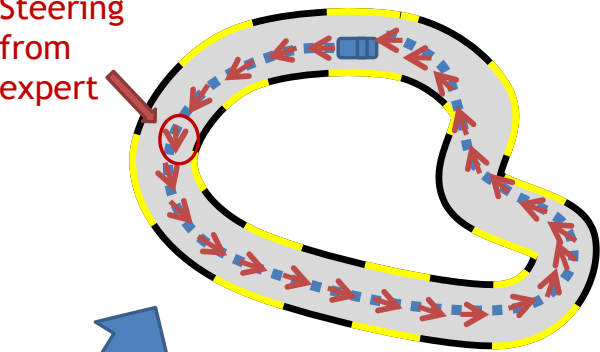


# Maximum Entropy IRL

# Dagger Recap

At iteration  $t$ , given  $\pi^t$

Steering from expert



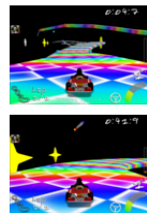
New policy

$\pi^{t+1}$

Aggregate Dataset

Supervised Learning

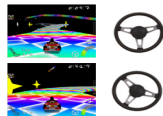
New Data



⋮



All previous data



⋮

$\pi^t$

$a \sim \pi^t(c, |s)$

Data Aggregation = Follow-the-Regularized-Leader Online Learner

## DAgger Performance Recap:

DAgger finds a policy  $\hat{\pi}$  such that it **matches to  $\pi^*$  under its own  $d_{\mu}^{\hat{\pi}}$**

$$\mathbb{E}_{s \sim d_{\mu}^{\hat{\pi}}} [\mathbf{1}\{\hat{\pi}(s) \neq \pi^*(s)\}] \leq \epsilon_{reg} = O(1/\sqrt{T})$$

$\uparrow$  # of iterations

## DAgger Performance Recap:

DAgger finds a policy  $\hat{\pi}$  such that it **matches to  $\pi^*$  under its own  $d_{\mu}^{\hat{\pi}}$**

$$\mathbb{E}_{s \sim d_{\mu}^{\hat{\pi}}} [\mathbf{1}\{\hat{\pi}(s) \neq \pi^*(s)\}] \leq \epsilon_{reg} = O(1/\sqrt{T})$$

If expert herself can **quickly recover from a deviation**, i.e.,  $|Q^{\pi^*}(s, a) - V^{\pi^*}(s)|$  is small for all  $s$ ,

$$V^{\pi^*} - V^{\hat{\pi}} \leq O\left(\frac{1}{1-\gamma} \cdot \epsilon_{reg}\right)$$

*linear*

BC:  $\left(\frac{1}{1-\gamma}\right)^2 \Sigma$

## DAgger Performance Recap:

DAgger finds a policy  $\hat{\pi}$  such that it **matches to  $\pi^*$  under its own  $d_{\mu}^{\hat{\pi}}$**

$$\mathbb{E}_{s \sim d_{\mu}^{\hat{\pi}}} [\mathbf{1}\{\hat{\pi}(s) \neq \pi^*(s)\}] \leq \epsilon_{reg} = O(1/\sqrt{T})$$

If expert herself can **quickly recover from a deviation**, i.e.,  $|Q^{\pi^*}(s, a) - V^{\pi^*}(s)|$  is small for all  $s$ ,

$$V^{\pi^*} - V^{\pi^t} \leq O\left(\frac{1}{1-\gamma} \cdot \epsilon_{reg}\right)$$

**This is a significant improvement over BC in both theory and practice**

## Plan for Today:

1. The principle of Maximum Entropy ✓

2. Constrained Optimization ✓

2. The Algorithm: Maximum Entropy Inverse RL ✓

## Setting

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

$\Delta$

cost  $C(s,a) \rightarrow [0, 1]$

$\tau$  explore policy

## Setting

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

- (1) Ground truth cost  $c(s, a)$  is unknown;
- (2) assume expert is the optimal policy  $\pi^\star$  of the cost  $c$
- (3) **transition  $P$  is known**



## Setting

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

- (1) Ground truth cost  $c(s, a)$  is unknown;
- (2) assume expert is the optimal policy  $\pi^\star$  of the cost  $c$
- (3) **transition  $P$  is known**

