} \left[ (w^t)^\top \phi(x,a) + \ln \pi(a \mid s) \right]}$ (# best response: $\pi^t = \arg\min_{\pi} \ell(\pi, w^t)$)

$w^{t+1} = w^t + \eta \left( \mathbb{E}_{s,a \sim d_\mu^\pi} \phi(s,a) - \mathbb{E}_{s,a \sim d_\mu^{\pi^\star}} \phi(s,a) \right)$

(# gradient update: $w^{t+1} = w^t + \eta \nabla_w \ell(\pi^t, w^t)$)

Return $\bar{\pi} = \underline{\text{Uniform}(\pi^0, \ldots, \pi^{T-1})}$

$\max_{w} \min_{\pi} \ell(\pi, w)$

# Maximum Entropy Inverse RL Algorithm framework

$$\max_{w \in \mathbb{R}^d} \min_{\pi} \underbrace{\mathbb{E}_{s,a \sim d_\mu^\pi} \ln \pi(a \,|\, s) + w^\top \left( \mathbb{E}_{s,a \sim d_\mu^\pi} \phi(s, a) - \mathbb{E}_{s,a \sim d_\mu^{\pi^\star}} \phi(s, a) \right)}_{:=\ell(\pi, w)}$$

Initialize $w^0 \in \mathbb{R}^d$

For $t = 0 \rightarrow T - 1$

This is like an RL problem w/ cost
$c(s, a) := (w^t)^\top \phi(s, a)$, but w/ an additional $\ln \pi(a \,|\, s)$

$$\pi^t = \arg\min_{\pi} \mathbb{E}_{s,a \sim d_\mu^\pi} \left[ (w^t)^\top \phi(s, a) + \ln \pi(a \,|\, s) \right] \quad \text{(\# best response: } \pi^t = \arg\min_{\pi} \ell(\pi, w^t))$$

$$w^{t+1} = w^t + \eta \left( \mathbb{E}_{s,a \sim d_\mu^{\pi^t}} \phi(s, a) - \mathbb{E}_{s,a \sim d_\mu^{\pi^\star}} \phi(s, a) \right)$$

(\# gradient update: $w^{t+1} = w^t + \eta \, \nabla_w \ell(\pi^t, w^t)$)

Return $\bar{\pi} = \text{Uniform}(\pi^0, \ldots, \pi^{T-1})$

# Plan for Today:

✅ 1. The Iterative Algorithm framework

2. How to compute best response: Soft Value Iteration (DP again)

$$\arg\min_{\pi} \mathbb{E}_{s,a \sim d_\mu^\pi} \left[ c(x,a) + \ln \pi(a \mid s) \right]$$

cost function; In Iter $\ell$: $c(s,a) = (w^\ell)^T \phi(s,a)$

3. The MaxEnt-IRL algorithm

# Maximum Entropy RL: Soft Value Iteration

$$\arg\min_{\pi} \mathbb{E}_{s,a \sim d_\mu^\pi} \left[ c(s,a) + \ln \pi(a \mid s) \right]$$

$$\Longleftrightarrow \quad \arg\min_{\pi} \quad \mathbb{E}_{s \sim d_\mu^\pi} \left[ C(s,a) - \text{Entropy}\left( \pi(\cdot \mid s) \right) \right]$$

# Maximum Entropy RL: Soft Value Iteration

$$\arg \min_{\pi} \mathbb{E}_{s,a \sim d_{\mu}^{\pi}} \left[ c(s,a) + \ln \pi(a \mid s) \right]$$

**Soft Value Iteration for finite horizon MDP (Dynamic Programming again):**

$$V_h^*(s) = \min_{\substack{\pi_h, \pi_{h+1} \\ \sim \pi_{H-1}}} \mathbb{E} \left[ \sum_{t=h}^{H-1} c(s_t, a_t) + \ln \pi_t(a_t \mid s_t) \mid \substack{s_h = s, \\ a_t \sim \pi_t} \right]$$

# Maximum Entropy RL: Soft Value Iteration

$$\arg \min_{\pi} \mathbb{E}_{s,a \sim d_\mu^\pi} \left[ c(s,a) + \ln \pi(a \mid s) \right]$$

