

# CS 277: Control and Reinforcement Learning

Winter 2022

## Lecture 13: Inverse RL

**Roy Fox**

Department of Computer Science

Bren School of Information and Computer Sciences

University of California, Irvine



# Logistics

---

assignments

- Assignment 3 is due **today**
- Assignment 4 to be published soon

quizzes

- Quiz 5 is due **Friday**

# Today's lecture

---

**Belief-state MDPs**

**RNNs**

**IRL, Feature Matching**

**MaxEnt IRL**

# Tiger domain

- 2 states: which door leads to a tiger (-100 reward) and which to \$\$\$ (+10)
- You can stop and listen:  $p(o_t = s_t | s_t) = 0.8$



$$p(s_0 = s_{\text{left}}) = 0.5$$

$$\mathbb{E}[r(s_0, a_{\text{left}})] = -45$$

$$o_1 = o_{\text{right}}$$

$$p(s_1 = s_{\text{left}}) = 0.2$$

$$\mathbb{E}[r(s_1, a_{\text{left}})] = -12$$

$$o_2 = o_{\text{left}}$$

$$p(s_2 = s_{\text{left}}) = 0.5$$

$$\mathbb{E}[r(s_2, a_{\text{left}})] = -45$$

$$o_3 = o_{\text{right}}$$

$$p(s_3 = s_{\text{left}}) = 0.2$$

$$\mathbb{E}[r(s_3, a_{\text{left}})] = -12$$

$$o_4 = o_{\text{right}}$$

$$p(s_4 = s_{\text{left}}) = \frac{0.04}{0.04 + 0.64} \approx 0.06 \quad \mathbb{E}[r(s_4, a_{\text{left}})] = -3.5$$

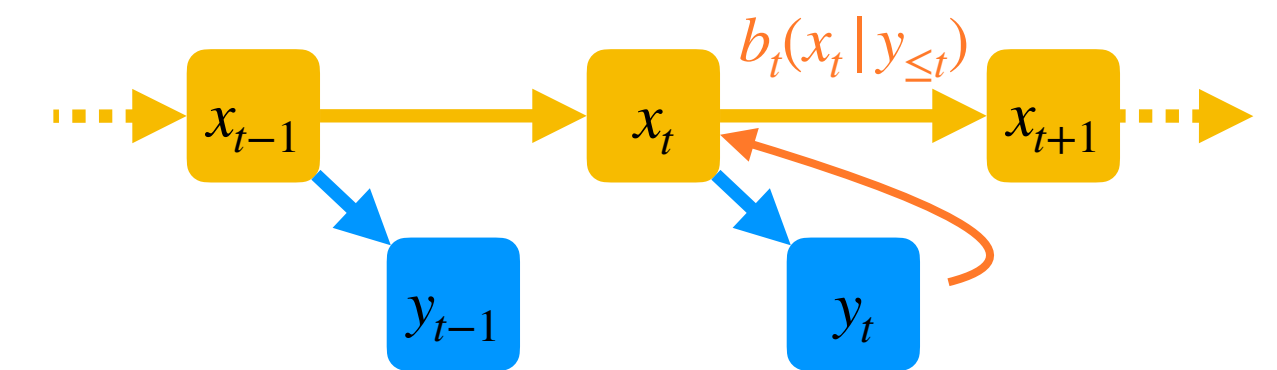
$$o_5 = o_{\text{right}}$$

$$p(s_5 = s_{\text{left}}) \approx 0.015$$

$$\mathbb{E}[r(s_4, a_{\text{left}})] = -8.3$$

# Belief

- **Belief** = distribution over the state  $b(s)$ 
  - If the agent reaches belief  $b$  after history  $h$ , that **does not imply**  $s \sim b$
- **Bayesian belief**  $b_h(s) = p(s | h)$ : a sufficient statistic of  $h$  for  $s$ 
  - For a Bayesian belief:  $s \sim b_h$  after history  $h$
- In the linear–Gaussian case: the **Kalman filter**
  - Bayesian belief is **Gaussian**  $p(x_t | h_t = y_{\leq t}) = \mathcal{N}(x_t; \hat{x}_t, \Sigma_t)$
  - Covariance can be **precomputed**  $\mathbb{V}(x_t | h_t) = \Sigma_t$  (independent of  $h_t$ )
  - Mean can be **updated linearly**:  $\hat{x}'_t = A\hat{x}_{t-1} + Bu_{t-1}$        $e_t = y_t - C\hat{x}'_t$        $\hat{x}_t = \hat{x}'_t + K_t e_t$



