

Without ground truth: how do we know which action is better or worse?

$G(s, a)$  = reward function. Which states/actions are better?

↳ Greedy system: always chooses high reward short-term. Not effective, we want long-term.

Markov decision process (MDP) - extension of Markov chain.

$M = \{S, A, T, r\}$

$S$  - state space, states  $s \in S$  (discrete/continuous)

$A$  - action space, actions  $a \in A$  (discrete/continuous)

$r$  - reward function  $r: S \times A \rightarrow \mathbb{R}$

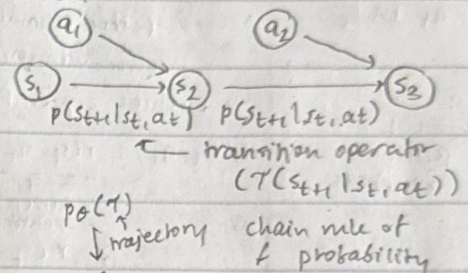
partially observed MDP -  $M = \{S, A, O, \gamma, E, r\}$

$O$  - observation space  $o \in O$  (discrete/continuous)

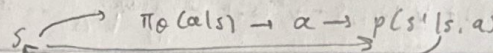
↳ refer to Bayes Net diagram in lecture 14.

$E$  - emission probability  $p(o_t | s_t)$

(short)



The goal of RL:



"J(theta)"  $\theta^* = \arg \max_{\theta} E_{\tau \sim p_{\theta}(T)} \left[ \sum_t r(s_t, a_t) \right]$  - pick the policy parameters that give highest expected reward.

↳ exponentially large, fix via unbiased estimator (sample)

$$J(\theta) \approx \frac{1}{N} \sum_{i=1}^N \sum_t r(s_{i,t}, a_{i,t})$$

sum over samples from  $\pi_{\theta}$

$$J(\theta) = E_{\tau \sim \pi_{\theta}(T)} [r(\tau)] = \int \pi_{\theta}(\tau) r(\tau) d\tau$$

sum of rewards

direct  
policy  
differentiation  
(gradient  
ascent)  
preparation  
↓

$$\nabla_{\theta} J(\theta) = \int \nabla_{\theta} \pi_{\theta}(\tau) r(\tau) d\tau$$

$$= \int \pi_{\theta}(\tau) \nabla_{\theta} \log \pi_{\theta}(\tau) r(\tau) d\tau$$

$$= E_{\tau \sim \pi_{\theta}(T)} [\nabla_{\theta} \log \pi_{\theta}(\tau) r(\tau)]$$

$$= E_{\tau \sim \pi_{\theta}(T)} \left[ \left( \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \right) \left( \sum_{t=1}^T r(s_t, a_t) \right) \right]$$

$$\approx \frac{1}{N} \sum_{i=1}^N \left( \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(a_{i,t} | s_{i,t}) \right) \left( \sum_{t=1}^T r(s_{i,t}, a_{i,t}) \right)$$

$$\approx \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(a_{i,t} | s_{i,t}) r(\tau_i)$$

"good stuff is made more likely" aka "assisted" trial/error.

Causality: policy at time  $t'$  cannot affect reward at time  $t$  when  $t < t'$ .

↳ small fix:  $\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(a_{i,t} | s_{i,t}) \left( \sum_{t'=t}^T r(s_{i,t'}, a_{i,t'}) \right)$  ← less "noise" (variance).

Baselines

→ we want to "normalize" trajectories

$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(a_{i,t} | s_{i,t}) [r(s_{i,t}, a_{i,t}) - b]$$

b =  $\frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T r(s_{i,t}, a_{i,t})$ . Proof:  $E[\nabla_{\theta} \log \pi_{\theta}(\tau) b] = \int \pi_{\theta}(\tau) \nabla_{\theta} \log \pi_{\theta}(\tau) b d\tau = b \int \nabla_{\theta} \pi_{\theta}(\tau) d\tau = b \nabla_{\theta} \int \pi_{\theta}(\tau) d\tau = b \nabla_{\theta} 1 = b \nabla_{\theta} 1 = 0$

Policy gradient = on-policy (we must generate new samples every time)

↳ if we don't have samples from  $\pi_{\theta}(T)$ , but  $\tilde{\pi}(T)$ : importance sampling

$$E_{\pi}(f(x)) = E_{\tilde{\pi}} \left[ f(x) \frac{\pi(x)}{\tilde{\pi}(x)} \right]$$

$$J(\theta) = E_{\tau \sim \tilde{\pi}(T)} \left[ \frac{\pi_{\theta}(\tau)}{\tilde{\pi}(\tau)} r(\tau) \right]$$

importance weight =  $\frac{\pi_{\theta}(a_t | s_t)}{\tilde{\pi}(a_t | s_t)}$

$$\nabla_{\theta} J(\theta) = E_{\tau \sim \tilde{\pi}(T)} \left[ \frac{\pi_{\theta}(\tau)}{\tilde{\pi}(\tau)} \nabla_{\theta} \log \pi_{\theta}(\tau) r(\tau) \right], \theta \neq \theta^*$$

$$\approx \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \frac{\pi_{\theta}(s_{i,t}, a_{i,t})}{\tilde{\pi}(s_{i,t}, a_{i,t})} \nabla_{\theta} \log \pi_{\theta}(a_{i,t} | s_{i,t}) \hat{Q}_{i,t}$$

$$\tilde{J}(\theta) \approx \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \log \pi_{\theta}(a_{i,t} | s_{i,t}) \hat{Q}_{i,t}$$

↑  
reward

for backprop: same NLL loss, except we weigh by  $\hat{Q}_{i,t}$ . Minimal code change.

off-policy  
policy  
gradient