**Soft Value Iteration for finite horizon MDP (Dynamic Programming again):**

$$V_H^\star(s) = 0 \quad , \quad \forall s$$

# Maximum Entropy RL: Soft Value Iteration

$$\arg \min_{\pi} \mathbb{E}_{s,a \sim d_{\mu}^{\pi}} \left[ c(s,a) + \ln \pi(a \mid s) \right]$$

## Soft Value Iteration for finite horizon MDP (Dynamic Programming again):

$$V_H^{\star}(s) = 0$$

$$Q_h^{\star}(s,a) = c(s,a) + \mathbb{E}_{s' \sim P(\cdot \mid s,a)} V_{h+1}^{\star}(s')$$

What's the $\pi_h^{\star}(\cdot \mid s)$ ??

$\in \Delta(A)$

$$\min_{P \in \Delta(A)} \left[ \sum_a P(a) \left[ c(s,a) + \ln P(a) + \mathbb{E}_{s' \sim P(s,a)} V_{h+1}^{\star}(s') \right] \right]$$

$$Q_h^{\star}(s,a)$$

# Maximum Entropy RL: Soft Value Iteration

$$\arg \min_{\pi} \mathbb{E}_{s,a \sim d_\mu^\pi} \left[ c(s,a) + \ln \pi(a \,|\, s) \right]$$

**Soft Value Iteration for finite horizon MDP (Dynamic Programming again):**

$$V_H^\star(s) = 0$$

$$Q_h^\star(s,a) = c(s,a) + \mathbb{E}_{s' \sim P(\cdot \,|\, s,a)} V_{h+1}^\star(s')$$

$$\pi_h^\star(\cdot \,|\, s) = \arg \min_{\rho \in \Delta(A)} \left[ \sum_a \rho(a) Q_h^\star(s,a) + \sum_a \rho(a) \ln \rho(a) \right]$$

$\Rightarrow -\text{Entropy} (\rho)$

$\text{S-t} \quad \sum_a \rho(a) = 1$

# Maximum Entropy RL: Soft Value Iteration

$$\arg\min_{\pi} \mathbb{E}_{s,a \sim d_{\mu}^{\pi}} \left[ c(s,a) + \ln \pi(a\,|\,s) \right]$$

**Soft Value Iteration for finite horizon MDP (Dynamic Programming again):**

$$V_H^{\star}(s) = 0$$

$$Q_h^{\star}(s,a) = c(s,a) + \mathbb{E}_{s' \sim P(\cdot|s,a)} V_{h+1}^{\star}(s')$$

$$\pi_h^{\star}(\,\cdot\,|\,s) = \arg\min_{\rho \in \Delta(A)} \left[ \sum_a \rho(a) Q_h^{\star}(s,a) + \sum_a \rho(a)\ln\rho(a) \right]$$

Use Lagrange (we have a constraint here $\sum_a \rho(a) = 1$), we can show:

# Maximum Entropy RL: Soft Value Iteration

$$\arg \min_{\pi} \mathbb{E}_{s,a \sim d_\mu^\pi} \left[ c(s,a) + \ln \pi(a \,|\, s) \right]$$

**Soft Value Iteration for finite horizon MDP (Dynamic Programming again):**

$$V_H^\star(s) = 0$$

$$Q_h^\star(s,a) = c(s,a) + \mathbb{E}_{s' \sim P(\cdot \,|\, s,a)} V_{h+1}^\star(s')$$

For s. at h:

$$\pi_h^\star(\cdot \,|\, s) = \arg \min_{\rho \in \Delta(A)} \left[ \sum_a \rho(a) Q_h^\star(s,a) + \sum_a \rho(a) \ln \rho(a) \right]$$

$\in \Delta(A)$

Use Lagrange (we have a constraint here $\sum_a \rho(a) = 1$), we can show:

$$\pi_h^\star(a \,|\, s) = \frac{\exp(-Q_h^\star(s,a))}{\sum_{a'} \exp(-Q_h^\star(s,a'))}$$