# Computing the Bayesian belief

- **Predict**  $s_{t+1}$  from  $h_t = (o_0, a_0, o_1, a_1, \dots, o_t)$  and  $a_t$ :

$$b'_t(s_{t+1} | h_t, a_t) = \sum_{s_t} p(s_t | h_t) p(s_{t+1} | s_t, a_t) = \sum_{s_t} b_t(s_t) p(s_{t+1} | s_t, a_t)$$

total probability over  $s_t$      
 previous belief  $b_t$      
 dynamics needs to be known

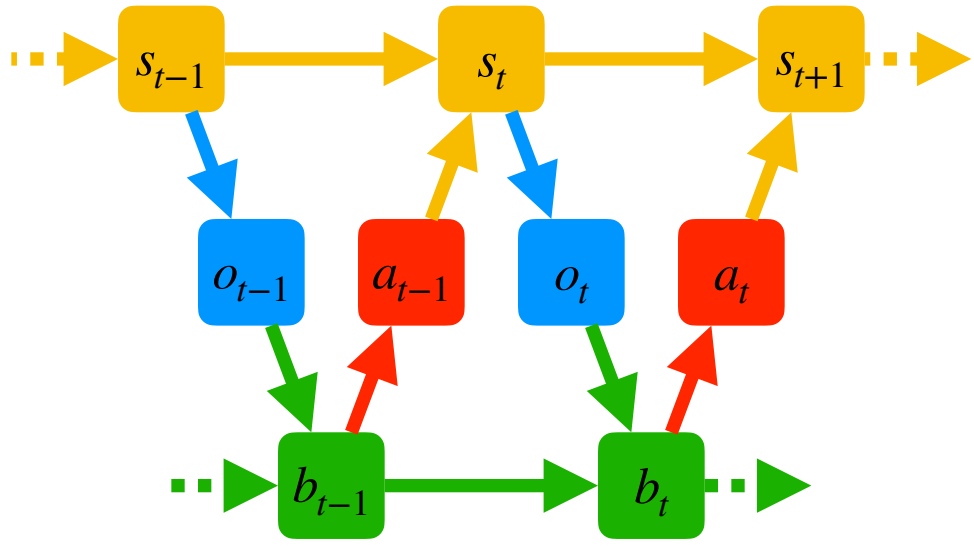
- **Update** belief of  $s_t$  after seeing  $h_t = (h_{t-1}, a_{t-1}, o_t)$ :

$$b_t(s_t | h_t) = \frac{p(s_t | h_{t-1}, a_{t-1}) p(o_t | s_t)}{p(o_t | h_{t-1}, a_{t-1})} = \frac{b'_{t-1}(s_t) p(o_t | s_t)}{\sum_{\bar{s}_t} b'_{t-1}(\bar{s}_t) p(o_t | \bar{s}_t)}$$

previous prediction     
 observation model  
Bayes' rule on  $o_t$      
  $o_t - s_t - (h_{t-1}, a_{t-1})$      
 normalizer

- A **deterministic, model-based** update:

▶  $b_{t-1}(s_{t-1}) \rightarrow$  use  $a_{t-1}$  to **predict**  $b'_{t-1}(s_t) \rightarrow$  use  $o_t$  to **update**  $b_t(s_t)$



