# Maximum Entropy
# Inverse Reinforcement Learning

## Sham Kakade and Wen Sun

**CS 6789: Foundations of Reinforcement Learning**

# Announcements

Project Presentation (Dec 8th and 10th):

Please sign up time slots
(see Piazza post for more details)

# Recap

**Offline IL Setting:**

# Recap

**Offline IL Setting:**

Ground truth reward $r(s, a) \in [0,1]$ is unknown; assume expert is a near optimal policy $\pi^\star$

# Recap

**Offline IL Setting:**

Ground truth reward $r(s, a) \in [0,1]$ is unknown; assume expert is a near optimal policy $\pi^\star$

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

# Recap

**Offline IL Setting:**

Ground truth reward $r(s,a) \in [0,1]$ is unknown; assume expert is a near optimal policy $\pi^\star$

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

**Offline IL Algorithm: Behavior Cloning (Maximum Likelihood)**

# Recap

**Offline IL Setting:**

Ground truth reward $r(s, a) \in [0,1]$ is unknown; assume expert is a near optimal policy $\pi^\star$

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

**Offline IL Algorithm: Behavior Cloning (Maximum Likelihood)**

$$\widehat{\pi} = \arg\max_{\pi \in \Pi} \sum_{i=1}^M \ln \pi(a_i^\star \mid s_i^\star)$$

# Recap

**Hybrid IL Setting:**

# Recap

**Hybrid IL Setting:**

We have a dataset $\mathcal{D} = (s_i^\star, a_i^\star)_{i=1}^M \sim d^{\pi^\star}$ **<u>and</u>** access to transition $P(\,\cdot\,|\,s, a), \forall s, a$

# Recap

**Hybrid IL Setting:**

We have a dataset $\mathcal{D} = (s_i^\star, a_i^\star)_{i=1}^M \sim d^{\pi^\star}$ **<u>and</u>** access to transition $P(\cdot \mid s, a), \forall s, a$

**Hybrid IL Algorithm: Distribution Matching**
**(Statistically efficient, but not computationally)**

Start from TV: $\mathcal{F} = \{ f : \|f\|_\infty \leq 1 \}$

**Recap**

$\forall \pi, \pi' \in \Pi$

$f_{\pi,\pi'} = \arg\max_{f \in \mathcal{F}} \mathbb{E}_{s,a \sim d^\pi} f(s,a) - \mathbb{E}_{s,a \sim d^{\pi'}} f(s,a)$

**Hybrid IL Setting:**

$\widetilde{\mathcal{F}} = \{ f_{\pi,\pi'} : \pi, \pi' \in \Pi \}$

We have a dataset $\mathcal{D} = (s_i^\star, a_i^\star)_{i=1}^M \sim d^{\pi^\star}$ **and** access to transition $P(\cdot \mid s, a), \forall s, a$

$\mathcal{F} \to \widetilde{\mathcal{F}} \quad |\widetilde{\mathcal{F}}| \leq |\Pi|^2$

**Hybrid IL Algorithm: Distribution Matching**
**(Statistically efficient, but not computationally)**

$$\hat{\pi} := \arg\min_{\pi \in \Pi} \left[ \max_{f \in \widetilde{\mathcal{F}}} \left[ \mathbb{E}_{s,a \sim d^\pi} f(s,a) - \frac{1}{M} \sum_{i=1}^M f(s_i^\star, a_i^\star) \right] \right]$$

$\underbrace{\qquad\qquad\qquad\qquad\qquad\qquad\qquad}_{}$ IPM($\widetilde{\mathcal{F}}$)

# Today: Hybrid Setting

Algorithm: Maximum Entropy Inverse Reinforcement Learning

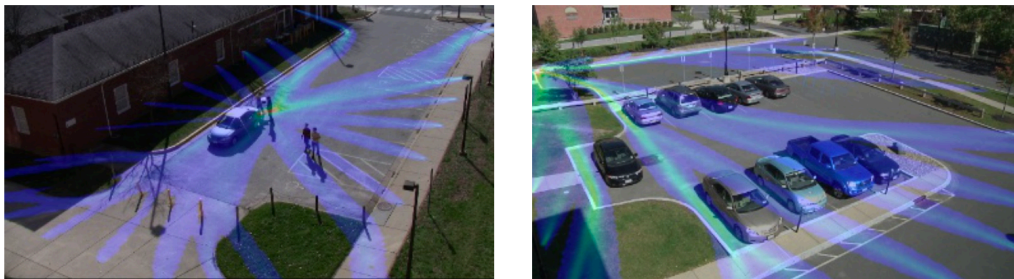# Running Example: Human trajectory forecasting

[Kitani, et al, ECCV 12]



**Fig. 1.** Given a single pedestrian detection, our proposed approach forecasts plausible paths and destinations from noisy vision-input