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

# Setting

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

- (1) Ground truth cost  $c(s, a)$  is unknown;
- (2) assume expert is the optimal policy  $\pi^\star$  of the cost  $c$
- (3) **transition  $P$  is known**

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

**Key Assumption on cost:**

$c(s, a) = \langle \theta^\star, \phi(s, a) \rangle$ , linear w.r.t feature  $\phi(s, a)$

$\Delta$   
learn  $\theta^\star$  ??  
 $\phi: S \times A \mapsto \mathbb{R}^d$

## Running Example: Define feature map

**Key Assumption on cost:**

$$c(s, a) = \langle \theta^*, \phi(s, a) \rangle, \text{ linear wrt feature } \phi(s, a)$$

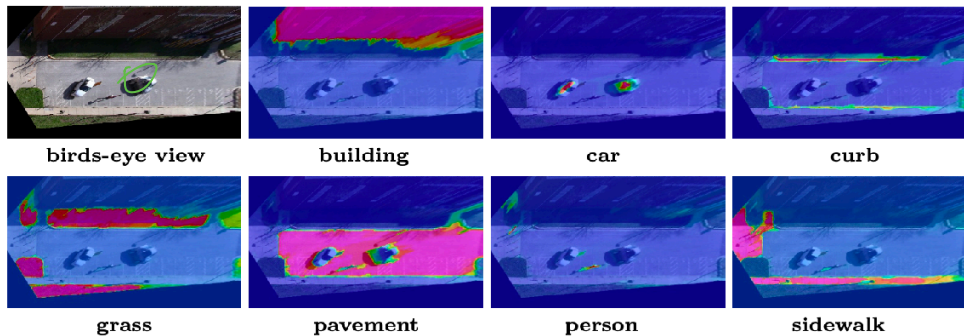


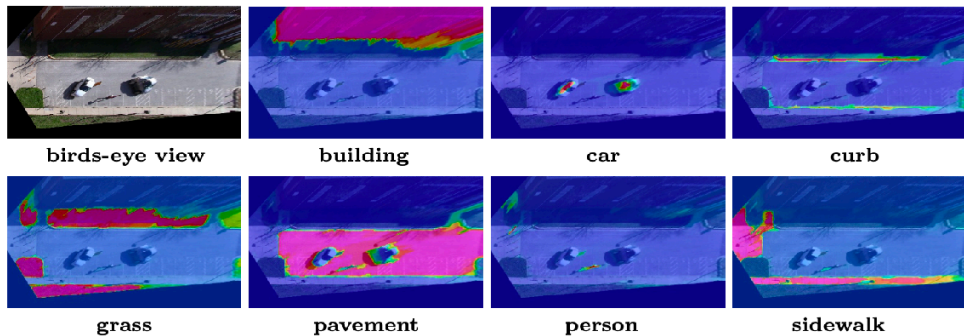
Fig. 4. Classifier feature response maps. Top left is the original image.

## Running Example: Define feature map

**Key Assumption on cost:**

$$c(s, a) = \langle \theta^*, \phi(s, a) \rangle, \text{ linear wrt feature } \phi(s, a)$$

State  $s$ : pixel or a group of neighboring pixels in image)



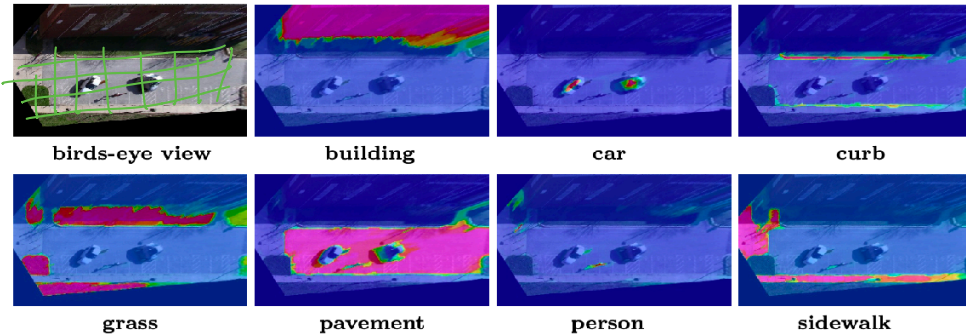
**Fig. 4.** Classifier feature response maps. Top left is the original image.