# Belief-state MDP

- In the linear–quadratic–Gaussian case: **certainty equivalence**

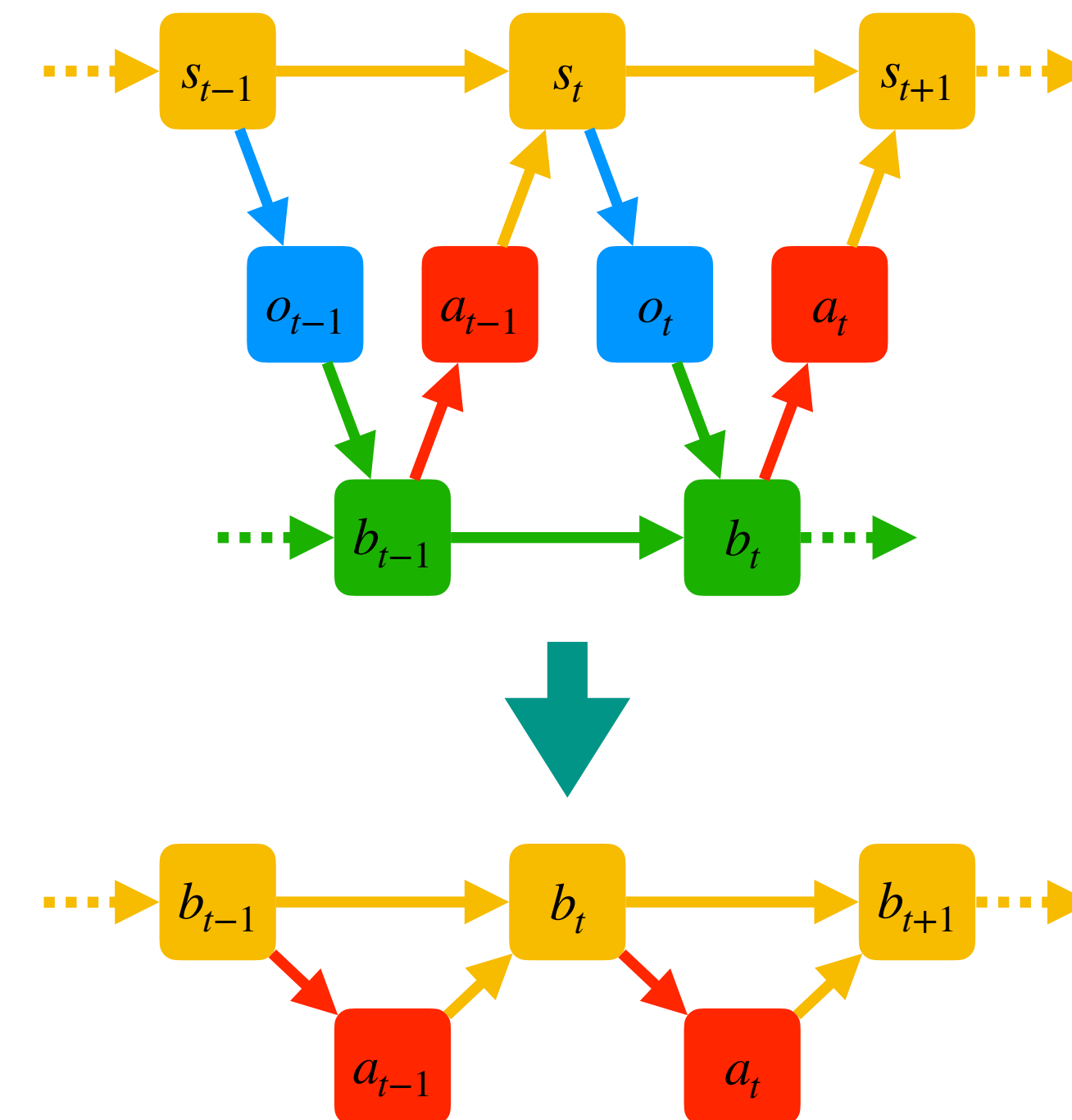
- Plan using  $\hat{x}_t$  as if it was  $x_t$

- More generally (though vastly less useful): **belief-state MDP**

- States:**  $\Delta(\mathcal{S})$    **Actions:**  $\mathcal{A}$    **Rewards:**  $r(b_t, a_t) = \sum_{s_t} b_t(s_t)r(s_t, a_t)$

- Transitions:** each possible observation  $o_{t+1}$  contributes its probability

$$p(o_{t+1} | b_t, a_t) = \sum_{s_t, s_{t+1}} b_t(s_t)p(s_{t+1} | s_t, a_t)p(o_{t+1} | s_{t+1})$$



to the total probability that the belief that follows  $(b_t, a_t, o_{t+1})$  is the **Bayesian belief**

$$b_{t+1}(s_{t+1}) = p(s_{t+1} | b_t, a_t, o_{t+1}) = \frac{\sum_{s_t} b_t(s_t)p(s_{t+1} | s_t, a_t)p(o_{t+1} | s_{t+1})}{p(o_{t+1} | b_t, a_t)}$$