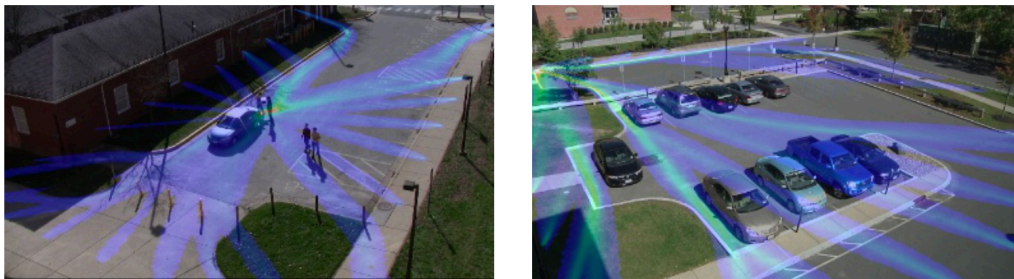# Running Example: Human trajectory forecasting

[Kitani, et al, ECCV 12]



**Fig. 1.** Given a single pedestrian detection, our proposed approach forecasts plausible paths and destinations from noisy vision-input

High-level assumptions:
(1) Experts may have some cost function regarding walking in their mind
(2) Experts are (approximately) optimizing the cost function

# Setting

Finite horizon MDP $\mathcal{M} = \{S, A, H, c, P, \mu_0, \pi^\star\}$

# Setting

Finite horizon MDP $\mathcal{M} = \{S, A, H, c, P, \mu_0, \pi^\star\}$

(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 $\mathcal{M} = \{S, A, H, c, P, \mu_0, \pi^\star\}$

(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       hybrid setting

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

# Setting

Finite horizon MDP $\mathcal{M} = \{S, A, H, c, P, \mu_0, \pi^\star\}$

(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** $\phi(s, a)$ $\in \mathbb{R}^d$

# Notation on Distributions

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

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

$\tau = \{s_0\, a_0\, s_1\, a_1 \cdots s_{H-1}, a_{H-1}, s_H\}$

$$\rho^\pi(\tau) := \mu_0(s_0)\pi(a_0 \,|\, s_0)P(s_1 \,|\, s_0, a_0)\pi(a_1 \,|\, s_1)\ldots\pi(a_{H-1} \,|\, s_{H-1})P(s_H \,|\, s_{H-1}, a_{H-1}):$$

Likelihood of the trajectory $\tau$ under $\pi$

# Running Example: Define feature map

**Key Assumption on cost:**
$$c(s, a) = \langle \theta^\star, \phi(s, a) \rangle\text{, linear wrt feature } \phi(s, a)$$



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

# Running Example: Define feature map

**Key Assumption on cost:**

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

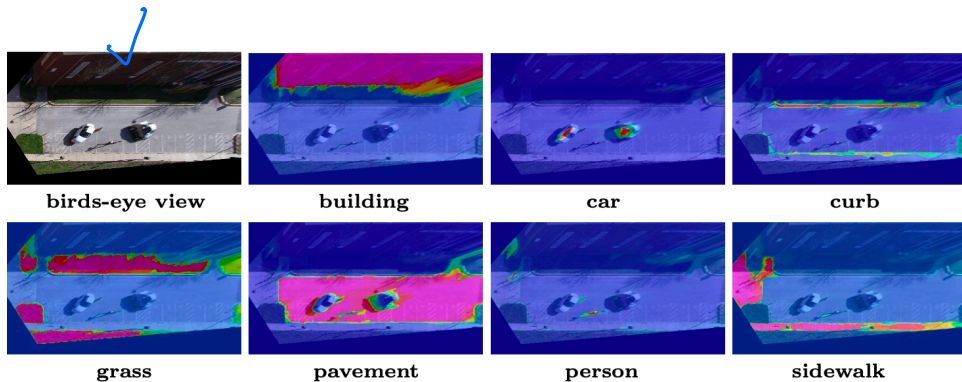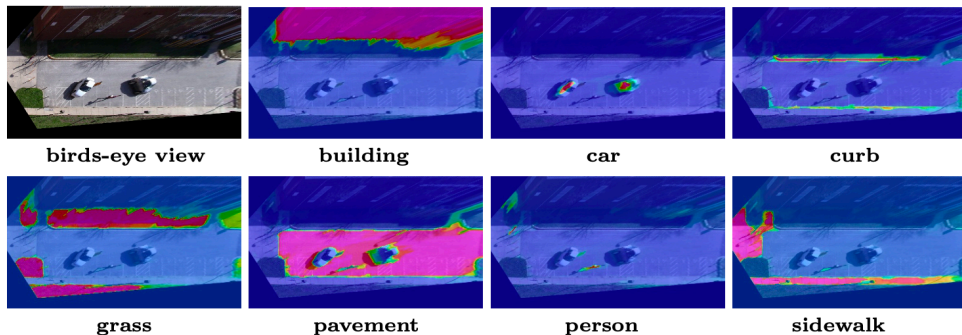State $s$: pixel or a group of neighboring pixels in image)