# Running Example: Define feature map

**Key Assumption on cost:**

$c(s, a) = \langle \theta^*, \phi(s, a) \rangle$ , linear wrt feature  $\phi(s, a)$

State  $s$ : pixel or a group of neighboring pixels in image



pixel  
 $\downarrow$   
 $\Delta$

$$\phi(s, a) = \begin{bmatrix} \mathbb{P}(\text{pixels being building}) \\ \mathbb{P}(\text{pixels being grass}) \\ \mathbb{P}(\text{pixels being sidewalk}) \\ \mathbb{P}(\text{pixels being car}) \\ \dots \end{bmatrix}$$

Fig. 4. Classifier feature response maps. Top left is the original image.

$$E_{s \sim d_{\mu}^{\pi}} \phi(s, a) \approx E_{s \sim d_{\mu}^{\pi}} [\phi(s, a)] \quad \checkmark$$

$$\theta^{*T} \left( E_{s \sim d_{\mu}^{\pi}} \phi(s, a) \right) \approx \theta^{*T} \left( E_{s \sim d_{\mu}^{\pi}} \phi(s, a) \right)$$

# Running Example: Define feature map

**Key Assumption on cost:**

$$c(s, a) = \langle \theta^*, \phi(s, a) \rangle, \text{ linear wrt feature } \phi(s, a)$$

State  $s$ : pixel or a group of neighboring pixels in image

$$\phi(s, a) = \begin{bmatrix} \mathbb{P}(\text{pixels being building}) \\ \mathbb{P}(\text{pixels being grass}) \\ \mathbb{P}(\text{pixels being sidewalk}) \\ \mathbb{P}(\text{pixels being car}) \\ \dots \end{bmatrix}$$

Maybe colliding with cars or buildings has **high** cost, but walking on sidewalk or grass has **low** cost

$$\theta^{*T} \phi(s, a)$$

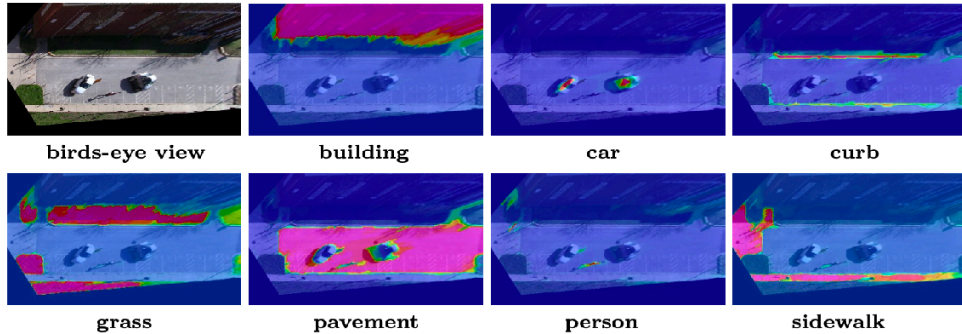


Fig. 4. Classifier feature response maps. Top left is the original image.

## Notation on Distributions

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

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

$\rho^\pi(\tau) := \mu_0(s_0)\pi(a_0 | s_0)P(s_1 | s_0, a_0)\pi(a_1 | s_1)\dots\pi(a_{H-1} | s_{H-1})P(s_H | s_{H-1}, a_{H-1})$ :  
Likelihood of the trajectory  $\tau$  under  $\pi$ , i.e., the prob of  $\pi$  generating  $\tau$

## Detour: Principle of Maximum Entropy

Definition of the Entropy of a distribution:



## Detour: Principle of Maximum Entropy

Definition of the Entropy of a distribution:

Given a distribution  $P \in \Delta(X)$ , the entropy is defined as:

$$\text{Entropy}(P) = - \sum_x P(x) \cdot \ln P(x) = \mathbb{E}_{x \sim P} [-\ln P(x)]$$

