

Review

What we covered this semester:

1. Basics of Markov Decision Process

2. Planning in MDP: VI and PI

3. Learning: Model-based RL, Policy Optimization, Bandit

4. Imitation

Basics of MDP

Understanding those widely used notations and terminologies:

$$\pi, V^\pi, Q^\pi, \pi^*, V^*, Q^*$$

$$\mathbb{P}_h^\pi(s; \mu), d_\mu^\pi(s)$$

Basics of MDP

Bellman Equation and Bellman Optimality:

$$\forall s, a : \quad Q^\pi(s, a) = r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} V^\pi(s')$$

policy evaluation

$$Q^*(s, a) = r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} \max_{a'} Q^*(s', a'), \forall s, a$$

TD DP

Planning in MDP

Q: When P and r are known, we can compute π^* via VI or PI

Algorithm 1: Value Iteration

$$\forall s, a : \underline{\underline{Q^{t+1}(s, a)}} = r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} \max_{a'} Q^t(s', a')$$

$$Q^{t+1} := \mathcal{T} Q^t$$

Planning in MDP

Q: When P and r are known, we can compute π^\star via VI or PI

Algorithm 1: Value Iteration

$$\forall s, a : Q^{t+1}(s, a) = r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} \max_{a'} Q^t(s', a')$$

Why it works?

Contraction + Q^\star being a fixed point

Planning in MDP

Q: When P and r are known, we can compute π^\star via VI or PI

Algorithm 1: Value Iteration

$$\forall s, a : Q^{t+1}(s, a) = r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} \max_{a'} Q^t(s', a')$$

Why it works?

Contraction + Q^\star being a fixed point

$$\|Q^{t+1} - Q^\star\|_\infty = \|\mathcal{T}Q^t - \underbrace{\mathcal{T}Q^\star}_{Q^\star = \mathcal{T}Q^\star}\|_\infty \leq \gamma \|Q^t - Q^\star\|_\infty$$

Planning in MDP

Q: When P and r are known, we can compute π^\star via VI or PI

Algorithm 2: Policy Iteration

Planning in MDP

Q: When P and r are known, we can compute π^\star via VI or PI

Algorithm 2: Policy Iteration

1. Policy Evaluation: $Q^{\pi^t}(s, a), \forall s, a$

Planning in MDP

Q: When P and r are known, we can compute π^\star via VI or PI

Algorithm 2: Policy Iteration

1. Policy Evaluation: $Q^{\pi^t}(s, a), \forall s, a$
2. Policy Improvement $\pi^{t+1}(s) := \arg \max_a Q^{\pi^t}(s, a), \forall s;$

Planning in MDP

Q: When P and r are known, we can compute π^* via VI or PI

Algorithm 2: Policy Iteration

1. Policy Evaluation: $Q^{\pi^t}(s, a), \forall s, a$

2. Policy Improvement $\pi^{t+1}(s) := \arg \max_a Q^{\pi^t}(s, a), \forall s;$

PDL

Key Properties:

Monotonic improvement + hit π^* in at most A^S many iterations (hw1)

What we covered this semester:

 1. Basics of Markov Decision Process

 2. Planning in MDP: VI and PI

3. Learning: Model-based RL, Policy Optimization, Bandit

4. Imitation

Learning:

Q: What we do when (P, r) are not known?

Learning:

Q: What we do when (P, r) are not known?

We considered two learning settings:

Learning:

Q: What we do when (P, r) are not known?

We considered two learning settings:

1. Generative model, i.e., we can reset to any (s, a)

Learning:

Q: What we do when (P, r) are not known?

We considered two learning settings:

1. Generative model, i.e., we can reset to any (s, a)
2. Reset from fixed initial state distribution μ ;

Learning:

Q: What we do when (P, r) are not known?

Under generative model setting, we learned a simple model-based RL alg:

Learning:

Q: What we do when (P, r) are not known?

Under generative model setting, we learned a simple model-based RL alg:

1. Model fitting:

$\forall s, a$: collect N next states, $s'_i \sim P(\cdot | s, a), i \in [N]$; set

$$\widehat{P}(s' | s, a) = \frac{\sum_{i=1}^N \mathbf{1}\{s'_i = s'\}}{N};$$

Learning:

Q: What we do when (P, r) are not known?