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

# Running Example: Define feature map

**Key Assumption on cost:**
$c(s, a) = \langle \theta^\star, \phi(s, a) \rangle$, **linear wrt feature** $\phi(s, a)$



birds-eye view     building     car     curb

grass     pavement     person     sidewalk

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

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

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

# Running Example: Define feature map

**Key Assumption on cost:**
$c(s, a) = \langle \theta^\star, \phi(s, a) \rangle$**, linear wrt feature** $\phi(s, a)$



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

birds-eye view    building    car    curb

grass    pavement    person    sidewalk

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

$$\phi(s, a) = \begin{bmatrix} \mathbb{P}(\text{pixels being building}) \\ \mathbb{P}(\text{pixels being grass}) \\ \mathbb{P}(\text{pixels being sidewalk}) \\ \mathbb{P}(\text{pixels being car}) \\ \dots \end{bmatrix}$$
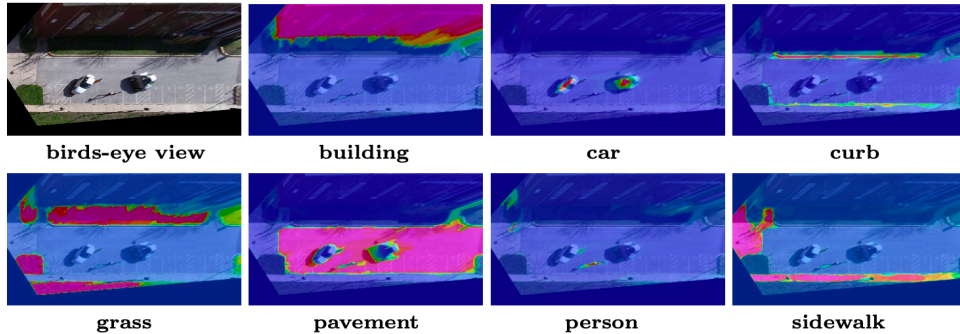
Maybe colliding with cars or buildings has **high** cost, but walking on sideway or grass has **low** cost

# 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

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

$Q^* ??$

# 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

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

Solution: $Q^\star = \mathcal{N}(\mu, \Sigma)$
(proof: use Lagrange multiplier)

$\exp\left(\text{quadratic form}\right)$

# Maximum Entropy Inverse RL:

# Maximum Entropy Inverse RL:

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

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

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

But there are potentially many such policies…

$$\theta^\star \cdot \phi(s,a)$$

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

# Maximum Entropy Inverse RL:

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

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

$\leftarrow$ *Traj — distribution*

Find a $\pi$ whose $\rho^\pi$ has the largest entropy,
subject to expected feature matching
$$\mathbb{E}_{s,a \sim d^\pi} \phi(s,a) = \mathbb{E}_{s,a \sim d^{\pi^\star}} \phi(s,a)$$

## Maximum Entropy Inverse RL:

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

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

Find a $\pi$ whose $\rho^\pi$ has the largest entropy,
subject to expected feature matching
$$\mathbb{E}_{s,a \sim d^\pi} \phi(s,a) = \mathbb{E}_{s,a \sim d^{\pi^\star}} \phi(s,a)$$

$$\max_{\pi} \text{entropy}[\rho^\pi]$$

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

# Maximum Entropy Inverse RL:

Let's simplify the objective $\max_{\pi} \text{entropy}[\rho^{\pi}]$:

# Maximum Entropy Inverse RL:

Let's simplify the objective $\max_{\pi}$ entropy$[\rho^{\pi}]$:

Recall the definition of trajectory distribution:
$$\rho^{\pi}(\tau) = \mu(s_0)\pi(a_0 \,|\, s_0)P(s_1 \,|\, s_0, a_0)\pi(a_1 \,|\, s_1)\ldots$$

# Maximum Entropy Inverse RL:

Let's simplify the objective $\max_{\pi}$ entropy$[\rho^\pi]$:

Recall the definition of trajectory distribution:

$$\rho^\pi(\tau) = \mu(s_0)\pi(a_0 \mid s_0)P(s_1 \mid s_0, a_0)\pi(a_1 \mid s_1)\ldots$$

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

# Maximum Entropy Inverse RL:

Let's simplify the objective $\max_{\pi} \text{entropy}[\rho^\pi]$:

Recall the definition of trajectory distribution:
$$\rho^\pi(\tau) = \mu(s_0)\pi(a_0 \mid s_0)P(s_1 \mid s_0, a_0)\pi(a_1 \mid s_1)\ldots$$

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

$$\arg\max_{\pi} \text{entroy}(\rho^\pi) = \arg\max_{\pi} -\sum_\tau \rho^\pi(\tau)\left[\sum_{h=0}^{H-1} \ln \pi(a_h \mid s_h)\right]$$