## Detour: Principle of Maximum Entropy

Definition of the Entropy of a distribution:

Given a distribution  $P \in \Delta(X)$ , the entropy is defined as:

$$\text{Entropy}(P) = - \sum_x P(x) \cdot \ln P(x)$$

**Higher entropy means more uncertainty, for instance:**

## Detour: Principle of Maximum Entropy

Definition of the Entropy of a distribution:

Given a distribution  $P \in \Delta(X)$ , the entropy is defined as:

$$\text{Entropy}(P) = - \sum_x P(x) \cdot \ln P(x)$$

**Higher entropy means more uncertainty, for instance:**

*$U(X)$  is uniform dist over  $X$*

Uniform distribution has the highest entropy,

i.e.,  $\text{Entropy}(U(X)) = - \sum_x \underbrace{(1/|X|)} \ln(1/|X|) = \ln(|X|)$  ✓

## Detour: Principle of Maximum Entropy

Definition of the Entropy of a distribution:

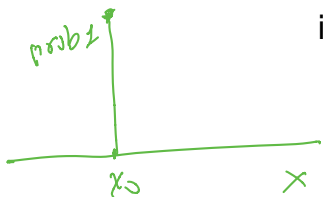
Given a distribution  $P \in \Delta(X)$ , the entropy is defined as:

$$\text{Entropy}(P) = - \sum_x P(x) \cdot \ln P(x)$$

Higher entropy means more uncertainty, for instance:

Uniform distribution has the highest entropy,  
i.e.,  $\text{Entropy}(U(X)) = - \sum_x (1/|X|) \ln(1/|X|) = \ln(|X|)$

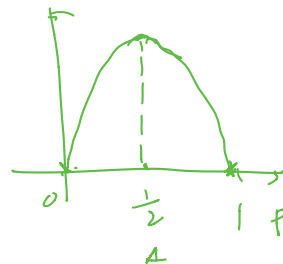
Deterministic distribution has zero entropy:  
i.e.,  $\text{Entropy}(\delta(x_0)) = - \underbrace{1 \cdot \ln 1} - \sum_{x \neq x_0} \underbrace{0 \ln 0} = 0$



with  
 $P \in [0,1], \pm 1$

with  
 $\downarrow P, -1$

entrop



## Detour: Principle of Maximum Entropy

We want to find a distribution whose mean and covariance matrix equal to  $\mu, \Sigma$ ,  
but there are infinitely many such distributions...

## Detour: Principle of Maximum Entropy

We want to find a distribution whose mean and covariance matrix equal to  $\mu, \Sigma$ ,  
but there are infinitely many such distributions...

Principle of Maximum Entropy:  
Entropy Maximization subject to Moment Matching constraints

## Detour: Principle of Maximum Entropy

We want to find a distribution whose mean and covariance matrix equal to  $\mu, \Sigma$ ,  
but there are infinitely many such distributions...

Principle of Maximum Entropy:

Entropy Maximization subject to Moment Matching constraints

$$\max_{P \in \Delta(X)} \text{entropy}(P), \quad \text{s.t.}, \quad \underbrace{\mathbb{E}_{x \sim P}[x]} = \mu, \quad \underbrace{\mathbb{E}_{x \sim P}[xx^T]} = \Sigma + \mu\mu^T$$

$\checkmark \mathbb{E}_{x \sim P} [(x - \mu)(x - \mu)^T]$

## Detour: Principle of Maximum Entropy

We want to find a distribution whose mean and covariance matrix equal to  $\mu, \Sigma$ ,  
but there are infinitely many such distributions...