# Maximum Entropy RL: Soft Value Iteration

$$\arg\min_{\pi} \mathbb{E}_{s,a \sim d_{\mu}^{\pi}} \left[ c(s,a) + \ln \pi(a \,|\, s) \right]$$

**Soft Value Iteration for finite horizon MDP (Dynamic Programming again):**

$$V_H^{\star}(s) = 0$$

$$Q_h^{\star}(s,a) = c(s,a) + \mathbb{E}_{s' \sim P(\cdot | s,a)} V_{h+1}^{\star}(s')$$

② $Q_h^{\star}(s,a) = Q_h^{\star}(s,a')$

$\pi_h^{\star}(\cdot | s) = $ uniform $(a,a')$ $\pi_h^{\star}(\cdot \,|\, s) = \arg\min_{\rho \in \Delta(A)} \left[ \sum_a \rho(a) Q_h^{\star}(s,a) + \sum_a \rho(a) \ln \rho(a) \right]$

Use Lagrange (we have a constraint here $\sum_a \rho(a) = 1$), we can show:

classiz

RL

① $Q_h^{\star}(s,a) = 0$

$Q_h^{\star}(s,a') = +\infty$

$\Rightarrow \pi_h^{\star}(a|s) = 1$

$$\pi_h^{\star}(a \,|\, s) = \frac{\exp(-Q_h^{\star}(s,a))}{\sum_{a'} \exp(-Q_h^{\star}(s,a'))}$$

(contrast this to $\arg\min_a Q^{\star}(s,a)$)

# Maximum Entropy RL: Soft Value Iteration

$$\arg \min_{\pi} \mathbb{E}_{s,a \sim d_\mu^\pi} \left[ c(s,a) + \ln \pi(a \mid s) \right]$$

## Soft Value Iteration for finite horizon MDP (continue)

$V_{-1}^{\star}$

$$Q_h^\star(s,a) = c(s,a) + \mathbb{E}_{s' \sim P(\cdot \mid s,a)} V_{h+1}^\star(s')$$

$$\pi_h^\star(a \mid s) = \frac{\exp(-Q_h^\star(s,a))}{\sum_{a'} \exp(-Q_h^\star(s,a'))}$$

$\Rightarrow V_h^{\star}(s)$

# Maximum Entropy RL: Soft Value Iteration

$$\arg \min_{\pi} \mathbb{E}_{s,a \sim d_\mu^\pi} \left[ c(s,a) + \ln \pi(a \mid s) \right]$$

## Soft Value Iteration for finite horizon MDP (continue)

$$Q_h^\star(s,a) = c(s,a) + \mathbb{E}_{s' \sim P(\cdot \mid s,a)} V_{h+1}^\star(s')$$

$$\pi_h^\star(a \mid s) = \frac{\exp(-Q_h^\star(s,a))}{\sum_{a'} \exp(-Q_h^\star(s,a'))}$$

$$V_h^\star(s) = \mathbb{E}_{a \sim \pi_h^\star(\cdot \mid s)} \left[ \ln \pi_h^\star(a \mid s) + Q_h^\star(s,a) \right] = -\ln \left( \sum_a \exp \left( -Q_h^\star(s,a) \right) \right)$$

# Maximum Entropy RL: Soft Value Iteration

$$\arg \min_{\pi} \mathbb{E}_{s,a \sim d_\mu^\pi} \left[ c(s,a) + \ln \pi(a \mid s) \right]$$

## Soft Value Iteration for finite horizon MDP (continue)

$$Q_h^\star(s,a) = c(s,a) + \mathbb{E}_{s' \sim P(\cdot \mid s,a)} V_{h+1}^\star(s')$$

$$\pi_h^\star(a \mid s) = \frac{\exp(-Q_h^\star(s,a))}{\sum_{a'} \exp(-Q_h^\star(s,a'))} \quad \checkmark$$

$$V_h^\star(s) = \mathbb{E}_{a \sim \pi_h^\star(\cdot \mid s)} \left[ \ln \pi_h^\star(a \mid s) + Q_h^\star(s,a) \right] = -\ln \left( \sum_a \exp \left( -Q_h^\star(s,a) \right) \right)$$

classic RL:

$$\left( V_h^\star(s) = \min_a Q_h^\pi(s a) \right)$$

(contrast this to $\min_a Q^\star(s,a)$ )

