

1) Infinite Horizon Discounted MDP

$$\mathcal{M} = \{ \mathcal{S}, \mathcal{A}, P, r, \gamma \}$$

\mathcal{S} : space of possible states $s \in \mathcal{S}$

\mathcal{A} : space of possible actions $a \in \mathcal{A}$

P : transition function

$$P: \mathcal{S} \times \mathcal{A} \rightarrow \Delta(\mathcal{S})$$

↑
probability
distribution

r : reward function

$$r: \mathcal{S} \times \mathcal{A} \rightarrow \Delta(\mathbb{R})$$

γ : discount factor

$$0 < \gamma < 1$$

$$[0, 1]$$

cumulative reward:

$$\sum_{t=0}^{\infty} \gamma^t r_t$$

Aside: deterministic $f: X \rightarrow Y$ $f(x)=y$
encode as stochastic $F: X \rightarrow \Delta(Y)$

$$\rightarrow F(x) = \begin{cases} f(x) & \text{w.p. } 1 \end{cases}$$

sometimes overload notation for
deterministic policies, reward fns.

ex. $a = \pi(s)$ $r_t = r(s_t, a_t)$

Adopt notation

$$F(y|x) = \mathbb{P}\{F(x) = y\}$$

e.g. $\pi(a|s)$

Goal of RL:

find a policy $\pi: \mathcal{S} \rightarrow \Delta(\mathcal{A})$
maximize the ^{expected} discounted
cumulative reward.

π^*

maximize $\mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \right]$

$$a_t \sim \pi(s_t)$$

$$s_{t+1} \sim P(s_{t+1} | s_t, a_t)$$

s_0 given

2) Value Function & Qfunction

allow us to reason about
long term effects of policies

$$V^{\pi}(s) = \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \mid s_0 = s \right]$$

$s_{t+1} \sim P(s_{t+1} | s_t, a_t)$
 $a_t \sim \pi(s_t)$

$$Q^\pi(s, a) = \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \mid s_0 = s, a_0 = a \right]$$

$$s_{t+1} \sim P(s_t, a_t)$$

$$a_t \sim \pi(s_t)$$

Bellman Equation

Notice that

$$\underbrace{\sum_{t=0}^{\infty} \gamma^t r_t}_{t=0} = r_0 + \underbrace{\sum_{t=1}^{\infty} \gamma^{t-1} r_t}_{t=1}$$

replace $t = t' + 1$

$$\underbrace{\sum_{t'=0}^{\infty} \gamma^{t'} r_{t'+1}}_{t'=0}$$

assume deterministic π & r

We can write

$$V^\pi(s) = r(s, \pi(s)) + \gamma \mathbb{E} [V^\pi(s') \mid s' \sim P(s, \pi(s))]$$

$$Q^\pi(s, a) = r(s, a) + \gamma \mathbb{E} [V^\pi(s') \mid s' \sim P(s, a)]$$

What property

$$\mathbb{E} [V^\pi(s') \mid s' \sim P(s, a)]$$

$$\begin{aligned}\sum_{t=0}^{\infty} \delta^t r_t &= r_0 + \delta r_1 + \delta^2 r_2 + \dots \\ &= r_0 + \delta \left(r_1 + \delta r_2 + \dots \right) \\ &\quad \underbrace{\hspace{10em}} \\ &\quad \sum_{t=0}^{\infty} \delta^t r_{t+1}\end{aligned}$$

3) Policy Evaluation

How good is a policy?
In terms of the value function.

Given MDP $M = \{S, A, P, \gamma, r\}$
and a policy π , what is V^π ?

Bellman equation:

$$\forall s \in S, \quad V^\pi(s) = r(s, \pi(s)) + \gamma \sum_{s' \in S} P(s'|s, \pi(s)) V(s')$$

$|S| = S$
states

S equations

S unknowns ($V^\pi(s)$)

in vector/matrix notation

$$V \in \mathbb{R}^S = r \in \mathbb{R}^S + \gamma \sum_{s' \in S} P(s'|s, \pi(s)) V(s')$$

$V \in \mathbb{R}^S$ $r \in \mathbb{R}^S$ $P \in \mathbb{R}^{S \times S}$ $V \in \mathbb{R}^S$

solving linear equation

$$\left. \begin{aligned} V &= R + \gamma P V \\ V - \gamma P V &= R \end{aligned} \right\} V = (I - \gamma P)^{-1} R$$

$$(I - \gamma P) V = R \quad \text{HWO}$$

Exact solution! $O(S^3)$

4) Approximate Policy Evaluation

$$V = R + \gamma P V$$

Algorithm

initialize V^0
for $t = 0, 1, \dots, T$

$$V^{t+1} \leftarrow R + \gamma P V^t$$

$O(T S^2)$

Lemma: $\|V^{t+1} - V^\pi\|_\infty \leq \gamma \|V^t - V^\pi\|_\infty$

\uparrow me

Proof

$$\|V^{t+1} - V^\pi\|_\infty = \|R + \gamma P V^t - V^\pi\|_\infty$$

$$V^\pi = R + \gamma P V^\pi$$

$$= \gamma \|P(V^t - V^\pi)\|_\infty$$

@ index s , $|\mathbb{E}_{s' \sim P(s, \pi(s))} [V^t(s') - V^\pi(s')]|$

$$\leq \mathbb{E}_{s' \sim P(s, \pi(s))} |V^t(s') - V^\pi(s')|$$

$$\|P(V^t - V^\pi)\|_\infty \leq \max_s \mathbb{E}_{s' \sim P} [|V^t(s') - V^\pi(s')|]$$

$$\leq \max_s \max_{s'} |V^t(s') - V^\pi(s')|$$

$$\leq \|V^t - V^\pi\|_\infty$$

Theorem (convergence) after T iterations,

$$\|V^T - V^{\pi}\|_{\infty} \leq \gamma^T \|V^0 - V^{\pi}\|_{\infty}$$

What T give ϵ -approximation?

$$T \geq \frac{\log\left(\frac{\|V^0 - V^{\pi}\|_{\infty}}{\epsilon}\right)}{\log(1/\gamma)}$$

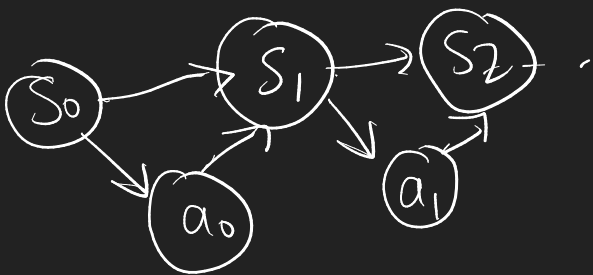
$$O(S^2 \log(1/\epsilon))$$

5) State-Action Distribution

Trajectory of MDP up to t ,
 $(s_0, a_0, s_1, a_1, \dots, s_t, a_t)$

Consider possible stochastic policies

$$P^{\pi}(s_0, a_0, \dots, s_t, a_t) = \pi(a_0 | s_0) P(s_1 | s_0, a_0) \\ \times \pi(a_1 | s_1) P(s_2 | s_1, a_1)$$



$$\times P(s_t | s_{t-1}, a_{t-1}) \\ \times \pi(a_t | s_t)$$

$$\mathbb{P}_t^\pi(s, a; s_0) = \sum_{\substack{a_{0:t-1} \in \mathcal{A}^t \\ s_{1:t-1} \in \mathcal{S}^{t-1}}} \mathbb{P}^\pi(s_0, a_0, \dots, s_t, a_t)$$