

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

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

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

# Recap

**Hybrid IL Setting:**



# Recap

## Hybrid IL Setting:

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

# Recap

## Hybrid IL Setting:

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

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

# Recap

## Hybrid IL Setting:

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

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

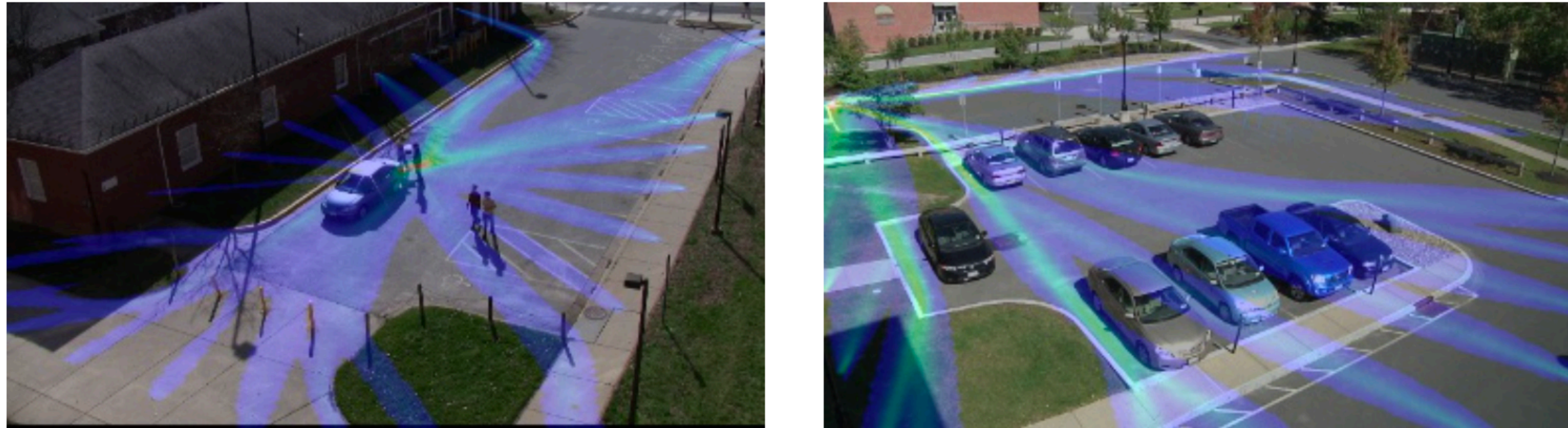
$$\hat{\pi} := \arg \min_{\pi \in \Pi} \left[ \max_{f \in \tilde{\mathcal{F}}} \left[ \mathbb{E}_{s, a \sim d^\pi} f(s, a) - \frac{1}{M} \sum_{i=1}^M f(s_i^\star, a_i^\star) \right] \right]$$

# **Today: Hybrid Setting**

Algorithm: Maximum Entropy Inverse Reinforcement Learning

# Running Example: Human trajectory forecasting

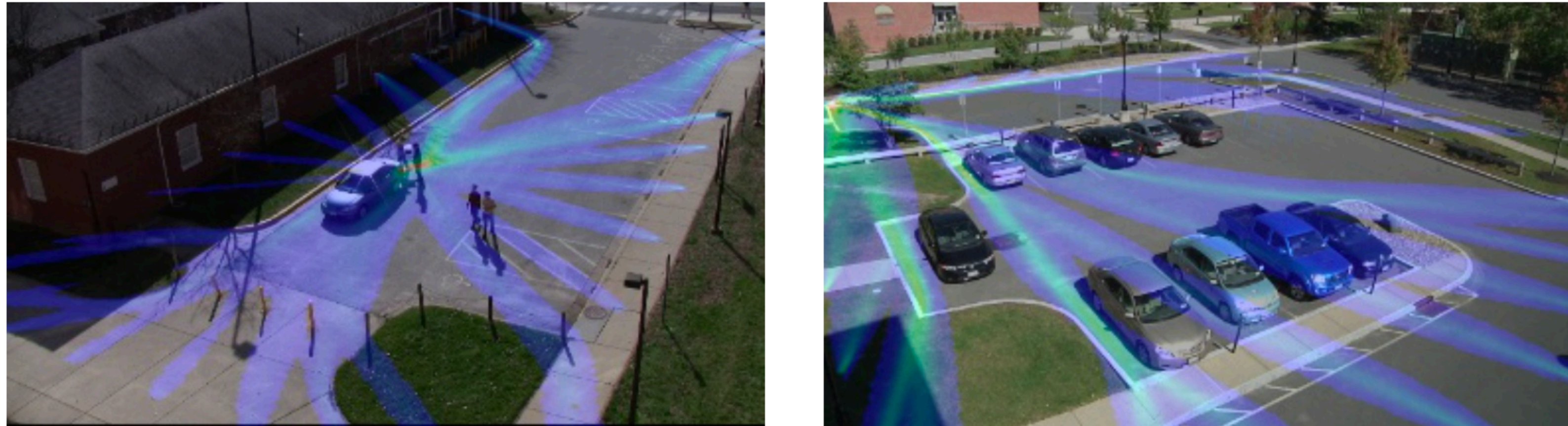
[Kitani, et al, ECCV 12]