① $Q_h^\star(s,a) = 0$

$Q_h^\star(s,a') = +\infty$

$\Rightarrow V_h^\star(s) = Q_h^\star(s,a) = 0$

② $Q_h^\star(s,a) = Q_h^\star(s,a')$

# Maximum Entropy RL: Soft Value Iteration

$$\arg \min_{\pi} \mathbb{E}_{s,a \sim d_{\mu}^{\pi}} \left[ c(s,a) + \ln \pi(a \mid s) \right]$$

## Soft Value Iteration for finite horizon MDP (continue)

$$Q_h^{\star}(s,a) = c(s,a) + \mathbb{E}_{s' \sim P(\cdot \mid s,a)} V_{h+1}^{\star}(s')$$

$$\pi_h^{\star}(a \mid s) = \frac{\exp(-Q_h^{\star}(s,a))}{\sum_{a'} \exp(-Q_h^{\star}(s,a'))}$$

$$V_h^{\star}(s) = \mathbb{E}_{a \sim \pi_h^{\star}(\cdot \mid s)} \left[ \ln \pi_h^{\star}(a \mid s) + Q_h^{\star}(s,a) \right] = -\ln \left( \sum_{a} \exp\left(-Q_h^{\star}(s,a)\right) \right)$$

(contrast this to $\min_{a} Q^{\star}(s,a)$ )

$$Q_{h-1}^{\star}(s,a) = c(s,a) + \mathbb{E}_{s' \sim P(\cdot \mid s,a)} V_h^{\star}(s')$$

$$\bullet \ \bullet \ \bullet$$

# Plan for Today:

✅ 1. The Iterative Algorithm framework

*finite Horizon*

$$\pi^* = \{\pi_0^*, \ldots, \pi_{M-1}^*\}$$

✅ 2. How to compute best response: Soft Value Iteration (DP again)

$$\arg\min_{\pi} \mathbb{E}_{s,a\sim d_\mu^\pi}\left[c(x,a) + \ln\pi(a\,|\,s)\right]$$

$$\Leftrightarrow \quad \min_{\pi} \mathbb{E}_{s\sim d_\mu^\pi}\left[c(s\text{-}a) - \text{Entropy}\,(\pi(\cdot|s))\right.$$

3. The MaxEnt-IRL algorithm

# Maximum Entropy Inverse RL Algorithm framework

$$\max_{w \in \mathbb{R}^d} \min_{\pi} \underbrace{\mathbb{E}_{s,a \sim d_\mu^\pi} \ln \pi(a \mid s) + w^\top \left( \mathbb{E}_{s,a \sim d_\mu^\pi} \phi(s,a) - \mathbb{E}_{s,a \sim d_\mu^{\pi^\star}} \phi(s,a) \right)}_{:= \ell(\pi, w)}$$

Initialize $w^0 \in \mathbb{R}^d$   $\left\{ \pi_0 \cdots \pi_{|u|} \right\}$

For $t = 0 \to T - 1$   $\overset{\infty}{\underset{\Delta}{DP}}$

$\pi^t =$ soft-VI $\left( c(s,a) := (w^t)^\top \phi(s,a) \right)$   $\widehat{(H)}$   (# best response: $\pi^t = \arg\min_\pi \ell(\pi, w^t)$)

$$w^{t+1} = w^t + \eta \left( \underbrace{\mathbb{E}_{s,a \sim d_\mu^{\pi^t}} \phi(s,a) - \mathbb{E}_{s,a \sim d_\mu^{\pi^\star}} \phi(s,a)} \right)$$

Return $\bar{\pi} = $ Uniform$(\pi^0, \ldots, \pi^{T-1})$   $\overset{\widetilde{\sim 0}}{\uparrow}$   (# gradient update: $w^{t+1} = w^t + \eta \nabla_w \ell(\pi^t, w^t)$)
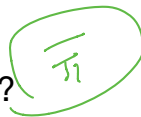
$$\left\| \mathbb{E}_{s,a \sim d_\mu^{\pi^t}} \phi(s,a) - \mathbb{E}_{s,a \sim d_\mu^{\pi^u}} \phi(s,a) \right\|_2 \leq \delta$$

# Maximum Entropy IRL: Calculate Trajectory Likelihood

Given a trajectory $\tau = \{s_0, a_0, \ldots, s_{H-1}, a_{H-1}\}$

What's the likelihood of $\tau$ being generated by expert?