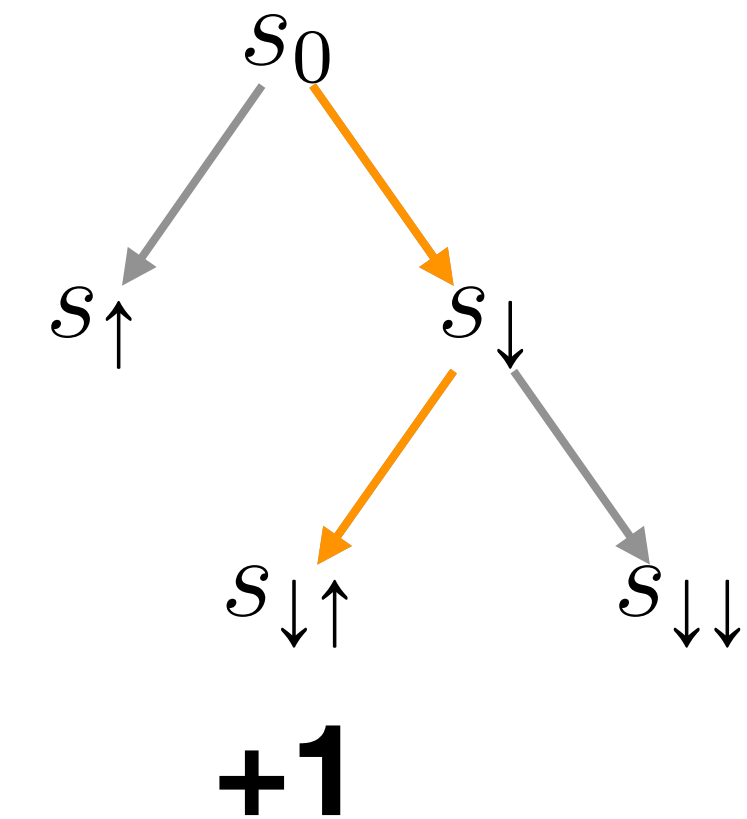
# Learning to use memory is hard

- Belief space  $b(s_t)$  is continuous and high-dimensional (dimension  $|\mathcal{S}|$ )
  - Curse of dimensionality
  - Beliefs are naturally multi-modal — how do we even represent them?
- The number of reachable beliefs may grow exponentially in  $t$  (one per  $h_t$ )
  - Curse of history
- Belief-value function can be very complex, hard to approximate
- There may not be optimal stationary deterministic policy  $\Rightarrow$  instability



# Stationary deterministic policy counterexample

- Assume **no observability**
- Stationary deterministic policies gets **no reward**
- **Non-stationary** policy:  $\downarrow, \uparrow$ ; expected return:  $+1$



▸ But non-stationary = observability of a clock  $t$

- Stationary **stochastic policy**:  $\downarrow / \uparrow$  with equal prob.; expected return:  $+0.25$

- Open problem: **Bellman optimality** is inherently stationary and deterministic

**no dependence on  $t$**  →  $V(s) = \max_a r(s, a) + \gamma \mathbb{E}_{(s'|s,a) \sim p} [V(s')]$  **maximum achieved for some action**