**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

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  $\phi(s, a)$

## Notation on Distributions

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

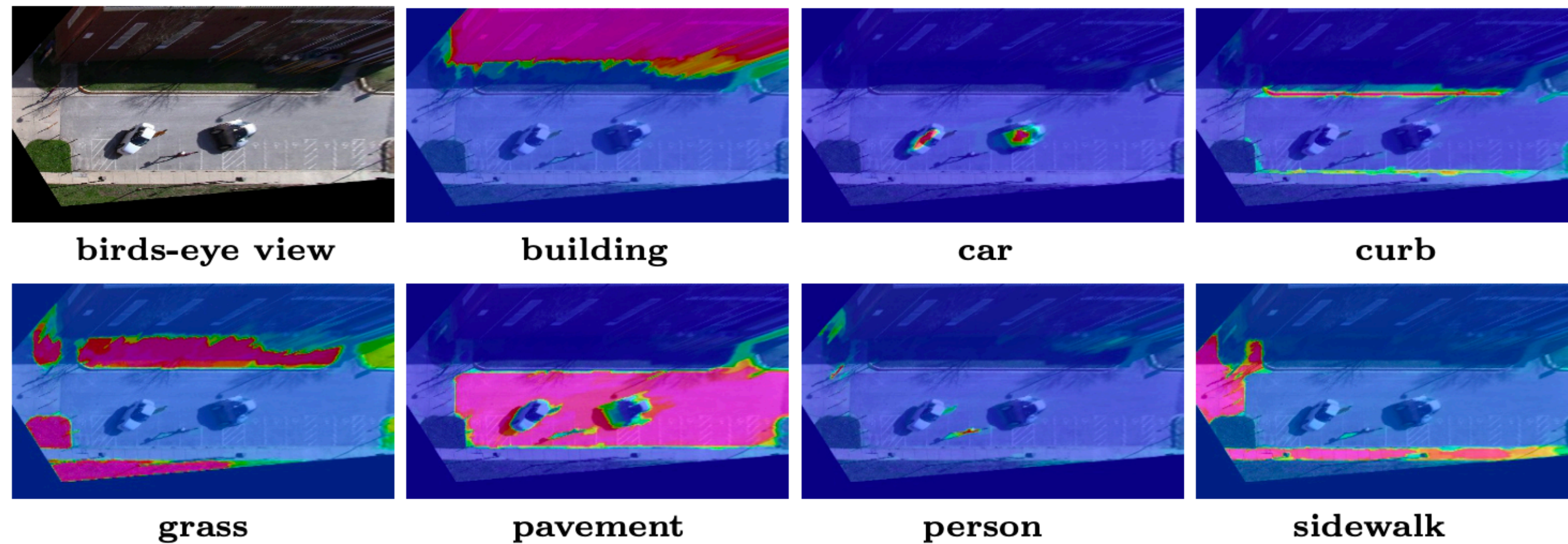
$$d^\pi(s, a) = \sum_{h=0}^{H-1} \mathbb{P}_h^\pi(s, a) / H: \text{average state-action visitation}$$