Under generative model setting, we learned a simple model-based RL alg:

1. Model fitting:

$\forall s, a$: collect N next states, $s'_i \sim P(\cdot | s, a), i \in [N]$; set

$$\widehat{P}(s' | s, a) = \frac{\sum_{i=1}^N \mathbf{1}\{s'_i = s'\}}{N};$$

$$\left(\widehat{r}(s, a) = \frac{\sum r_i}{N} \right)$$

2. Planning w/ the learned model:

$$\widehat{\pi}^* = \text{PI} \left(\widehat{P}, r \right)$$

Assume r is known

Learning:

Q: What we do when (P, r) are not known?

An Important Lemma that is widely used in model-based approach

Simulation Lemma:

$$\begin{aligned} \widehat{V}^\pi(s_0) - V^\pi(s_0) &= \frac{\gamma}{1-\gamma} \mathbb{E}_{s,a \sim d_{s_0}^\pi} \left[\mathbb{E}_{s' \sim \widehat{P}(\cdot | s, a)} \widehat{V}^\pi(s') - \mathbb{E}_{s' \sim P(\cdot | s, a)} \widehat{V}^\pi(s') \right] \\ &\leq \frac{\gamma}{(1-\gamma)^2} \mathbb{E}_{s,a \sim d_{s_0}^\pi} \left\| \underbrace{\widehat{P}(\cdot | s, a) - P(\cdot | s, a)} \right\|_1 \end{aligned}$$

Learning:

Q: What we do when (P, r) are not known?

Under resetting from μ , we learned policy gradient algorithms

τ
Initial Distribution

Learning:

Q: What we do when (P, r) are not known?

Under resetting from μ , we learned policy gradient algorithms

Given a differentiable parameterized policy $\pi_{\theta}(a | s)$, w/ $\theta \in \mathbb{R}^d$:

Learning:

Q: What we do when (P, r) are not known?

Under resetting from μ , we learned policy gradient algorithms

Given a differentiable parameterized policy $\pi_\theta(a | s)$, w/ $\theta \in \mathbb{R}^d$:

REINFORCE:

$$\nabla_\theta J(\pi_\theta) = \mathbb{E}_{\tau \sim \rho^{\pi_\theta}} \left[R(\tau) \underbrace{\sum_{h=0}^{H-1} \nabla_\theta \ln \pi_\theta(a_h | s_h)} \right]$$

Learning:

Q: What we do when (P, r) are not known?

Under resetting from μ , we learned policy gradient algorithms

Given a differentiable parameterized policy $\pi_\theta(a | s)$, w/ $\theta \in \mathbb{R}^d$:

REINFORCE: ✓

$$\nabla_\theta J(\pi_\theta) = \mathbb{E}_{\tau \sim \rho^{\pi_\theta}} \left[\cancel{R(\tau)} \sum_{h=0}^{H-1} \nabla_\theta \ln \pi_\theta(a_h | s_h) \right] = \mathbb{E}_{\tau \sim \rho^{\pi_\theta}} \left[\sum_{h=0}^{H-1} \nabla_\theta \ln \pi_\theta(a_h | s_h) \left(\sum_{\tau=h}^{H-1} r(s_\tau, a_\tau) \right) \right]$$

starting from h

Learning:

Q: What we do when (P, r) are not known?

Under resetting from μ , we learned policy gradient algorithms

Given a differentiable parameterized policy $\pi_\theta(a | s)$, w/ $\theta \in \mathbb{R}^d$:

The Q-version:

$$\nabla_\theta J(\pi_\theta) = \frac{1}{1 - \gamma} \mathbb{E}_{s, a \sim d_\mu^{\pi_\theta}} \left[\nabla_\theta \ln \pi_\theta(a | s) (Q^{\pi_\theta}(s, a) - b(s)) \right]$$

π : Actor

Q_w : critic \rightarrow Approximate

$Q^{\pi_\theta}(s, a)$

Learning:

Q: What we do when (P, r) are not known?

Under resetting from μ , we learned policy gradient algorithms

Given a differentiable parameterized policy $\pi_{\theta}(a | s)$, w/ $\theta \in \mathbb{R}^d$:

The Natural Policy Gradient:

$$\theta_{t+1} = \theta_t + \eta F_{\theta_t}^{-1} \nabla_{\theta} J(\pi_{\theta_t})$$

\uparrow Fisher Info-Matrix

$$KL(p^{\pi_{\theta^*}} || p^{\pi_{\theta_t}})$$

$$\mathbb{D} = (\theta - \theta_t)^{\top} F_{\theta} (\theta - \theta_t)$$

\uparrow
PSD

$$\text{if } \max_{s,a} \frac{d\pi^*}{J(s,a)} \leq c$$

Learning:

Q: What we do when (P, r) are not known?

