

AVERAGE COST MDPs

Goal : Minimize

$$\Phi_{\pi}(i) = \limsup_{n \rightarrow \infty} \frac{\mathbb{E}_{\pi} \left[\sum_{t=0}^n C(X_t, A_t) \mid X_0 = i \right]}{n+1}$$

(S, A, P, C)

$$\Phi^*(i) = \inf_{\pi} \Phi_{\pi}(i)$$



S : countable

A : finite

π^* is optimal if $\Phi_{\pi^*}(i) = \Phi^*(i) \forall i$

Eg $S = \{1, 1', 2, 2', 3, 3', \dots\}$

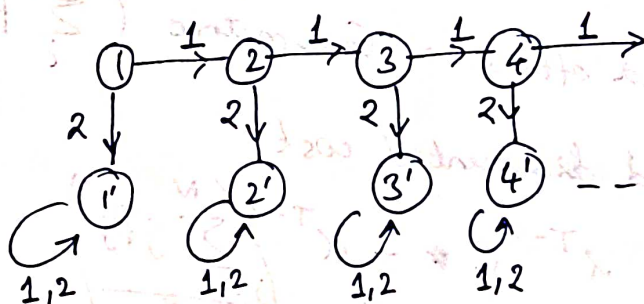
$A = \{1, 2\}$

$P_{i, i+1}(1) = P_{i, i'}(2) = 1$

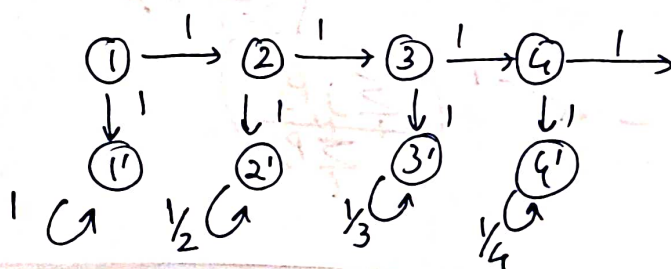
$P_{i' i'}(1) = P_{i' i}(2) = 1$

$C(i, \cdot) = 1$

$C(i', \cdot) = 1/i$



transitions are represented over arrows
The number denotes action chosen



costs associated with the transitions are on the arrows

Claim: This MDP has no optimal policy

For any policy π , $\phi_\pi(1) > 0$

But can make $\phi_\pi(1)$ arbitrarily close to zero

Note that if we throw in a discount factor, there will exist a stationary optimal policy

Ex $S = N$
 $A = \{1, 2\}$

$$P_{i,i+1}(1) = P_{i,i}(2) = 1$$

$$C(i,1) = 1; C(i,2) = 1/i$$

Under any stationary policy π , $\phi_\pi(1) > 0$

In state i , play action 2 i times, then action 1

Cost seq. $1, 1, \frac{1}{2}, \frac{1}{2}, 1, \frac{1}{3}, \frac{1}{3}, \frac{1}{3}, 1, \dots$

This policy will have a time-average of 0 as $t \rightarrow \infty$

$$\phi_\pi(1) = 0$$

Ex Thm can also find a randomized stationary policy which is optimal.

Thm: If there exists $h \in B(S)$ and constant g s.t.

$$g + h(i) = \min_a [C(i,a) + \sum_j P_{ij}(a)h(j)] \quad \forall i \quad (1)$$

then there exists an optimal stationary policy π^* , where

$$\pi^*(i) = \arg \min_a [C(i,a) + \sum_j P_{ij}(a)h(j)],$$

$$\text{and } \phi_{\pi^*}(i) = g \quad \forall i$$

Note: The correct interpretation of g is the optimal time-average cost starting from the state i associated with the optimal policy

(1) becomes the Bellman Equation for a Time-average cost MDP.

The g has to be unique.

h is not unique, because we can add an additive constant to h and it will still satisfy Bellman equation. This was not the case for discounted MDP, because the α acted as a counterweight.

Proof Under any policy π ,

$$\mathbb{E}_\pi[h(X_t) | H_{t-1}] = \sum_j h(j) P_{X_{t-1},j}(A_{t-1})$$

$$= C(X_{t-1}, A_{t-1}) + \sum_j h(j) P_{X_{t-1},j}(A_{t-1}) - C(X_{t-1}, A_{t-1})$$

$$\geq \min_a [C(x_{t-1}, a) + \sum h(j) P_{x_{t-1}, j}(a)] - C(x_{t-1}, A_{t-1})$$

$$= g + h(x_{t-1}) - C(x_{t-1}, A_{t-1})$$

system

$$\Rightarrow 0 \leq \mathbb{E}_\pi \left[\sum_1^n (h(x_t) - h(x_{t-1}) - g + C(x_{t-1}, A_{t-1})) \right]$$

unity

$$\Rightarrow g \leq \frac{\mathbb{E}_\pi[h(x_n)]}{n} - \frac{\mathbb{E}_\pi[h(x_0)]}{n} + \frac{\mathbb{E}_\pi \left[\sum_1^n C(x_{t-1}, A_{t-1}) \right]}{n}$$

Letting n tend to infinity and noting that h is a bounded function

$$g \leq \Phi_\pi(x_0)$$

(5)

recall $\Phi_\pi(i) = \limsup_{n \rightarrow \infty} \frac{\mathbb{E}_\pi \left[\sum_0^n C(x_t, A_t) \mid x_0 = i \right]}{n+1}$

There exists only one place in the whole proof, where there was a inequality, hence if we choose the policy π^* , all inequalities turn into equalities and hence we have

$$g = \Phi_{\pi^*}(x_0)$$

St

α -optimal policy (optimal with discount factor α) "solves" (when)

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

Idea: Consider

$$\lim_{\alpha \uparrow 1} [C(i, a) + \alpha \sum P_{ij}(a) V_\alpha^*(j)]$$

An issue is that the limit may not exist

$h_\alpha(i) = V_\alpha^*(i) - V_\alpha^*(0)$ where 0 is some reference state

$$(1-\alpha)V_\alpha^*(0) + h_\alpha(i) = \min_a [C(i, a) + \alpha \sum P_{ij} h_\alpha(j)]$$

It is more likely that things will work out better using this format

The "h" are now called not value functions, but

DIFFERENTIAL VALUE FUNCTIONS

The hope is that $(1-\alpha)V_\alpha^*(0) \rightarrow g$ as $\alpha \rightarrow 1$

Thm X: If $\exists N < \infty$ s.t. $|V_\alpha^*(i) - V_\alpha^*(o)| < N \quad \forall \alpha, i$, then

- i) There exist g, h satisfying $(1) \rightarrow$ Bellman equation for time-avg MDP]
- ii) For some seq. $\alpha_n \rightarrow 1$,

$$h(i) = \lim_{n \rightarrow \infty} [V_{\alpha_n}^*(i) - V_{\alpha_n}^*(o)]$$
- iii) $\lim_{\alpha \rightarrow 1} (1 - \alpha) V_\alpha^*(o) = g$

Def

An MDP is unichain if it is irreducible under any stationary policy.

↓
it is possible to
reach every other
state from every state

Thm On a finite unichain MDP, hypothesis of Thm X holds.