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

# Running Example: Define feature map

**Key Assumption on cost:**

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



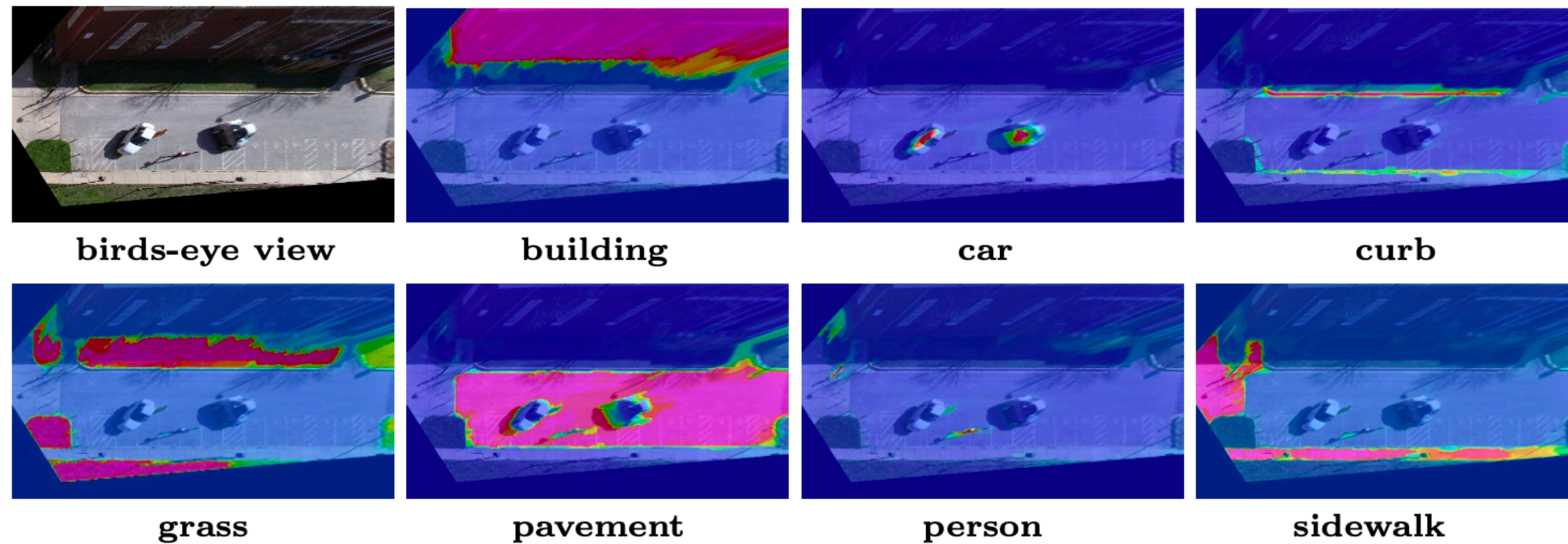
**Fig. 4.** Classifier feature response maps. Top left is the original image.