Principle of Maximum Entropy:  
Entropy Maximization subject to Moment Matching constraints

$$\max_{P \in \Delta(X)} \text{entropy}(P), \quad \text{s.t.}, \quad \mathbb{E}_{x \sim P}[x] = \mu, \quad \mathbb{E}_{x \sim P}[xx^T] = \Sigma + \mu\mu^T$$

Solution:  $P^* = \mathcal{N}(\mu, \Sigma)$   
(proof: out of scope)

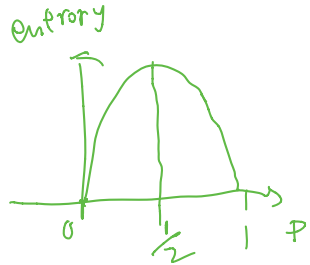


# Detour: Principle of Maximum Entropy

In summary:

Maximum Entropy Principle says that:

Among the distributions that satisfy pre-defined constraints (mean & variance),  
let's pick the one that is the most uncertain  
(uncertainty measured in entropy)



## Plan for Today:



1. The principle of Maximum Entropy

2. Constrained Optimization

3. The Algorithm: Maximum Entropy Inverse RL

## Constrained Optimization:

Consider the following constrained optimization problem:

$$\begin{aligned} & \min_x f(x) \\ & \text{s.t.}, \underline{g_1(x) = 0, \quad g_2(x) = 0} \end{aligned}$$

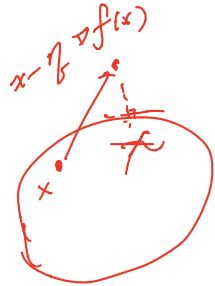
$$\begin{aligned} & \min_{x,y} x+y \\ & \text{s.t.} \quad x^2+y^2-1=0 \end{aligned}$$

## Constrained Optimization:

Consider the following constrained optimization problem:

$$\begin{aligned} \min_x f(x) \\ \text{s.t.}, g_1(x) = 0, \quad g_2(x) = 0 \end{aligned}$$

Denote  $x^*$  as the optimal solution here.



## Constrained Optimization:

Consider the following constrained optimization problem:

$$\begin{aligned} & \min_x f(x) \\ & \text{s.t.}, g_1(x) = 0, \quad g_2(x) = 0 \end{aligned}$$

Denote  $x^*$  as the optimal solution here.

How to solve such constrained optimization problem?

## Constrained Optimization:

$g_1$   
 $g_2$

Define two Lagrange multiplier  $w_1, w_2 \in \mathbb{R}$ , we consider the following Lagrange formulation:

$$\min_x \left[ \max_{w_1, w_2} \underbrace{f(x) + w_1 g_1(x) + w_2 g_2(x)} \right]$$

$\leftarrow$   $\min f(x)$   
s.t.  $g_1(x) \leq 0$   
 $g_2(x) \leq 0$

## Constrained Optimization:

Define two Lagrange multiplier  $w_1, w_2 \in \mathbb{R}$ , we consider the following Lagrange formulation:

$$\min_x \left[ \max_{w_1, w_2} f(x) + w_1 g_1(x) + w_2 g_2(x) \right]$$

For any  $x$  that does not satisfy constraints, i.e.,  $g_1(x) \neq 0$  or  $g_2(x) \neq 0$ ,

we must have:  $\max_{w_1, w_2} f(x) + w_1 g_1(x) + w_2 g_2(x) = +\infty$

$$g_1(x) = 5$$

$$w_1 = +\infty$$

## Constrained Optimization:

Define two Lagrange multiplier  $w_1, w_2 \in \mathbb{R}$ , we consider the following Lagrange formulation:

$$\min_x \left[ \max_{w_1, w_2} f(x) + w_1 g_1(x) + w_2 g_2(x) \right]$$

*↙ infeasible*

For any  $x$  that does not satisfy constraints, i.e.,  $g_1(x) \neq 0$  or  $g_2(x) \neq 0$ ,  
we must have:  $\max_{w_1, w_2} f(x) + w_1 g_1(x) + w_2 g_2(x) = +\infty$

*↙ feasible*

For any  $x$  that satisfies constraints, i.e.,  $g_1(x) = 0$  and  $g_2(x) = 0$ ,  
we must have:  $\max_{w_1, w_2} f(x) + w_1 g_1(x) + w_2 g_2(x) = f(x)$