Under resetting from μ , we learned policy gradient algorithms

Given a differentiable parameterized policy $\pi_\theta(a | s)$, w/ $\theta \in \mathbb{R}^d$:

The Natural Policy Gradient:

$$\theta_{t+1} = \theta_t + \eta F_{\theta_t}^{-1} \nabla_{\theta} J(\pi_{\theta_t})$$

Regular GD

↓

Instead of using Euclidean distance metric, we use local geometry metric

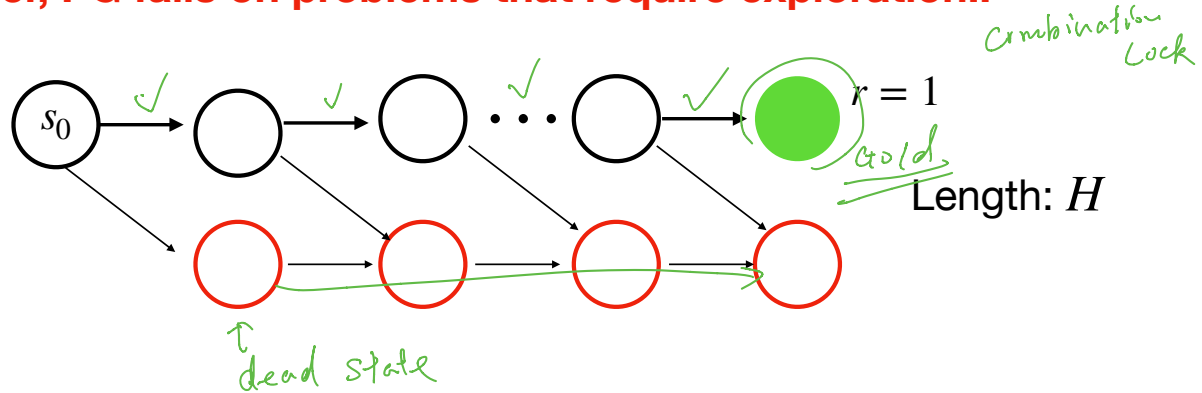
$$\underline{d(\theta, \theta_t) := (\theta - \theta_t)^\top F_{\theta_t} (\theta - \theta_t)}$$

$$d(\theta, \theta_t) = (\theta - \theta_t)^\top F_{\theta_t} (\theta - \theta_t)$$

Learning:

Q: What we do when (P, r) are not known?

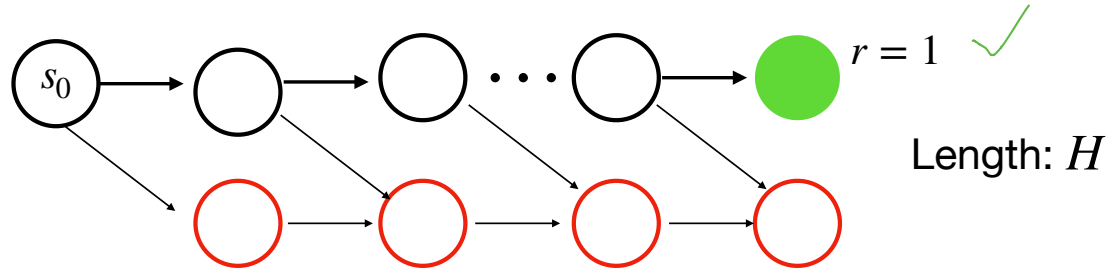
However, PG fails on problems that require exploration..



Learning:

Q: What we do when (P, r) are not known?

However, PG fails on problems that require exploration..



