

Markov Decision Processes

In bandits, we were considering a static environment. MDPs generalize this setting.

- State space S , assumed to be countable
- Action space A , assumed to be finite
- Say $X_t = s$ (state), $A_t = a$
- cost $C(s, a)$, may be stochastic or deterministic. In former case, use expected values
- $X_{t+1} = s'$ w.p. $P_{ss'}(a)$ $\left\{ \sum_{s'} P_{ss'}(a) = 1 \right\}$

$$P(X_{t+1} = s' | X_0 = x_0, A_0 = a_0, \dots, X_t = s, A_t = a) = P_{ss'}(a)$$

- Policies are mappings from histories to actions
- An MAB is a trivial example of MDPs with only one state

Stationary policy: A_t is a deterministic function of S_t . $\pi: S \rightarrow A$

Two cost metrics:

- Discounted average cost $\mathbb{E}_\pi \left[\sum_{t=0}^{\infty} \alpha^t C(X_t, A_t) \right]$

$\{ \alpha \in (0, 1) \text{ is a discount factor} \}$

The summation may diverge when $\alpha = 1$

We also need $|C(s, a)| \leq M \quad \forall s, a$

- Average cost

$$\lim_{t \rightarrow \infty} \mathbb{E}_\pi \left[\frac{1}{t+1} \sum_{0}^t C(X_t, A_t) \right]$$

Discounted Average Cost MDP

(S, A, P, C, α)

P : transition kernel

For any policy π
 $V_\pi(i) = \mathbb{E}_\pi \left[\sum_0^\infty \alpha^t C(X_t, A_t) \mid X_0 = i \right] \rightarrow$ Value function for policy π

Define $V^*(i) = \inf_\pi V_\pi(i) \rightarrow$ Optimal value function

- This is a pointwise infimum. We first have to fix the starting state

Def: Policy π is optimal if $V_\pi(i) = V^*(i) \forall i$

Note that there is no adversary in this setting, this is a generalization of the stochastic MAB setting

Theorem: Optimal value function V^* satisfies

$$V^*(i) = \min_a [C(i, a) + \alpha \sum P_{ij}(a) V^*(j)]$$

Bellman Equations

Note that this is just a necessity

Proof: (LB)

Consider π

$$V_\pi(i) = \sum_a P_{ia} [C(i, a) + \alpha \sum P_{aj}(\pi) W_\pi(j)]$$

$P_{\pi, a}$: starting at $X_0 = i$, $P(\text{action } a \text{ is played by } \pi)$

Claim: $W_\pi(j) \geq V^*(j)$

Q: Is $W_\pi(j) = V^*(j)$?

A: No, unless we have a Markovian policy. In general we can have policies that are time-dependent

$$\begin{aligned} V_\pi(i) &\geq \sum_a P_{ia} [C(i, a) + \alpha \sum P_{aj}(\pi) V^*(j)] \\ &\geq \min_a [C(i, a) + \alpha \sum P_{aj}(\pi) V^*(j)] \end{aligned}$$

$$\Rightarrow \inf_\pi V_\pi(i) = V^*(i) \geq \underline{\hspace{2cm}}$$

(UB)

let a_0 be s.t.

$$a_0 = \operatorname{argmin} [C(i, a) + \alpha \sum P_{aj}(a) V^*(j)]$$

Fix $\epsilon > 0$. π : Play a_0 at time 0, then from state j at time 1, play π_j

$$\text{s.t. } V_{\pi_j}(j) \leq V^*(j) + \epsilon, \text{ for } \epsilon > 0$$

because $V^*(j)$ is the infimum

$$V_{\pi}(i) = C(i, a_0) + \alpha \sum P_{ij}(a_0) V_{\pi}(j) \\ \leq C(i, a) + \alpha \sum P_{ij}(a_0) V^*(j) + \alpha \epsilon$$

$$\therefore V^*(i) \leq V^{\pi}(i) \leq \min_a [C(i, a) + \alpha \sum P_{ij}(a) V^*(j)] + \alpha \epsilon \\ \therefore V^*(i) \leq \min_a [C(i, a) + \alpha \sum P_{ij}(a) V^*(j)] + \alpha \epsilon \forall \epsilon$$

\Rightarrow upper bound follows letting $\epsilon \downarrow 0$

$\therefore LB + UB \Rightarrow V^*$ satisfies Bellman equation

$BC(s)$: set of bounded value functions on S .

Def: For stationary policy f , $T_f : BC(s) \rightarrow BC(s)$

$$(T_f u)(i) = C(i, f(i)) + \alpha \sum_{a=f(i)} P_{ij}(a) u(j)$$

u is a mapping from $S \rightarrow \mathbb{R}$

$(T_f u)$ is also a mapping from $S \rightarrow \mathbb{R}$

Lemma For $u, v \in BC(s)$, stationary policy f ,

$$(i) u \leq v \Rightarrow T_f u \leq T_f v$$

$$(ii) T_f v_f = v_f$$

$$(iii) T_f^n u \rightarrow v_f \forall u$$

EXERCISE

Note (i) is an elementwise inequality. The intuition is that u, v are like terminal costs

(ii) v_f is the value function associated with the policy f . Hence, v_f is a fixed point of T_f .

(iii) Every element of $BC(s)$ can be thought out of as a vector. If we take any vector $u \in BC(s)$ and keep on applying T_f to it, there will be convergence towards v_f , which is the fixed point of T_f .

* $u_n \rightarrow u \Leftrightarrow u_n(i) \rightarrow u(i)$ uniformly in i

$$(T_f^2 u)(i) = C(i, f(i)) + \alpha \sum P_{ij}(f(i)) (T_f u)(j)$$

$$= C(i, f(i)) + \alpha \sum P_{ij}(f(i)) [C(j, f(j)) + \alpha \sum P_{jk}(f(j)) u(k)]$$

$$= C(i, f(i)) + \alpha \sum P_{ij}(f(i)) C(j, f(j)) + \alpha^2 \dots$$

expected cost at the second time-step

Thm: let π^* be stationary policy that attains the $\min_a [C(i,a) + \alpha \sum_j P_{ij}(a) V^*(j)]$

Then $V_{\pi^*} = V^*$ (i.e. π^* is optimal)

The optimality is over the space of all policies

Proof sketch

Apply T_{π^*} to V^*

Claim: $T_{\pi^*} V^* = V^*$

$$(T_{\pi^*})^n V^* = V^*$$

$\downarrow n \uparrow \infty$

$$V_{\pi^*}$$