# Maximum Entropy Inverse RL:

Let's simplify the objective $\max_\pi \text{entropy}[\rho^\pi]$:

Recall the definition of trajectory distribution:
$$\rho^\pi(\tau) = \mu(s_0)\pi(a_0 \mid s_0)P(s_1 \mid s_0, a_0)\pi(a_1 \mid s_1)\ldots$$

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

$$\arg\max_\pi \text{entroy}(\rho^\pi) = \arg\max_\pi -\sum_\tau \rho^\pi(\tau)\left[\sum_{h=0}^{H-1} \ln \pi(a_h \mid s_h)\right]$$

$$= \arg\max_\pi -\sum_{h=0}^{H-1} \mathbb{E}_{s,a\sim \overset{P_h^\pi}{\Delta}} \ln \pi(a \mid s)$$

# Maximum Entropy Inverse RL:

Reformulating the optimization program:

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

$$\mathbb{E}_{s \sim d^{\pi}} \, C(s,a)$$

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

← feature-matching

# Maximum Entropy Inverse RL:

Reformulating the optimization program:

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

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

Using Lagrange formulation (Lagrange multiplier $\theta$), we get:

$$\min_{\pi} \mathbb{E}_{s,a \sim d^\pi} \ln \pi(a \mid s) + \max_{\theta} \left( \mathbb{E}_{s,a \sim d^\pi} \theta^\top \phi(s, a) - \mathbb{E}_{s,a \sim d^{\pi\star}} \theta^\top \phi(s, a) \right)$$

# Maximum Entropy Inverse RL:

Reformulating the optimization program:

$$\min_{\pi} \mathbb{E}_{s,a \sim d^{\pi}} \ln \pi(a \,|\, s)$$

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

$$\pi(\cdot \,|\, s) \in \Delta(A)$$

Using Lagrange formulation (Lagrange multiplier $\theta$), we get:

$$\min_{\pi} \mathbb{E}_{s,a \sim d^{\pi}} \ln \pi(a \,|\, s) + \max_{\theta} \left( \mathbb{E}_{s,a \sim d^{\pi}} \theta^{\top} \phi(s,a) - \mathbb{E}_{s,a \sim d^{\pi\star}} \theta^{\top} \phi(s,a) \right)$$

$$\min_{\pi} \mathbb{E}_{s \sim d^{\pi}} \left[ \text{Entropy}(\pi(\cdot \,|\, s)) + c(s,a) \right]$$

Using minimax theorem (John von Neumann), we can swap the order of min-max:

$$\max_{\theta} \min_{\pi} \left[ \mathbb{E}_{s,a \sim d^{\pi}} \ln \pi(a \,|\, s) + \mathbb{E}_{s,a \sim d^{\pi}} \theta^{\top} \phi(s,a) - \mathbb{E}_{s,a \sim d^{\pi\star}} \theta^{\top} \phi(s,a) \right] \checkmark$$

# Maximum Entropy Inverse RL: Final Algorithm

We get the final formulation:

$$\max_{\theta} \min_{\pi} \left[ \mathbb{E}_{s,a \sim d^\pi} \theta^\top \phi(s,a) - \mathbb{E}_{s,a \sim d^{\pi^\star}} \theta^\top \phi(s,a) + \mathbb{E}_{s,a \sim d^\pi} \ln \pi(a \mid s) \right]$$

$$\underbrace{\phantom{\mathbb{E}_{s,a \sim d^\pi} \ln \pi(a \mid s)}}_{\text{Regularization}}$$

fixed $\theta$

$\min_{\pi} \mathbb{E}_{s,a \sim d^\pi} \theta^\top \phi(s,a)$      $\mathbb{E}_{s,a \sim d^{\pi^\star}} \theta^\top \phi(s,a)$

$\textcircled{1}$          $\textcircled{2}$

$\textcircled{1} < \textcircled{2}$

for $\theta^\star$      $\textcircled{2} \subseteq \textcircled{1}$   ($\pi^\star$ is optimal under $\theta^\star$)

# Maximum Entropy Inverse RL: Final Algorithm

*incremental apdee on θ*

*Best Response on π*

We get the final formulation:

$$\max_\theta \min_\pi \left[ \mathbb{E}_{s,a\sim d^\pi}\theta^\top\phi(s,a) - \mathbb{E}_{s,a\sim d^{\pi^\star}}\theta^\top\phi(s,a) + \mathbb{E}_{s,a\sim d^\pi}\ln\pi(a\,|\,s) \right]$$

Algorithm: gradient ascent on $\theta$ (w/ fixed $\pi$),
and exact computation (e.g, planning, VI) for $\pi$ (w/ fixed $\theta$

# Maximum Entropy Inverse RL: Final Algorithm

We get the final formulation:

$$\max_{\theta} \min_{\pi} \left[ \mathbb{E}_{s,a \sim d^{\pi}} \theta^{\top} \phi(s,a) - \mathbb{E}_{s,a \sim d^{\pi^{\star}}} \theta^{\top} \phi(s,a) + \mathbb{E}_{s,a \sim d^{\pi}} \ln \pi(a \mid s) \right]$$