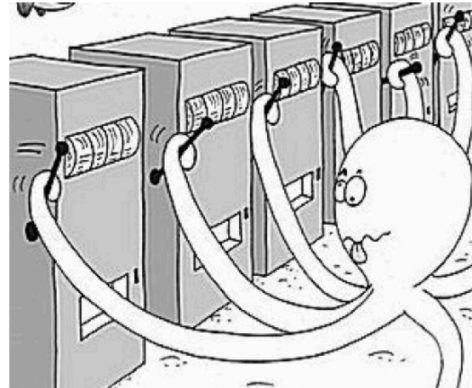
What is the probability of a random policy generating a trajectory that hits the goal?

$$2^{-H}$$

Learning:

Q: How to learn efficient (i.e., balance explore and exploit) in Multi-armed Bandit setting:

We have K many arms (or actions): a_1, \dots, a_K



Learning:

Q: How to learn efficient (i.e., balance explore and exploit) in Multi-armed Bandit setting:

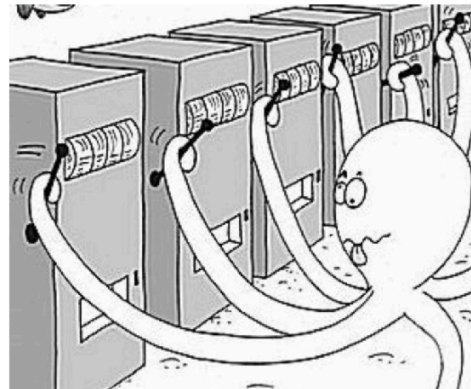
We have K many arms (or actions): a_1, \dots, a_K

Each arm has a unknown reward distribution,

i.e., $\nu_i \in \Delta([0,1])$,

w/ mean $\mu_i = \mathbb{E}_{r \sim \nu_i}[r]$

$$I^* = \operatorname{argmax}_{i \in [K]} \mu_i$$



Learning:

Q: How to learn efficient in Multi-armed Bandit setting:

1. Explore and ~~Committee~~ algorithm:

Learning:

Q: How to learn efficient in Multi-armed Bandit setting:

1. Explore and Committee algorithm:

1. For the first NK rounds, try each arm N times, compute its average mean $\hat{\mu}_i$

Learning:

Q: How to learn efficient in Multi-armed Bandit setting:

1. Explore and Committee algorithm:

1. For the first NK rounds, try each arm N times, compute its average mean $\hat{\mu}_i$
2. For all future $T-KN$ rounds, play the best empirical arm $\hat{I} = \arg \max_{i \in [K]} \hat{\mu}_i$

$$\text{Regret}_T = \mathcal{O}(T^{2/3} K^{1/3})$$

Learning:

Q: How to learn efficient in Multi-armed Bandit setting:

2. The Upper Confidence Bound Algorithm

Learning:

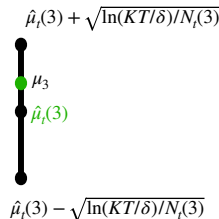
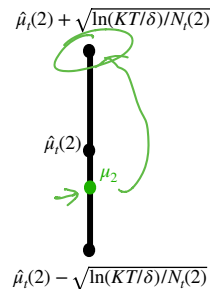
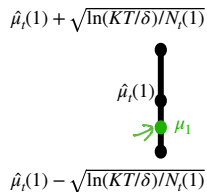
Q: How to learn efficient in Multi-armed Bandit setting:

2. The Upper Confidence Bound Algorithm

For $t = 0 \rightarrow T - 1$:

$$I_t = \arg \max_{i \in [K]} \left(\hat{\mu}_t(i) + \sqrt{\frac{\ln(KT/\delta)}{N_t(i)}} \right)$$

(# Upper-conf-bound of arm i)



$$\text{Regret}_T = O(\sqrt{KT})$$

Learning:

Q: How to learn efficient in Contextual Bandit setting:

1. Explore and Commit (or ϵ -greedy)

Importance weighting + Reward-sensitive Classification

Learning:

Q: How to learn efficient in Contextual Bandit setting:

1. Explore and Commit (or ϵ -greedy)

Importance weighting + Reward-sensitive Classification

1. For the first N rounds, randomly try actions to construct a classification dataset:

$$\{x_i, \hat{\mathbf{r}}_i\}_{i=0}^{N-1}$$

$\hat{\mathbf{r}}_i$
 $x_i \sim \mu$

IW \rightarrow $E[\hat{\mathbf{r}}] = r$

$$\hat{\mathbf{r}} := \begin{bmatrix} 0 \\ 0 \\ \dots \\ r_i/p(a_i) \\ 0, \\ \dots \\ 0 \end{bmatrix}$$

Learning:

Q: How to learn efficient in Contextual Bandit setting:

1. Explore and Commit (or ϵ -greedy)

Importance weighting + Reward-sensitive Classification

1. For the first N rounds, randomly try actions to construct a classification dataset:

$$\{x_i, \hat{\mathbf{r}}_i\}_{i=0}^{N-1}$$

2. Call RSC oracle: $\hat{\pi} = \arg \max_{\pi \in \Pi} \sum_{i=0}^{N-1} \hat{\mathbf{r}}_i[\pi(x_i)]$

$$\hat{\mathbf{r}} := \begin{bmatrix} 0 \\ 0 \\ \dots \\ r_i/p(a_i) \\ 0, \\ \dots \\ 0 \end{bmatrix}$$

$O(T^{2/3} K^{1/3})$

What we covered this semester:

 1. Basics of Markov Decision Process

 2. Planning in MDP: VI and PI

 3. Learning: Model-based RL, Policy Optimization, Bandit

4. Imitation

Imitation:

Q: what we do when r is not available but we have an expert $\pi^e (\approx \pi^\star)$?

1. Offline IL: only expert data $\{s_i^\star, a_i^\star\}_{i=1}^N$ is available (no other interaction)

Imitation:

Q: what we do when r is not available but we have an expert $\pi^e (\approx \pi^\star)$?

1. Offline IL: only expert data $\{s_i^\star, a_i^\star\}_{i=1}^N$ is available (no other interaction)

BC: a Reduction to Supervised Learning:

Imitation:

Q: what we do when r is not available but we have an expert $\pi^e (\approx \pi^*)$?

1. Offline IL: only expert data $\{s_i^*, a_i^*\}_{i=1}^N$ is available (no other interaction)

BC: a Reduction to Supervised Learning:

$$\hat{\pi} = \arg \min_{\pi \in \Pi} \sum_{i=1}^M \ell(\pi, s^*, a^*)$$

loss function for classification/regression

Imitation:

Q: what we do when r is not available but we have an expert $\pi^e (\approx \pi^\star)$?

1. Offline IL: only expert data $\{s_i^\star, a_i^\star\}_{i=1}^N$ is available (no other interaction)

BC: a Reduction to Supervised Learning:

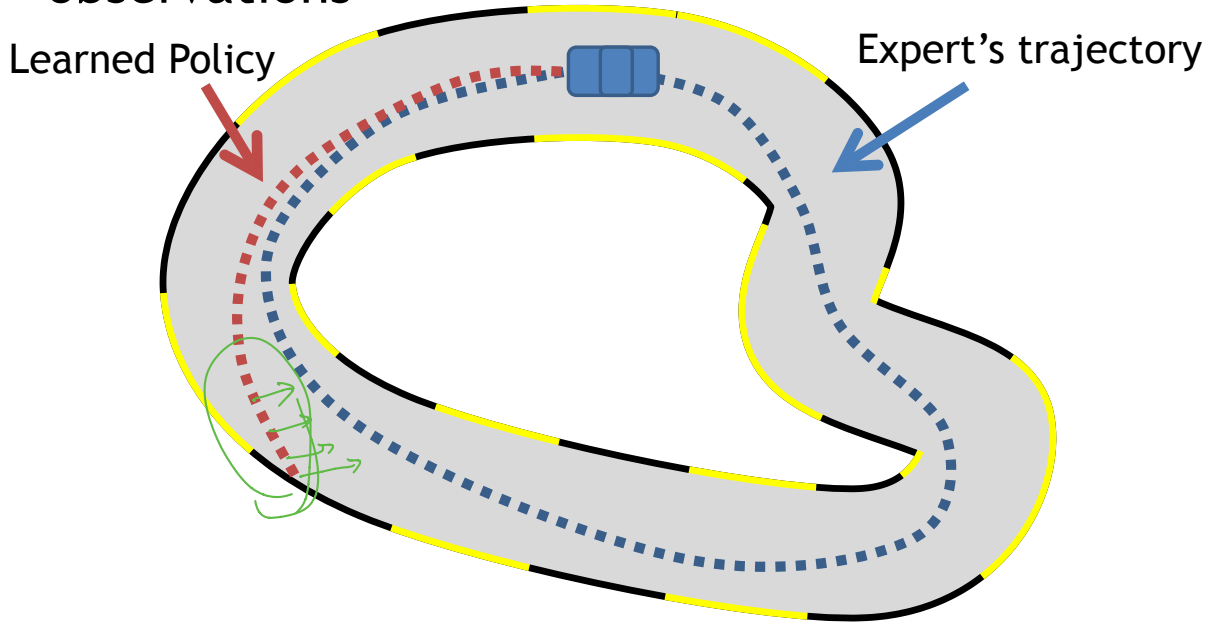
$$\hat{\pi} = \arg \min_{\pi \in \Pi} \sum_{i=1}^M \ell(\pi, s_i^\star, a_i^\star)$$

e.g., Negative log-likelihood (NLL): $\ell(\pi, s, a^\star) = -\ln \pi(a^\star | s^\star)$ (used in AlphaGo)

Distribution shift!