# 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.

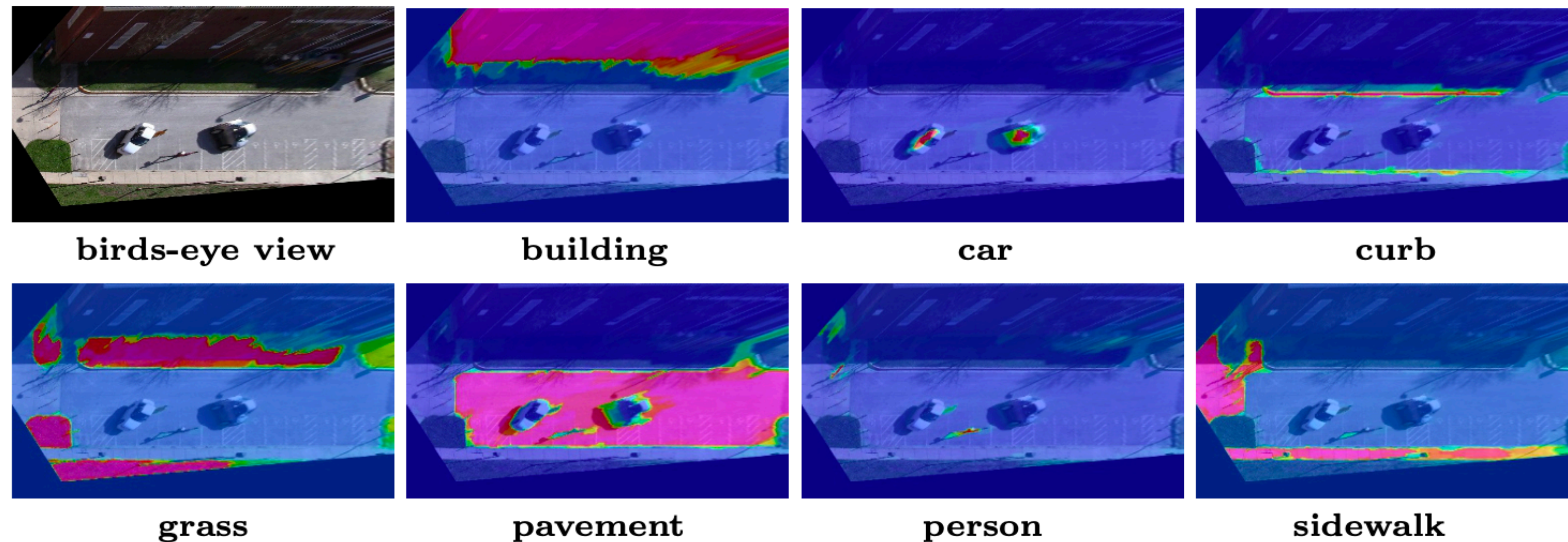
# Running Example: Define feature map

**Key Assumption on cost:**

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

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

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



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

# Running Example: Define feature map

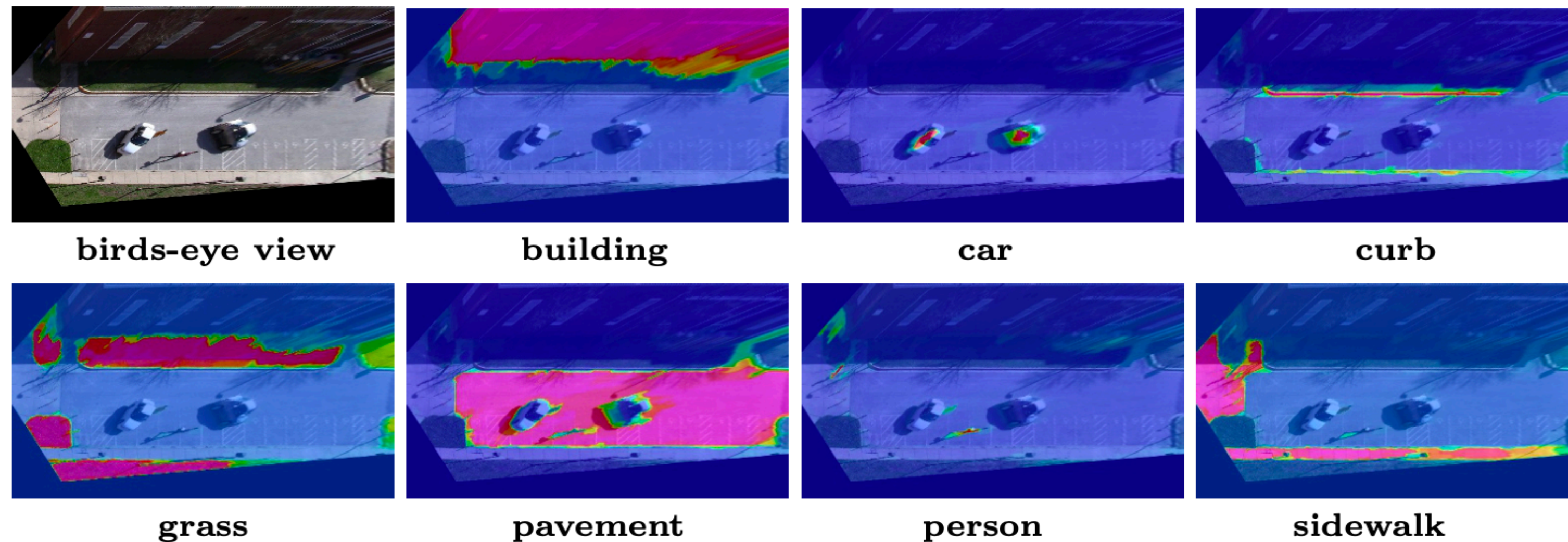
**Key Assumption on cost:**

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

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

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

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



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

## 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

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

# 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)

# Maximum Entropy Inverse RL:

# Maximum Entropy Inverse RL:

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

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

But there are potentially many such policies...

# Maximum Entropy Inverse RL:

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

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

But there are potentially many such policies...

Find a  $\pi$  whose  $\rho^\pi$  has the largest entropy,  
subject to expected feature matching

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

# Maximum Entropy Inverse RL:

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

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

But there are potentially many such policies...