Algorithm: gradient ascent on $\theta$ (w/ fixed $\pi$),
and exact computation (e.g, planning, VI) for $\pi$ (w/ fixed $\theta$

$$\max_{\theta} \min_{\pi}$$

Initializing $\theta_0$:

For $t = 0, \ldots,$

*Best Response to $\theta_t$*

$$\pi_t = \arg\underset{\pi}{\min} \max \mathbb{E}_{s,a \sim d^{\pi}} \left[ \theta_t^{\top} \phi(s,a) + \pi(a \mid s) \right]$$

$$\theta_{t+1} = \theta_t + \eta \left( \mathbb{E}_{s,a \sim d^{\pi_t}} \phi(s,a) - \mathbb{E}_{s,a \sim d^{\pi^{\star}}} \phi(s,a) \right)$$

*Gradient of $\theta$ wrt $\pi_t$*

Return $\theta_T$

# Maximum Entropy Inverse RL: Final Algorithm

We get the final formulation:

$$\max_{\theta} \min_{\pi} \left[ \mathbb{E}_{s,a \sim d^\pi} \theta^\top \phi(s,a) - \mathbb{E}_{s,a \sim d^{\pi^\star}} \theta^\top \phi(s,a) + \mathbb{E}_{s,a \sim d^\pi} \ln \pi(a \mid s) \right]$$

Algorithm: gradient ascent on $\theta$ (w/ fixed $\pi$),
and exact computation (e.g, planning, VI) for $\pi$ (w/ fixed $\theta$

*( hybrid setting. $P$ is known)*

Initializing $\theta_0$:

For $t = 0, \ldots,$

(Maximum Entropy RL)

$$\pi_t = \arg\max_{\pi} \mathbb{E}_{s,a \sim d^\pi} \left[ \theta_t^\top \phi(s,a) + \pi(a \mid s) \right]$$

$$\theta_{t+1} = \theta_t + \eta \left( \mathbb{E}_{s,a \sim d^{\pi_t}} \phi(s,a) - \mathbb{E}_{s,a \sim d^{\pi^\star}} \phi(s,a) \right)$$

Return $\theta_T$

# Maximum Entropy Inverse RL: Final Algorithm

We get the final formulation:

$$\max_{\theta} \min_{\pi} \left[ \mathbb{E}_{s,a\sim d^{\pi}} \theta^{\top}\phi(s,a) - \mathbb{E}_{s,a\sim d^{\pi\star}} \theta^{\top}\phi(s,a) + \mathbb{E}_{s,a\sim d^{\pi}} \ln \pi(a\,|\,s) \right]$$

Algorithm: gradient ascent on $\theta$ (w/ fixed $\pi$),
and exact computation (e.g, planning, VI) for $\pi$ (w/ fixed $\theta$

Initializing $\theta_0$:

For $t = 0,\ldots,$

*we have $\not{p}$ P (planning)*

(Maximum Entropy RL)

$$\pi_t = \arg\max_{\pi} \mathbb{E}_{s,a\sim d^{\pi}} \left[ \theta_t^{\top}\phi(s,a) + \pi(a\,|\,s) \right]$$

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

$$\theta_{t+1} = \theta_t + \eta \left( \mathbb{E}_{s,a\sim d^{\pi_t}}\phi(s,a) - \mathbb{E}_{s,a\sim d^{\pi\star}}\phi(s,a) \right)$$

(Gradient equal to the difference of expected features)

Return $\theta_T$

# Maximum Entropy RL: Soft ~~Policy~~ *Value* Iteration

Maximum Entropy RL: what we do when our "cost" depends on policy $\pi$?

$$\arg \max_{\pi} \mathbb{E}_{s,a \sim d^{\pi}} \left[ c(s,a) + \pi(a \mid s) \right]$$

# Maximum Entropy RL: Soft Policy Iteration

Maximum Entropy RL: what we do when our "cost" depends on policy $\pi$?

$$\arg\max_{\pi} \mathbb{E}_{s,a \sim d^{\pi}} \left[ c(s,a) + \pi(a \mid s) \right]$$

*Value*

**Soft Policy Iteration:**

# Maximum Entropy RL: Soft Policy Iteration

Maximum Entropy RL: what we do when our "cost" depends on policy $\pi$?