[Pomerleau89, Daume09]

- Predictions affect future inputs/ observations



Imitation:

Q: what we do when r is not available but we have an expert $\pi^e (\approx \pi^\star)$?

1. Interactive IL: expert is available for query during the learning process

Imitation:

Q: what we do when r is not available but we have an expert $\pi^e (\approx \pi^\star)$?

1. Interactive IL: expert is available for query during the learning process

The DAgger Algorithm (Data Aggregation):

Imitation:

Q: what we do when r is not available but we have an expert $\pi^e (\approx \pi^*)$?

1. Interactive IL: expert is available for query during the learning process

The DAgger Algorithm (Data Aggregation):

1. W/ π^t , generate dataset $\mathcal{D}^t = \{s_i, a_i^*\}$, $s_i \sim d_\mu^{\pi^t}$, $a_i^* = \pi^*(s_i)$
← Expert action
from our policy π^t

Imitation:

Q: what we do when r is not available but we have an expert $\pi^e (\approx \pi^\star)$?

1. Interactive IL: expert is available for query during the learning process

The DAgger Algorithm (Data Aggregation):

1. W/ π^t , generate dataset $\mathcal{D}^t = \{s_i, a_i^\star\}$, $s_i \sim d_\mu^{\pi^t}$, $a_i^\star = \pi^\star(s_i)$
2. **Data aggregation:** $\mathcal{D} = \mathcal{D} + \mathcal{D}^t$

Imitation:

Q: what we do when r is not available but we have an expert $\pi^e (\approx \pi^\star)$?

1. Interactive IL: expert is available for query during the learning process

The DAgger Algorithm (Data Aggregation):

1. W/ π^t , generate dataset $\mathcal{D}^t = \{s_i, a_i^\star\}$, $s_i \sim d_\mu^{\pi^t}$, $a_i^\star = \pi^\star(s_i)$
2. **Data aggregation:** $\mathcal{D} = \mathcal{D} + \mathcal{D}^t$
3. **Update policy via Supervised-Learning:** $\pi^{t+1} = \text{SL}(\mathcal{D})$

Imitation:

Q: what we do when r is not available but we have an expert $\pi^e (\approx \pi^\star)$?

1. Interactive IL: expert is available for query during the learning process

The DAgger Algorithm (Data Aggregation):

1. W/ π^t , generate dataset $\mathcal{D}^t = \{s_i, a_i^\star\}$, $s_i \sim d_\mu^{\pi^t}$, $a_i^\star = \pi^\star(s_i)$
2. **Data aggregation:** $\mathcal{D} = \mathcal{D} + \mathcal{D}^t$
3. **Update policy via Supervised-Learning:** $\pi^{t+1} = \text{SL}(\mathcal{D})$

(Recap the connection to online learning and how it avoids distribution shift..)

Imitation:

Q: what we do when r is not available but we have an expert $\pi^e (\approx \pi^\star)$?

3. Hybrid IL: expert data $\{s_i^\star, a_i^\star\}_{i=1}^N$, and access to real transition $P(\cdot | s, a)$

Imitation:

Q: what we do when r is not available but we have an expert $\pi^e (\approx \pi^\star)$?

3. Hybrid IL: expert data $\{s_i^\star, a_i^\star\}_{i=1}^N$, and access to real transition $P(\cdot | s, a)$

Formulation of Maximum Entropy Inverse RL:

Imitation:

Q: what we do when r is not available but we have an expert $\pi^e (\approx \pi^*)$?

3. Hybrid IL: expert data $\{s_i^*, a_i^*\}_{i=1}^N$, and access to real transition $P(\cdot | s, a)$

Formulation of Maximum Entropy Inverse RL:

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

Imitation:

Q: what we do when r is not available but we have an expert $\pi^e (\approx \pi^\star)$?