Find a  $\pi$  whose  $\rho^\pi$  has the largest entropy,  
subject to expected feature matching

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

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

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

# Maximum Entropy Inverse RL:

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



# Maximum Entropy Inverse RL:

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

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

# Maximum Entropy Inverse RL:

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

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

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

# Maximum Entropy Inverse RL:

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

Recall the definition of trajectory distribution:

$$\rho^{\pi}(\tau) = \mu(s_0)\pi(a_0 | s_0)P(s_1 | s_0, a_0)\pi(a_1 | s_1)\dots$$

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

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

# Maximum Entropy Inverse RL:

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

Recall the definition of trajectory distribution:

$$\rho^{\pi}(\tau) = \mu(s_0)\pi(a_0 | s_0)P(s_1 | s_0, a_0)\pi(a_1 | s_1)\dots$$

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

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

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

# Maximum Entropy Inverse RL:

Reformulating the optimization program:

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

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

# Maximum Entropy Inverse RL:

Reformulating the optimization program:

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

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

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

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

# Maximum Entropy Inverse RL:

Reformulating the optimization program:

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

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

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

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

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

# Maximum Entropy Inverse RL: Final Algorithm

We get the final formulation:

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



# Maximum Entropy Inverse RL: Final Algorithm

We get the final formulation:

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

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:

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

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

Initializing  $\theta_0$ :

For  $t = 0, \dots,$

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

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

Return  $\theta_T$

# Maximum Entropy Inverse RL: Final Algorithm

We get the final formulation:

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

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

Initializing  $\theta_0$ :

For  $t = 0, \dots,$

(Maximum Entropy RL)

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

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

Return  $\theta_T$

# Maximum Entropy Inverse RL: Final Algorithm

We get the final formulation:

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

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

Initializing  $\theta_0$ :

For  $t = 0, \dots,$

(Maximum Entropy RL)

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

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

(Gradient equal to the difference of expected features)

Return  $\theta_T$

# Maximum Entropy RL: Soft Value Iteration

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

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

# Maximum Entropy RL: Soft Value Iteration

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

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

**Soft Value Iteration:**

# Maximum Entropy RL: Soft Value Iteration

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

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

**Soft Value Iteration:**

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

# Maximum Entropy RL: Soft Value Iteration

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

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

**Soft Value Iteration:**

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

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



# Maximum Entropy RL: Soft Value Iteration

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

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

**Soft Value Iteration:**

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

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

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

# Maximum Entropy RL: Soft Value Iteration

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

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

**Soft Value Iteration:**

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

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

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

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

# Maximum Entropy RL: Soft Value Iteration

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

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

**Soft Value Iteration:**

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

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

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

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

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

# Maximum Entropy RL: Soft Value Iteration

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

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

**Soft Value Iteration:**

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

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

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

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

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

Derivation: DP!

# Maximum Entropy RL: Calculate Trajectory Likelihood

Initializing  $\theta_0$ :

For  $t = 0, \dots,$

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

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

Return  $\theta_T$

$$\hat{\pi} := \text{soft VI} (\theta_T^{\top} \phi(s, a))$$

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

# Maximum Entropy RL: Calculate Trajectory Likelihood

Initializing  $\theta_0$ :

For  $t = 0, \dots,$

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

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

Return  $\theta_T$

$$\hat{\pi} := \text{soft VI} (\theta_T^{\top} \phi(s, a))$$

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

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

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

# Maximum Entropy RL: Calculate Trajectory Likelihood

Initializing  $\theta_0$ :

For  $t = 0, \dots,$

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

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

Return  $\theta_T$

$$\hat{\pi} := \text{soft VI} (\theta_T^{\top} \phi(s, a))$$

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

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

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

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

# Maximum Entropy RL: Calculate Trajectory Likelihood

Initializing  $\theta_0$ :

For  $t = 0, \dots,$

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

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

Return  $\theta_T$

$$\hat{\pi} := \text{soft VI} (\theta_T^{\top} \phi(s, a))$$

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

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

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

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

**Special case: deterministic MDP and state-dependent cost:**



# Maximum Entropy RL: Calculate Trajectory Likelihood

Initializing  $\theta_0$ :