# Maximum Entropy IRL: Calculate Trajectory Likelihood

Given a trajectory $\tau = \{s_0, a_0, \ldots, s_{H-1}, a_{H-1}\}$

What's the likelihood of $\tau$ being generated by expert?

$$\ln\left(\rho^{\bar{\pi}}(\tau)\right) = \sum_{h=0}^{H-1} \left[\ln P(s_{h+1} \mid s_h, a_h) + \ln \bar{\pi}(a_h \mid s_h)\right]$$

known        our policy from MaxEnt-IRL

$\mu(s_0) \, \bar{\pi}(a_0 \mid s_0) \, P(s_1 \mid s_0, a_0) \cdots$

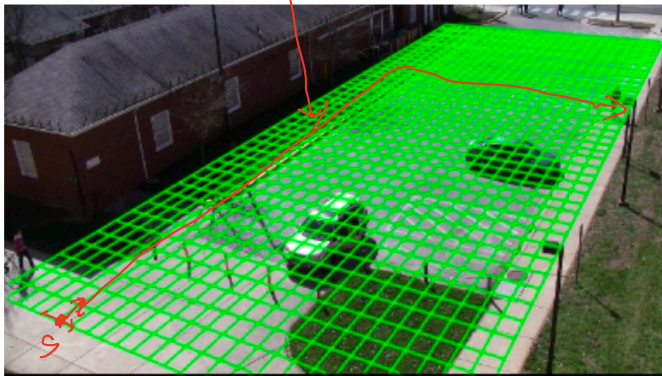# Maximum Entropy IRL: Calculate Trajectory Likelihood

Given a trajectory $\tau = \{s_0, a_0, \ldots, s_{H-1}, a_{H-1}\}$

What's the likelihood of $\tau$ being generated by expert?

$$\ln\left(\rho^{\bar{\pi}}(\tau)\right) = \sum_{h=0}^{H-1} \left[\ln P(s_{h+1} \,|\, s_h, a_h) + \ln \bar{\pi}(a_h \,|\, s_h)\right]$$

State space: grid,
action space: 4 actions

$P(s_{h+1} \,|\, s_h, a_h)$ is deterministic

# Summary for Today:

1. Maximum Entropy IRL framework

$$\max_{\pi} \mathbb{E}_{s \sim d_{\mu}^{\pi}} \left[ \text{entropy} \left( \pi(\cdot \mid s) \right) \right]$$

$$s.t, \mathbb{E}_{s,a \sim d_{\mu}^{\pi}} \phi(s,a) = \mathbb{E}_{s,a \sim d_{\mu}^{\pi^{\star}}} \phi(s,a)$$

feature Matching

# Summary for Today:

1. Maximum Entropy IRL framework

$$\max_{\pi} \mathbb{E}_{s \sim d_{\mu}^{\pi}} \left[ \text{entropy} \left( \pi( \cdot \mid s) \right) \right]$$

$$s.t, \mathbb{E}_{s,a \sim d_{\mu}^{\pi}} \phi(s,a) = \mathbb{E}_{s,a \sim d_{\mu}^{\pi^{\star}}} \phi(s,a)$$

← Best Response

← Soft Value Iteration

2. Inside MaxEnt-IRL, we perform Maximum Entropy RL:

Note

$$\min_{\pi_0, \ldots, \pi_{H-1}} \mathbb{E} \left[ \sum_{h=0}^{H-1} \left( c(s_h, a_h) - \text{entropy}(\pi_h( \cdot \mid s_h)) \right) \mid s_0 \sim \mu, a_h \sim \pi_h( \cdot \mid s_h) \right]$$