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

$$\pi^{+}_{H-1}(s) = \arg\min_{a\in\pi_s} Q^{*}(s,a)$$

**Soft Policy Iteration:**

$$Q^{\star}_{H-1}(s,a) = \underbrace{c(s,a)}_{} \quad \boxed{\pi^{\star}_{H-1}(a\,|\,s) \propto \exp\left(-Q^{\star}_{H-1}(s,a)\right) \propto \exp(-A^{\star}_{H-1}(s,a))}$$

$$V^{\pi}_{H-1}(s) = \overset{\text{(*)}}{\mathbb{E}}_{a\sim\pi(\cdot|s)} \left[ c(s,a) + \underbrace{\ln\pi(a|s)}_{} \right]$$

$$\min_{\pi} V^{\pi}_{H-1}(s)$$

# Maximum Entropy RL: Soft Policy Iteration

Maximum Entropy RL: what we do when our "cost" depends on policy $\pi$?

$$\arg \max_{\pi} \mathbb{E}_{s,a \sim d^{\pi}} \left[ c(s,a) + \pi(a \mid s) \right]$$

## Soft Policy Iteration:

$$Q^{\star}_{H-1}(s,a) = c(s,a) \quad \pi^{\star}_{H-1}(a \mid s) \propto \exp \left( -Q^{\star}_{H-1}(s,a) \right) \propto \exp(-A^{\star}_{H-1}(s,a))$$

$$V^{\star}_{H-1}(s) = \mathbb{E}_{a \sim \pi^{\star}_{H-1}(a \mid s)} \left[ \ln \pi^{\star}_{H-1}(a \mid s) + Q^{\star}_{H-1}(s,a) \right] = -\ln \left( \sum_{a} \exp \left( -Q^{\star}_{H-1}(s,a) \right) \right)$$

$$V^{*}_{H-1}(s) = \min_{a \in A} Q^{*}_{H-1}(s,a)$$

# Maximum Entropy RL: Soft Policy Iteration

Maximum Entropy RL: what we do when our "cost" depends on policy $\pi$?

$$\arg \max_{\pi} \mathbb{E}_{s,a \sim d^{\pi}} \left[ c(s,a) + \pi(a \mid s) \right]$$

**Soft Policy Iteration:**

$$Q_{H-1}^{\star}(s,a) = c(s,a) \quad \pi_{H-1}^{\star}(a \mid s) \propto \exp\left(-Q_{H-1}^{\star}(s,a)\right) \propto \exp(-A_{H-1}^{\star}(s,a))$$

$$V_{H-1}^{\star}(s) = \mathbb{E}_{a \sim \pi_{H-1}^{\star}(a \mid s)} \left[\ln \pi_{H-1}^{\star}(a \mid s) + Q_{H-1}^{\star}(s,a)\right] = -\ln\left(\sum_{a} \exp\left(-Q_{H-1}^{\star}(s,a)\right)\right)$$

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

# Maximum Entropy RL: Soft Policy Iteration

Maximum Entropy RL: what we do when our "cost" depends on policy $\pi$?

$$\arg \max_{\pi} \mathbb{E}_{s,a \sim d^{\pi}} \left[ c(s,a) + \pi(a \mid s) \right]$$

## Soft Policy Iteration:

$$Q_{H-1}^{\star}(s,a) = c(s,a) \quad \pi_{H-1}^{\star}(a \mid s) \propto \exp \left( -Q_{H-1}^{\star}(s,a) \right) \propto \exp(-A_{H-1}^{\star}(s,a))$$

$$V_{H-1}^{\star}(s) = \mathbb{E}_{a \sim \pi_{H-1}^{\star}(a \mid s)} \left[ \ln \pi_{H-1}^{\star}(a \mid s) + Q_{H-1}^{\star}(s,a) \right] = - \ln \left( \sum_{a} \exp \left( -Q_{H-1}^{\star}(s,a) \right) \right)$$

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

$$\pi_{h}^{\star}(a \mid s) \propto \exp(-Q_{h}^{\star}(s,a)) \propto \exp(-A_{h}^{\star}(s,a))$$

# Maximum Entropy RL: Soft Policy Iteration

Maximum Entropy RL: what we do when our "cost" depends on policy $\pi$?

$$\arg \max_{\pi} \mathbb{E}_{s,a \sim d^{\pi}} \left[ c(s,a) + \pi(a \mid s) \right]$$

## Soft Policy Iteration:

$$Q^{\star}_{H-1}(s,a) = c(s,a) \quad \pi^{\star}_{H-1}(a \mid s) \propto \exp\left(-Q^{\star}_{H-1}(s,a)\right) \propto \exp(-A^{\star}_{H-1}(s,a))$$

$$V^{\star}_{H-1}(s) = \mathbb{E}_{a \sim \pi^{\star}_{H-1}(a \mid s)} \left[ \ln \pi^{\star}_{H-1}(a \mid s) + Q^{\star}_{H-1}(s,a) \right] = -\ln\left( \sum_{a} \exp\left(-Q^{\star}_{H-1}(s,a)\right) \right)$$