For  $t = 0, \dots,$

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

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

Return  $\theta_T$

$$\hat{\pi} := \text{soft VI} (\theta_T^{\top} \phi(s, a))$$

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

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

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

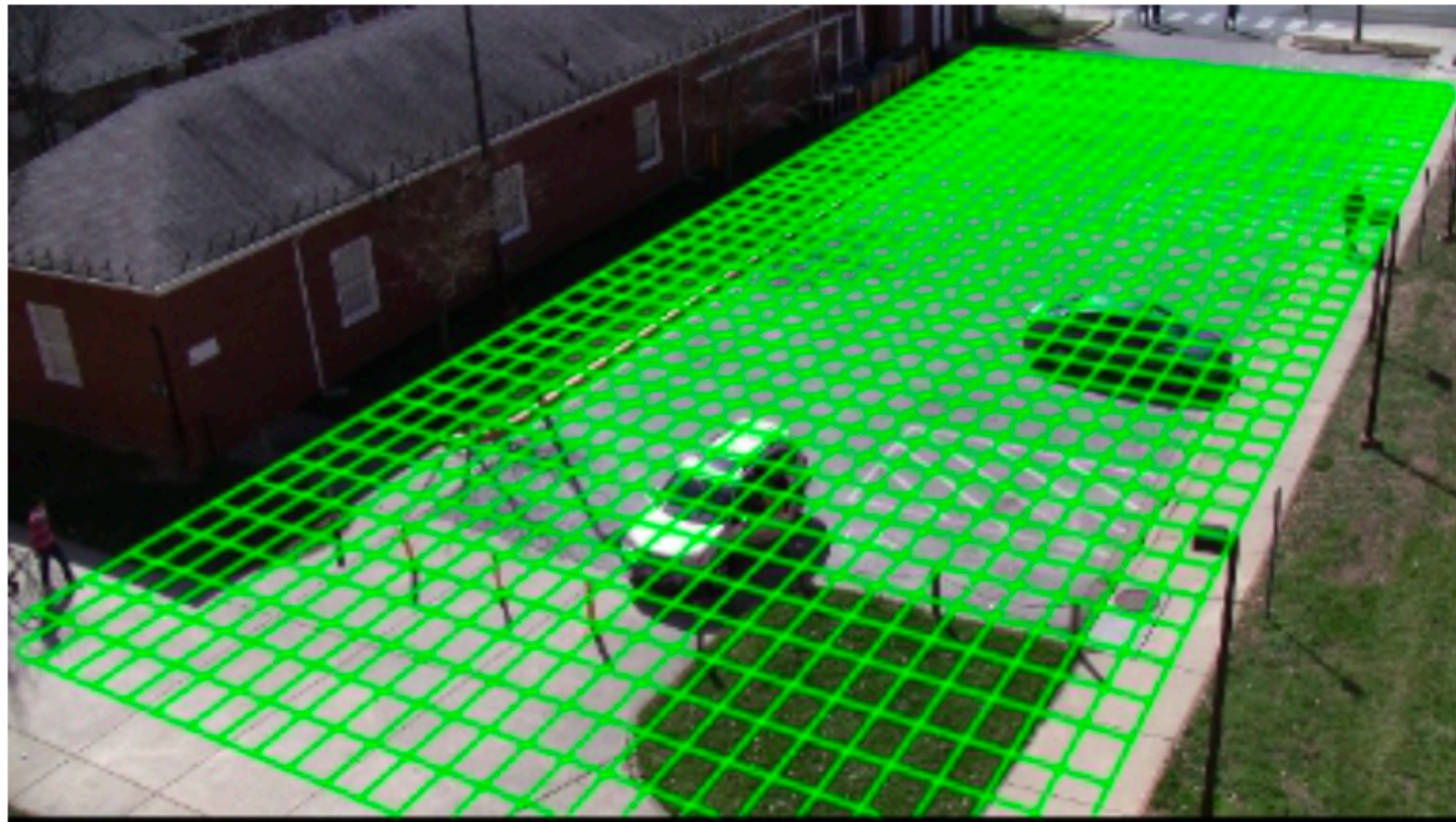
$$\ln \left( \rho^{\hat{\pi}}(\tau) \right) = \sum_{h=0}^{H-1} [\ln P(s_{h+1} | s_h, a_h) + \ln \hat{\pi}(a_h | s_h)]$$

**Special case: deterministic MDP and state-dependent cost:**

For a state trajectory, we have:

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

# Running Example: Human Trajectory Forecasting



State space: grid,  
action space: 4 actions



We predict that we are more likely to use  
sidewalk