# Today's lecture

---

Belief-state MDPs

**RNNs**

IRL, Feature Matching

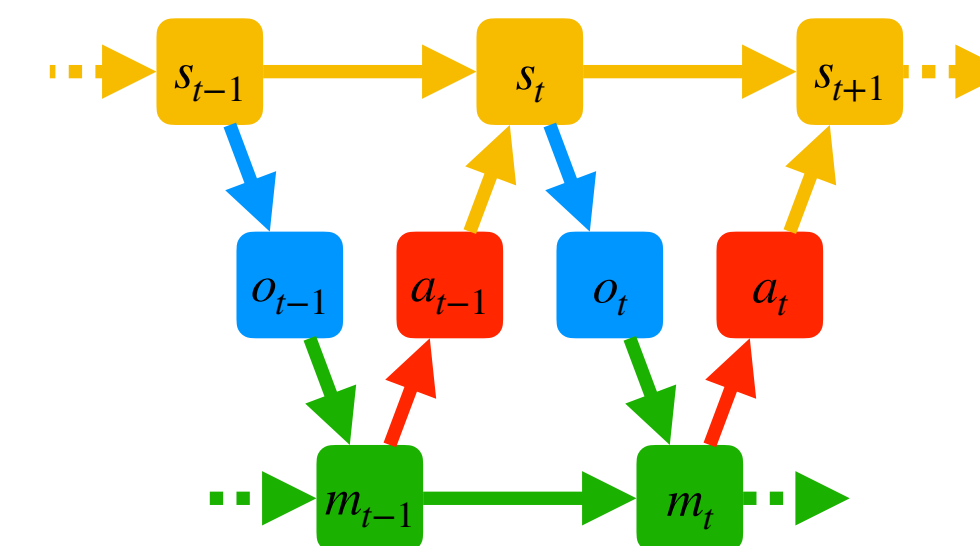
MaxEnt IRL

# Filtering with function approximation

- Instead of Bayesian belief: **memory update**  $m_t = f_\theta(m_{t-1}, o_t)$  ( $a_{t-1}$  optional)

- ▶ **Action policy**:  $\pi_\theta(a_t | m_t)$

- ▶ Sequential structure = **Recurrent Neural Network (RNN)**



- **Training**: back-propagate gradients through the whole sequence

- ▶ **Back-propagation through time (BPTT)**

- Unfortunately, gradients tend to **vanish**  $\rightarrow 0$  / **explode**  $\rightarrow \infty$

- ▶ **Long term coordination** of memory updates + actions is challenging

- ▶ RNN **can't use** information not remembered, but backup **no gradient** unless used

# RNNs in on-policy methods

- Training RNNs with **on-policy methods** is straightforward (and backward)

- ▶ **Roll out policy**: parameters of  $a_t$  distribution are determined by  $\pi_\theta(m_t)$  with

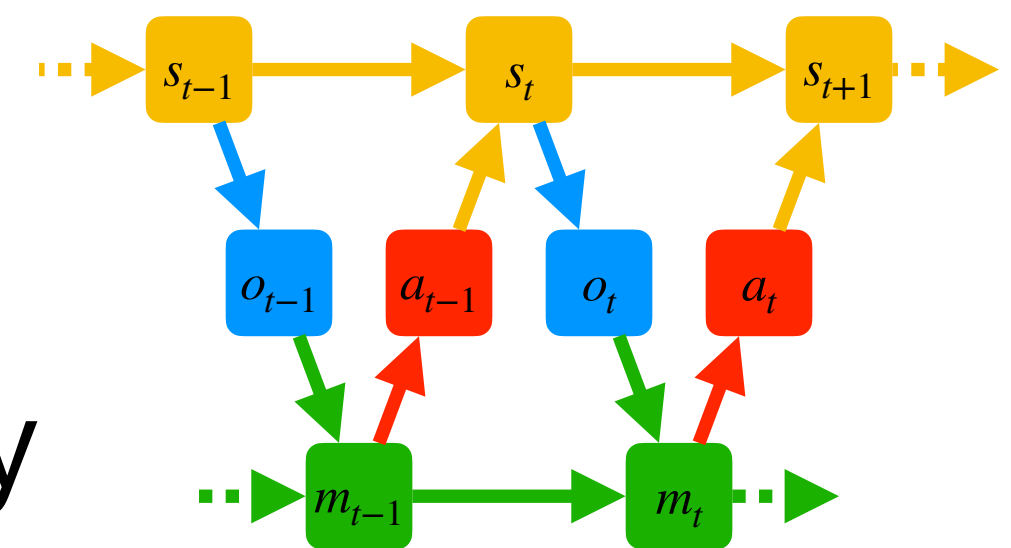
$$m_t = f_\theta(\dots f_\theta(f_\theta(o_0), o_1), \dots o_t)$$

- ▶ Compute  $\nabla_\theta \log \pi_\theta(a_t | m_t)$  with **BPTT** all the way to initial observation  $o_0$