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

$$\pi^{\star}_{h}(a \mid s) \propto \exp(-Q^{\star}_{h}(s,a)) \propto \exp(-A^{\star}_{h}(s,a))$$

$$V^{\star}_{h}(s) = -\ln\left( \sum_{a} \exp(-Q^{\star}_{h}(s,a)) \right)$$

# Maximum Entropy RL: Soft Policy Iteration

Maximum Entropy RL: what we do when our "cost" depends on policy $\pi$?

$$\arg \max_{\pi} \mathbb{E}_{s,a \sim d^{\pi}} \left[ c(s,a) + \pi(a \mid s) \right]$$

**Soft Policy Iteration:**

$$Q_{H-1}^{\star}(s,a) = c(s,a) \quad \pi_{H-1}^{\star}(a \mid s) \propto \exp\left(-Q_{H-1}^{\star}(s,a)\right) \propto \exp(-A_{H-1}^{\star}(s,a))$$

$$V_{H-1}^{\star}(s) = \mathbb{E}_{a \sim \pi_{H-1}^{\star}(a \mid s)} \left[ \ln \pi_{H-1}^{\star}(a \mid s) + Q_{H-1}^{\star}(s,a) \right] = -\ln\left( \sum_a \exp\left(-Q_{H-1}^{\star}(s,a)\right) \right)$$

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

$$\pi_h^{\star}(a \mid s) \propto \exp(-Q_h^{\star}(s,a)) \propto \exp(-A_h^{\star}(s,a))$$

$$V_h^{\star}(s) = -\ln\left( \sum_a \exp(-Q_h^{\star}(s,a)) \right)$$

Derivation: DP!

## Maximum Entropy RL: Calculate Trajectory Likelihood

Initializing $\theta_0$:

For $t = 0, \ldots,$

$$\pi_t = \arg\max_{\pi} \mathbb{E}_{s,a \sim d^{\pi}} \left[ \theta_t^{\top} \phi(s,a) + \pi(a \mid s) \right]$$

$$\theta_{t+1} = \theta_t + \eta \left( \mathbb{E}_{s,a \sim d^{\pi_t}} \phi(s,a) - \mathbb{E}_{s,a \sim d^{\pi^\star}} \phi(s,a) \right)$$

Return $\theta_T$

$$\hat{\pi} := \text{soft } \forall \! \left( \theta_T^{\top} \phi(s,a) \right)$$

$$\hat{\pi}_h(a \mid s) \propto \exp\left( -Q_h^{\star}(s,a;\theta_T) \right)$$

$$c(s,a) = \theta_T^{\top} \cdot \phi(s,a)$$

## Maximum Entropy RL: Calculate Trajectory Likelihood

Initializing $\theta_0$:

For $t = 0, \ldots,$

$\quad \pi_t = \arg\max_{\pi} \mathbb{E}_{s,a \sim d^\pi} \left[ \theta_t^\top \phi(s,a) + \pi(a \mid s) \right]$

$\quad \theta_{t+1} = \theta_t + \eta \left( \mathbb{E}_{s,a \sim d^{\pi_t}} \phi(s,a) - \mathbb{E}_{s,a \sim d^{\pi^\star}} \phi(s,a) \right)$

Return $\theta_T$

$\widehat{\pi} := \text{soft PI} \left( \theta_T^\top \phi(s,a) \right)$

$\widehat{\pi}_h(a \mid s) \propto \exp\left( -Q_h^\star \left( s, a; \theta_T \right) \right)$

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 RL: Calculate Trajectory Likelihood

Initializing $\theta_0$:

For $t = 0, \ldots,$

$\qquad \pi_t = \arg\max_{\pi} \mathbb{E}_{s,a \sim d^{\pi}} \left[ \theta_t^{\top} \phi(s, a) + \pi(a \mid s) \right]$

$\qquad \theta_{t+1} = \theta_t + \eta \left( \mathbb{E}_{s,a \sim d^{\pi_t}} \phi(s, a) - \mathbb{E}_{s,a \sim d^{\pi^\star}} \phi(s, a) \right)$

Return $\theta_T$

$\qquad \hat{\pi} := \text{soft PI} \left( \theta_T^{\top} \phi(s, a) \right)$

$\qquad \hat{\pi}_h(a \mid s) \propto \exp \left( -Q_h^{\star}(s, a; \theta_T) \right)$

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?

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

# Maximum Entropy RL: Calculate Trajectory Likelihood

Initializing $\theta_0$:

For $t = 0, \ldots,$

$\qquad \pi_t = \arg\max_{\pi} \mathbb{E}_{s,a \sim d^{\pi}} \left[ \theta_t^{\top} \phi(s,a) + \pi(a \mid s) \right]$

$\qquad \theta_{t+1} = \theta_t + \eta \left( \mathbb{E}_{s,a \sim d^{\pi_t}} \phi(s,a) - \mathbb{E}_{s,a \sim d^{\pi^{\star}}} \phi(s,a) \right)$

Return $\theta_T$

$\hat{\pi} := \text{softVI} \left( \theta_T^{\top} \phi(s,a) \right)$

$\hat{\pi}_h(a \mid s) \propto \exp \left( -Q_h^{\star}(s,a; \theta_T) \right)$

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?