3. Hybrid IL: expert data $\{s_i^\star, a_i^\star\}_{i=1}^N$, and access to real transition $P(\cdot | s, a)$

Formulation of Maximum Entropy Inverse RL:

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

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

Imitation:

Q: what we do when r is not available but we have an expert $\pi^e (\approx \pi^\star)$?

3. Hybrid IL: expert data $\{s_i^\star, a_i^\star\}_{i=1}^N$, and access to real transition $P(\cdot | s, a)$

Formulation of Maximum Entropy Inverse RL:

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

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

Assume the ground truth reward $r(s, a) = (\theta^\star)^\top \phi(s, a)$

Imitation:

Q: what we do when r is not available but we have an expert $\pi^e (\approx \pi^*)$?

3. Hybrid IL: expert data $\{s_i^*, a_i^*\}_{i=1}^N$, and access to real transition $P(\cdot | s, a)$

Formulation of Maximum Entropy Inverse RL:

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

Lagrangian multipliers

Imitation:

Q: what we do when r is not available but we have an expert $\pi^e (\approx \pi^*)$?

3. Hybrid IL: expert data $\{s_i^*, a_i^*\}_{i=1}^N$, and access to real transition $P(\cdot | s, a)$

Incremental update

Formulation of Maximum Entropy Inverse RL:

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

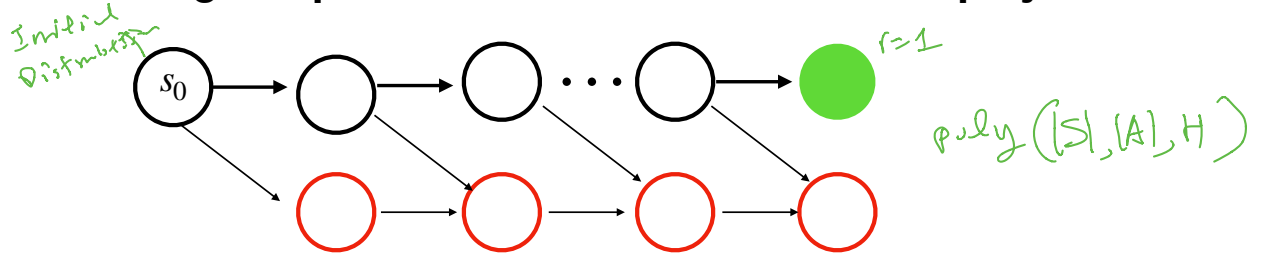
Base Response

Iterate: $w_{t+1} = w_t + \eta \nabla_w \ell(w_t, \pi_t)$, $\pi^{t+1} = \arg \min_{\pi} \ell(w_{t+1}, \pi)$ via Soft-VI

Gradient Ascent on w *Base Response on π*

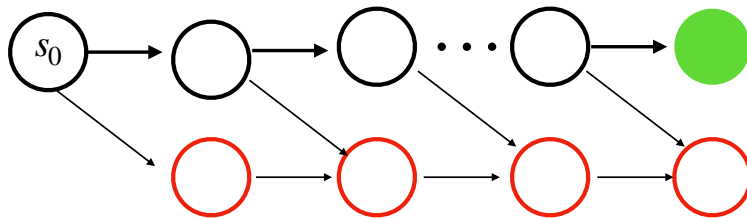
What we did not cover:

1. How to do strategic exploration in RL? Can we do it in poly time?



What we did not cover:

1. How to do strategic exploration in RL? Can we do it in poly time?

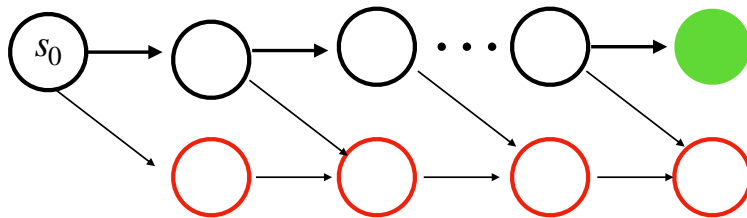


2. When does Policy Gradient guarantee Global optimality?

Though the RL objective function is non-convex wrt policy, under some cases, PG provably converges to global optimal policies!

What we did not cover:

1. How to do strategic exploration in RL? Can we do it in poly time?



2. When does Policy Gradient guarantee Global optimality?

Though the RL objective function is non-convex wrt policy, under some cases, PG provably converges to global optimal policies!

3. Deep Reinforcement Learning

Most of the time, it is Deep nets (e.g., policies) + RL