- **Problems**: computation graph > **RAM**; **vanishing / exploding** grads

- ▶ **Solutions**: **stop gradients** every  $k$  steps; use **attention**

- **Problem**: cannot learn **longer memory** — but that's hard anyway



# RNNs in off-policy methods

- **Problem:** RNN states in replay buffer disagree with current RNN params
- **Solution 1:** use  $n$ -step rollouts to reduce mismatch effect

$$Q_{\theta}(o_t, m_t, a_t) \rightarrow r_t + \gamma r_{t+1} + \dots + \gamma^{n-1} r_{t+n-1} + \gamma^n \max_{a'} Q_{\theta}(o_{t+n}, m_{t+n}, a')$$

- **Solution 2:** “burn in”  $m_t$  from even earlier stored steps
  - Same target, but  $m_t$  is initialized from  $(o_{t-k}, \dots, o_{t-1})$
- **In practice:** RNNs rarely used
  - **Stacking  $k$  frames** every step  $(o_{t-k+1}, \dots, o_t)$  may help with short-term memory

# Deep RL as partial observability

---

- Memory-based policies fail us in **Deep RL**, where we need them most:
  - Deep RL is inherently **partially observable**
- Consider what **deeper layers** get as input:
  - High-level / action-relevant state features are **not Markov!**
- **Memory management** is a huge open problem in Deep RL
  - Actually, in other areas of ML too: NLP, time-series analysis, video processing, ...

# Recap and further considerations

---

- Let policies depend on **observable history** through **memory**
- **Memory update**: Bayesian, approximate, or learned
  - **Learning to update memory** is one of the biggest open problems in all of ML
- Let policy be **stochastic**
  - Should memory be stochastic? interesting research question...
- Let policies be **non-stationary** if possible, otherwise learning may be unstable
  - **Time-dependent** policies for finite-horizon tasks
  - **Periodic** policies for periodic tasks

# Today's lecture

---

Belief-state MDPs

RNNs

**IRL, Feature Matching**

MaxEnt IRL



# Learning rewards from demonstrations

- RL: rewards  $\rightarrow$  policy; IL: demonstrations  $\rightarrow$  policy
- Inverse Reinforcement Learning (IRL): demonstrations  $\rightarrow$  reward function
  - Better understand agents (humans, animals, users, markets)
    - Preference elicitation, teleology (the “what for” of actions), theory of mind, language
  - First step toward Apprenticeship Learning: demos  $\rightarrow$  rewards  $\rightarrow$  policy
    - Infer the teacher's goals and learn to achieve them; overcome suboptimal demos
    - Partly model-based (learn  $r$  but not  $p$ ); may be easier to learn, generalize, transfer
    - Teacher and learner can have different action spaces (e.g., human  $\rightarrow$  robot)

# Inverse Reinforcement Learning (IRL)

- Given a dataset of **demonstration trajectories**  $\mathcal{D} = \{\xi_i\}$
- Find teacher's **reward function**  $r : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ 
  - ▶ **Principle**: demonstrated actions should achieve high expected return
- IRL is **ill-defined**
  - ▶ How low is the reward for states and actions **not in**  $\mathcal{D}$ ?
  - ▶ How is the reward **distributed** along the trajectory?
    - Sparse rewards = identify “**subgoal**” states; dense = score **each step**, as hard as IL
  - ▶ Demonstrator can be **fallible** = take suboptimal actions; how much?

# Feature matching

- Assume **linear reward**  $r_\theta(s) = \theta^\top f_s$  in given **state features**  $f_s \in \mathbb{R}^d$ 
  - ▶ **Value**  $= J_\theta^\pi = \sum_t \gamma^t \mathbb{E}_{s_t \sim p_\pi} [\theta^\top f_{s_t}] = \mathbb{E}_{s \sim p_\pi} [\theta^\top f_s]$ , with  $p_\pi(s) \propto \sum_t \gamma^t p_\pi(s_t)$ 
    - $t \sim \text{Geom}(1 - \gamma)$   
**missing const:**  $(1 - \gamma)$
- **Teacher optimality:** expert value  $J_\theta^{\pi^*}$  higher than any other policy's value  $J_\theta^\pi$ 
  - ▶ Find  $\theta$  that maximizes the **gap**  $J_\theta^{\pi^*} - J_\theta^\pi$ ; but for which  $\pi$ ?
  - ▶ **Apprenticeship Learning:** find  $\pi$  that maximizes  $J_\theta^\pi$ ; but for which  $\theta$ ?
- **Solve:**  $\max_\theta \min_\pi \{ J_\theta^{\pi^*} - J_\theta^\pi \} = \max_\theta \min_\pi \{ \mathbb{E}_{s \sim p^*} [\theta^\top f_s] - \mathbb{E}_{s \sim p_\pi} [\theta^\top f_s] \}$ 
  - ▶ **Approximate**  $s \sim p^*$  with  $s \sim \mathcal{D}$