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

Special case: deterministic MDP and state-dependent cost:

$\min_{\pi} \ \mathbb{E}_{s \sim d^{\pi}} \left[ c(s,a) \right] \quad \Leftarrow NPG$

$\min_{\pi} \ \mathbb{E}_{s \sim d^{\pi}} \left[ c(s,a) \right] + \frac{1}{\beta} \left( -\text{entropy}(\rho^{\pi}) \right)$

# Maximum Entropy RL: Calculate Trajectory Likelihood

Initializing $\theta_0$:

For $t = 0, \ldots,$

$\quad \pi_t = \arg\max\limits_{\pi} \mathbb{E}_{s,a \sim d^\pi} \left[ \theta_t^\top \phi(s,a) + \pi(a \mid s) \right]$

$\quad \theta_{t+1} = \theta_t + \eta \left( \mathbb{E}_{s,a \sim d^{\pi_t}} \phi(s,a) - \mathbb{E}_{s,a \sim d^{\pi^\star}} \phi(s,a) \right)$

Return $\theta_T$

$\quad \hat{\pi} := \text{soft PI} \left( \theta_T^\top \phi(s,a) \right)$

$\quad \hat{\pi}_h(a \mid s) \propto \exp\left( -Q_h^\star(s, a; \theta_T) \right)$

$\qquad\qquad \theta_T \cdot \phi(s)$

Given a trajectory $\tau = \{ s_0, a_0, \ldots, s_{H-1}, a_{H-1} \}$
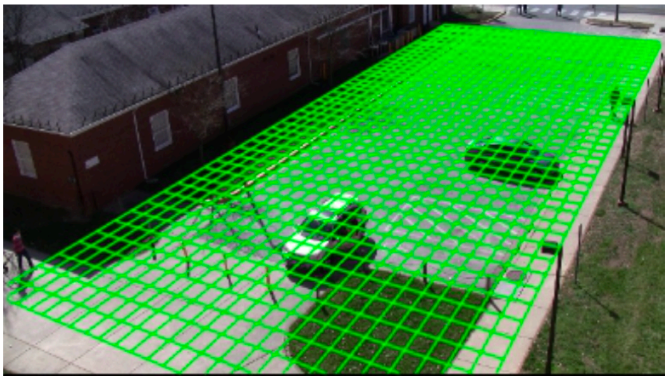
What's the likelihood of $\tau$ being generated by expert?

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

Special case: deterministic MDP and state-dependent cost: $\theta^* \phi(s)$

For a state trajectory, we have:

$$\rho^\pi(s_0, s_1, \ldots, s_H) \propto \exp\left( -\sum_h \theta_T^\top \phi(s_h) \right)$$

# Running Example: Human Trajectory Forecasting



State space: grid,
action space: 4 actions

We predict that we are more likely to use
sidewalk

$$\hat{C}(s) = \theta_J^T \phi(s)$$
$$p^\pi(z) \propto \exp\left(-\sum_{h=0}^{H} \hat{C}(s_h)\right)$$

# MaxEnt – IRL

$$\max_{\pi} \text{Entropy}(\rho^{\pi})$$

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

---

Soft VI:

(SAC)

$$V_h^*(s) = -\ln\left[ \sum_a \exp(-Q_h^*(s,a)) \right]$$

---

$$\pi^*(a|s) \propto \exp(-Q_h^*(s,a))$$

$$\left\| \mathop{\mathbb{E}}_{s,a \sim \pi} \phi(s,a) - \mathop{\mathbb{E}}_{s,a \sim \pi^*} \phi(s,a) \right\|_2^2$$

$$\|x\|_2 \iff \max_{\theta : \|\theta\|_2 \leq 1} \theta^T \cdot x$$

$$\max_{\theta \in \Theta} \left[ \mathop{\mathbb{E}}_{s,a \sim \pi} \theta^T \cdot \phi(s,a) - \mathop{\mathbb{E}}_{s,a \sim \pi^*} \theta^T \phi(s,a) \right]$$

$\uparrow$

unit-Ball

$$\mathcal{F} = \left\{ \theta^T \cdot \phi(s,a) ; \theta \in \text{Unit-Ball} \right\}$$

$p^\pi$

$d^\pi$

$\text{Entropy}(p^\pi) \leftarrow$ strongly

convex
function

Deterministic & $\phi(s)$

$$p_\theta(\tau) \propto \exp\left(-\sum_{h=0}^{H-1} \theta^T \cdot \phi(s)\right)$$

$\tau_1^*, \tau_2^* \dots \tau_N^* \leftarrow p^{\pi^*}$

$$\max_\theta \max \sum_{i=1}^{N} \ln\left(p_\theta(\tau_i^*)\right)$$

$\uparrow$

original - MaxEnt IRL