*$g_1(x) = 0$        $g_2(x) = 0$*



## Constrained Optimization:

Define two Lagrange multiplier  $w_1, w_2 \in \mathbb{R}$ , we consider the following Lagrange formulation:

$$\min_x \left[ \max_{w_1, w_2} f(x) + w_1 g_1(x) + w_2 g_2(x) \right]$$

In other words,

$$\max_{w_1, w_2} f(x) + w_1 g_1(x) + w_2 g_2(x) = \begin{cases} +\infty & g_1(x) \neq 0 \text{ or } g_2(x) \neq 0 \text{ i.e., infeasible} \\ f(x) & g_1(x) = g_2(x) = 0 \text{ i.e., feasible} \end{cases}$$

## Constrained Optimization:

Define two Lagrange multiplier  $w_1, w_2 \in \mathbb{R}$ , we consider the following Lagrange formulation:

$$\min_x \left[ \max_{w_1, w_2} f(x) + w_1 g_1(x) + w_2 g_2(x) \right]$$

In other words,

$$\max_{w_1, w_2} f(x) + w_1 g_1(x) + w_2 g_2(x) = \begin{cases} +\infty & g_1(x) \neq 0 \text{ or } g_2(x) \neq 0 \text{ i.e., infeasible} \\ f(x) & g_1(x) = g_2(x) = 0 \text{ i.e., feasible} \end{cases}$$

$-\infty$  infeasible  
 $f(x)$  feasible

Thus, solving the Lagrange formulation is equivalent to the original formulation:

$$\arg \min_x \left[ \max_{w_1, w_2} f(x) + w_1 g_1(x) + w_2 g_2(x) \right] = x^*$$

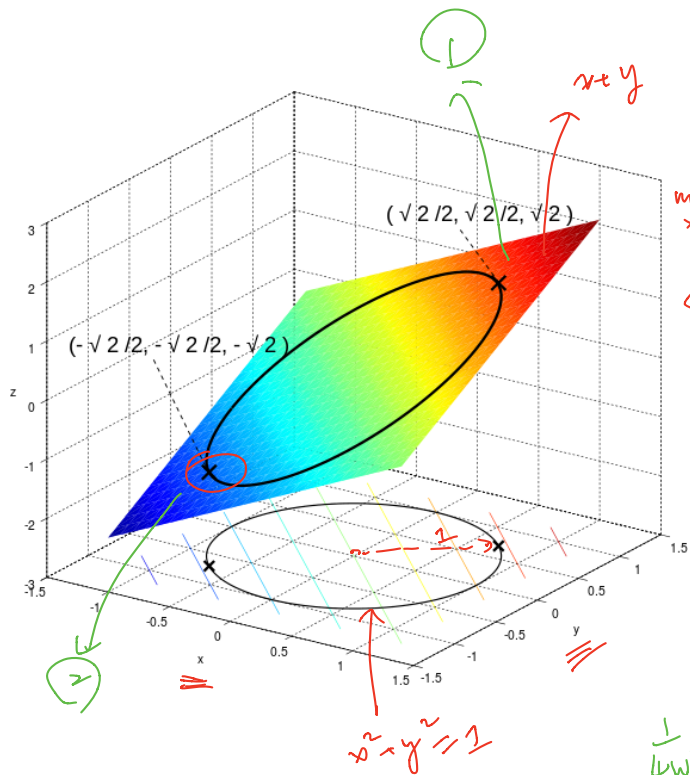
## Constrained Optimization:

In summary, we have that

$$\arg \min_x \left[ \max_{w_1, w_2} f(x) + w_1 g_1(x) + w_2 g_2(x) \right] = x^*$$

Where  $x^*$  is the optimal solution of the original constrained program:

$$\begin{aligned} & \min_x f(x) \\ & s.t., g_1(x) = 0, \quad g_2(x) = 0 \end{aligned}$$



$$x, y \in \mathbb{R}$$

Example:

Introduce WEP

$$\min_{x,y} x + y, \text{ s.t. } x^2 + y^2 = 1$$

$$\min_{x,y} \max_w (x + y + w(x^2 + y^2 - 1))$$

$$\Leftrightarrow \min_{x,y} \max_w l(x, y, w)$$

$$\checkmark \frac{dl}{dx}(x, y, w) = 1 + 2wx = 0 \Rightarrow x = -\frac{1}{2w}$$

$$\checkmark \frac{dl}{dy}(x, y, w) = 1 + 2wy = 0 \Rightarrow y = -\frac{1}{2w}$$

$$\checkmark \frac{dl}{dw}(x, y, w) = x^2 + y^2 - 1 = 0 \Rightarrow x^2 + y^2 = 1$$

$$\frac{1}{4w^2} + \frac{1}{4w^2} = 1 \Rightarrow w = \pm \sqrt{\frac{1}{2}} \Rightarrow (x, y) = \begin{cases} (\frac{1}{\sqrt{2}}, \frac{1}{\sqrt{2}}) \text{ (1)} \\ (-\frac{1}{\sqrt{2}}, -\frac{1}{\sqrt{2}}) \text{ (2)} \end{cases}$$

## Constrained Optimization:

We will often be interested in solving the dual version, i.e.,

$$\max_{w_1, w_2} \min_x \underbrace{f(x) + w_1 g_1(x) + w_2 g_2(x)}_{:=\ell(x, w)}$$

$$\min_x f(x)$$

$$\text{s.t. } g_1(x) \leq 0$$

$$g_2(x) \leq 0$$

## Constrained Optimization:

We will often be interested in solving the dual version, i.e.,

$$\max_{w_1, w_2} \min_x \underbrace{f(x) + w_1 g_1(x) + w_2 g_2(x)}_{:=\ell(x, w)}$$

**And one procedure to solve a max min is the following iterative algorithm:**

## Constrained Optimization:

We will often be interested in solving the dual version, i.e.,

$$\max_{w_1, w_2} \min_x \underbrace{f(x) + w_1 g_1(x) + w_2 g_2(x)}_{:=\ell(x, w)}$$

**And one procedure to solve a max min is the following iterative algorithm:**

**Initialize** Lagrange multipliers  $w_1^0, w_2^0$

## Constrained Optimization:

We will often be interested in solving the dual version, i.e.,

$$\max_{w_1, w_2} \min_x \underbrace{f(x) + w_1 g_1(x) + w_2 g_2(x)}_{:=\ell(x, w)}$$

**And one procedure to solve a max min is the following iterative algorithm:**

**Initialize** Lagrange multipliers  $w_1^0, w_2^0$

For  $t = 0 \rightarrow T - 1$



## Constrained Optimization:

We will often be interested in solving the dual version, i.e.,

$$\max_{w_1, w_2} \min_x \underbrace{f(x) + w_1 g_1(x) + w_2 g_2(x)}_{:= \ell(x, w)}$$

And one procedure to solve a max min is the following iterative algorithm:

**Initialize** Lagrange multipliers  $w_1^0, w_2^0$

For  $t = 0 \rightarrow T - 1$

$$x^t = \arg \min_x f(x) + \underbrace{w_1^t}_{\text{green}} g_1(x) + \underbrace{w_2^t}_{\text{green}} g_2(x) \quad (\# \text{ best response: } \underbrace{\arg \min_x \ell(x, w^t)}_{\text{green}})$$

$$w^t = \begin{bmatrix} w_1^t \\ w_2^t \end{bmatrix}$$

## Constrained Optimization:

We will often be interested in solving the dual version, i.e.,

$$\max_{w_1, w_2} \min_x \underbrace{f(x) + w_1 g_1(x) + w_2 g_2(x)}_{:=\ell(x, w)}$$

**And one procedure to solve a max min is the following iterative algorithm:**

**Initialize** Lagrange multipliers  $w_1^0, w_2^0$

For  $t = 0 \rightarrow T - 1$

$$x^t = \arg \min_x f(x) + w_1^t g_1(x) + w_2^t g_2(x) \quad (\# \text{ best response: } \arg \min_x \ell(x, w^t))$$

$$w_1^{t+1} = w_1^t + \eta g_1(x^t)$$

$$w_2^{t+1} = w_2^t + \eta g_2(x^t)$$

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

$w = w_1, w_2$

## Constrained Optimization:

We will often be interested in solving the dual version, i.e.,

$$\max_{w_1, w_2} \min_x \underbrace{f(x) + w_1 g_1(x) + w_2 g_2(x)}_{:=\ell(x, w)}$$

And one procedure to solve a max min is the following iterative algorithm:

**Initialize** Lagrange multipliers  $w_1^0, w_2^0$

For  $t = 0 \rightarrow T - 1$

$$x^t = \arg \min_x f(x) + w_1^t g_1(x) + w_2^t g_2(x) \quad (\# \text{ best response: } \arg \min_x \ell(x, w^t))$$

$$w_1^{t+1} = w_1^t + \eta g_1(x^t)$$

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

$$w_2^{t+1} = w_2^t + \eta g_2(x^t)$$

**Return:**  $\bar{x} = \sum_{t=0}^{T-1} x_t / T$

## Constrained Optimization:

We will often be interested in solving the dual version, i.e.,

$$\max_{w_1, w_2} \min_x \underbrace{f(x) + w_1 g_1(x) + w_2 g_2(x)}_{:= \ell(x, w)}$$

And one procedure to solve a max min is the following iterative algorithm:

**Initialize** Lagrange multipliers  $w_1^0, w_2^0$

For  $t = 0 \rightarrow T - 1$

$$x^t = \arg \min_x f(x) + w_1^t g_1(x) + w_2^t g_2(x) \quad (\# \text{ best response: } \arg \min_x \ell(x, w^t))$$

$$w_1^{t+1} = w_1^t + \eta g_1(x^t)$$

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

$$w_2^{t+1} = w_2^t + \eta g_2(x^t)$$

**Return:**  $\bar{x} = \sum_{t=0}^{T-1} x_t / T$     Informal theorem: when  $f, g$  are convex,  $\bar{x} \rightarrow x^*$ , as  $T \rightarrow \infty$

## Plan for Today:

✓ 1. The principle of Maximum Entropy

✓ 2. Constrained Optimization

3. The Algorithm: Maximum Entropy Inverse RL

$$\min f(x)$$

$$\text{s.t. } g(x) \leq 0$$

$$h(x) = 0$$

## Maximum Entropy Inverse RL:

## Maximum Entropy Inverse RL:

$$(s_i^*, a_i^*) \sim d_{\pi^*}$$
$$\approx \frac{1}{N} \sum_{i=1}^N \phi(s_i^*, a_i^*)$$

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

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

But there are potentially many such policies...

$$\mathbb{E}_{s,a \sim d_{\pi}} \left( \theta^T \phi(s,a) \right) - \mathbb{E}_{s,a \sim d_{\pi^*}} \left( \theta^T \phi(s,a) \right)$$
$$= \theta^T \left[ \underbrace{\mathbb{E}_{s,a \sim d_{\pi}} \phi(s,a) - \mathbb{E}_{s,a \sim d_{\pi^*}} \phi(s,a)}_{=0} \right]$$

## Maximum Entropy Inverse RL:

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

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

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:



## Maximum Entropy Inverse RL:

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

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

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:

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

*Encourage diversity*

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

## Maximum Entropy Inverse RL:

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

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

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:

$$\begin{aligned} & \max_{\pi} \mathbb{E}_{s \sim d_{\mu}^{\pi}} \left[ \text{entropy}(\pi(\cdot | s)) \right] \\ & s.t., \mathbb{E}_{s,a \sim d_{\mu}^{\pi}} \phi(s, a) = \mathbb{E}_{s,a \sim d_{\mu}^{\pi^*}} \phi(s, a) \end{aligned}$$

This can be estimated using expert data:

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

## Maximum Entropy Inverse RL:

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

## Maximum Entropy Inverse RL:

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

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

## Maximum Entropy Inverse RL:

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

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

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

## Maximum Entropy Inverse RL:

We arrive at the following constraint optimization problem:

$$\begin{aligned} & \arg \min_{\pi} \mathbb{E}_{s,a \sim d_{\mu}^{\pi}} \ln \pi(a | s) \\ & s . t, \mathbb{E}_{s,a \sim d_{\mu}^{\pi}} \phi(s, a) = \mathbb{E}_{s,a \sim d_{\mu}^{\pi^*}} \phi(s, a) \end{aligned}$$

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

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

## Maximum Entropy Inverse RL:

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

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

Next lecture,  
we will design algorithm (in high level, it is the iterative algorithm framework)  
for this max – min problem