# Feature matching

- Solving  $\max_{\theta} \min_{\pi} \{ \mathbb{E}_{s \sim p^*} [\theta^\top f_s] - \mathbb{E}_{s \sim p_\pi} [\theta^\top f_s] \}$

---

## Algorithm Feature Matching

---

Initialize policy set  $\Pi = \{\pi_0\}$

**repeat**

Solve Quadratic Program:  $\max_{\eta, \|\theta\|_2 \leq 1} \eta$  ←  $\theta$  must be bounded, or solution at  $\infty$

s.t.  $\mathbb{E}_{s \sim \mathcal{D}} [\theta^\top f_s] \geq \mathbb{E}_{s \sim p_\pi} [\theta^\top f_s] + \eta \quad \forall \pi \in \Pi$

$\pi \leftarrow$  optimal policy for  $r_\theta(s) = \theta^\top f_s$

Add  $\pi$  to  $\Pi$

---

- On convergence:  $\pi$  optimal for  $\theta$  (no gap), can't find  $\theta$  with gap

feature matching

▶  $\Rightarrow \mathbb{E}_{s \sim \mathcal{D}} [\theta^\top f_s] \approx \mathbb{E}_{s \sim p_\pi} [\theta^\top f_s]$  for all  $\theta \Rightarrow \mathbb{E}_{s \sim \mathcal{D}} [f_s] \approx \mathbb{E}_{s \sim p_\pi} [f_s]$  ←

# Today's lecture

---

Belief-state MDPs

RNNs

IRL, Feature Matching

**MaxEnt IRL**

# Modeling bounded teachers

- An **expert** teacher maximizes the value  $J_{\theta}^{\pi^*} = \sum_t \gamma^t \mathbb{E}_{s_t \sim p^*} [\theta^\top f_{s_t}] = \mathbb{E}_{\xi \sim p^*} [\theta^\top f_{\xi}]$ 
  - ▶ With trajectory-summed features  $f_{\xi} = \sum_t \gamma^t f_{s_t}$
- Assume teacher has unintentional / uninformed **prior policy**  $\pi_0$ 
  - ▶ **Bounded rationality**: cost to intentionally diverge  $\mathbb{D}[\pi^* \parallel \pi_0]$  (with  $\pi_0$  uniform:  $\mathbb{H}[\pi^*]$ )
  - ▶ Total cost:  $\sum_t \mathbb{E}_{(s_t, a_t) \sim p^*} \left[ \log \frac{\pi^*(a_t | s_t)}{\pi_0(a_t | s_t)} \right] = \mathbb{E}_{\xi \sim p^*} \left[ \log \frac{p^*(\xi)}{p_0(\xi)} \right] = \mathbb{D}[p^*(\xi) \parallel p_0(\xi)]$
- **Bounded optimality**:  $\max_{\pi^*} \mathbb{E}_{\xi \sim p^*} [\theta^\top f_{\xi}] - \tau \mathbb{D}[p^* \parallel p_0]$

# Bounded optimality: naïve solution

- Bounded optimality:  $\max_{\pi^* p^*} \mathbb{E}_{\xi \sim p^*}[\theta^\top f_\xi] - \mathbb{D}[p^* || p_0]$ 
  - ▶ Naïve solution: allow **any** distribution  $p^*$  over trajectories
  - ▶ No need to be consistent with **dynamics**  $p(s' | s, a) \Rightarrow p^*$  may be **unachievable**

- Add the **constraint**  $\sum_{\xi} p^*(\xi) = 1$  with Lagrange multiplier  $\lambda$

- **Differentiate** by  $p^*(\xi)$  and  $= 0$  to optimize

$$\theta^\top f_\xi - \log p^*(\xi) + \log p_0(\xi) - 1 + \lambda = 0 \implies p^*(\xi) = \frac{p_0(\xi) \exp(\theta^\top f_\xi)}{\sum_{\bar{\xi}} p_0(\bar{\xi}) \exp(\theta^\top f_{\bar{\xi}})}$$

# IRL with bounded teacher

- Assume that **demonstrations** are distributed  $p_{\theta}(\xi) = \frac{1}{Z_{\theta}} p_0(\xi) \exp(\theta^{\top} f_{\xi})$ 
  - ▶ With **partition function**  $Z_{\theta} = \mathbb{E}_{\bar{\xi} \sim p_0} [\exp(\theta^{\top} f_{\bar{\xi}})]$
- Find  $\theta$  that **minimizes NLL** of demonstrations

$$\begin{aligned} \nabla_{\theta} \log p_{\theta}(\xi) &= \nabla_{\theta} (\theta^{\top} f_{\xi} - \log Z_{\theta}) = f_{\xi} - \frac{1}{Z_{\theta}} \nabla_{\theta} Z_{\theta} \\ &= f_{\xi} - \frac{1}{Z_{\theta}} \mathbb{E}_{\bar{\xi} \sim p_0} [\exp(\theta^{\top} f_{\bar{\xi}}) f_{\bar{\xi}}] = f_{\xi} - \mathbb{E}_{\bar{\xi} \sim p_{\theta}} [f_{\bar{\xi}}] \end{aligned}$$

- ▶ To compute gradient, we need  $p_{\theta}$ , but how to **compute**  $Z_{\theta}$ ?



# Computing $Z_\theta$ : backward recursion

- Partition function:  $Z_\theta = \mathbb{E}_{\xi \sim p_0}[\exp(\theta^\top f_\xi)]$

- Compute  $Z_\theta$  recursively **backward**: like a value function, but + becomes ·

$$Z_\theta(s_t, a_t) = \mathbb{E}_{p_0}[\exp(\theta^\top f_{\xi \geq t}) \mid s_t, a_t] = \exp(\theta^\top f_{s_t}) \mathbb{E}_{(s_{t+1} \mid s_t, a_t) \sim p}[Z_\theta(s_{t+1})]$$

$$Z_\theta(s_t) = \mathbb{E}_{p_0}[\exp(\theta^\top f_{\xi \geq t}) \mid s_t] = \mathbb{E}_{(a_t \mid s_t) \sim \pi_0}[Z_\theta(s_t, a_t)]$$

- How to get a policy from  $Z_\theta$ ?

▶ **Marginalize**:  $\pi_\theta(a_t \mid s_t) = \frac{p_\theta(\xi \mid s_t, a_t)}{p_\theta(\xi \mid s_t)} = \pi_0(a_t \mid s_t) \frac{Z_\theta(s_t, a_t)}{Z_\theta(s_t)}$

consistent  $\pi$  may not even exist

- ▶ This  $\pi_\theta$  is not globally **consistent**  $p_\theta(\xi) \neq p_{\pi_\theta}(\xi)$ ,  $p_\theta(\xi)$  ignores the **dynamics**

# MaxEnt IRL

- For each sample  $\xi \sim \mathcal{D}$ :

## Limitations:

- ▶ Compute  $Z_\theta = \mathbb{E}_{\xi \sim p_0} [\exp(\theta^\top f_\xi)]$  recursively **backward**
  - ▶ Compute  $\mathbb{E}_{\bar{\xi} \sim p_{\pi_\theta}} [f_{\bar{\xi}}]$  recursively **forward**
  - ▶ Take a gradient step to **improve**  $\theta$ :  $\nabla_\theta \log p_\theta(\xi) \approx f_\xi - \mathbb{E}_{\bar{\xi} \sim p_{\pi_\theta}} [f_{\bar{\xi}}]$
- Requires dynamics  $p$
  - Assumes  $p_\theta = p_{\pi_\theta}$
  - Assumes  $\mathcal{D} = p^*$

- At the optimum: **feature matching**  $\mathbb{E}_{\xi \sim \mathcal{D}} [f_\xi] = \mathbb{E}_{\xi \sim p_{\pi_\theta}} [f_\xi]$

- ▶ **MaxEnt IRL** approximates  $\max_{\theta} \mathbb{H}[\pi_\theta] \quad \text{s.t.} \quad \mathbb{E}_{\xi \sim \mathcal{D}} [f_\xi] = \mathbb{E}_{\xi \sim p_{\pi_\theta}} [f_\xi]$