

(s)

EE6106 LECTURE 20

Dated: 10th April 2024Discounted Average Cost MDPs (S, A, P, C, α) - V^* satisfies

$$V^*(i) = \min_a \left[C(i, a) + \alpha \sum_j P_{ij}(a) V^*(j) \right] \quad (1)$$

Step- Given stationary policy f , $T_f: B(S) \rightarrow B(S)$

$$T_f^n u \xrightarrow{n \uparrow \infty} V_f \quad (2)$$

- Stationary policy that minimizes RHS of (1) is optimal.

The minimization is over all policies.

↳ (1) can be thought of as a certain fixed point operation

In (2), u can be thought of somewhat like a termination cost

We are not in the learning setting right now, but in the "PLANNING" setting

Def A mapping $T: B(S) \rightarrow B(S)$ is a contraction mapping if, for $u, v \in B(S)$

$$\|Tu - Tv\| \leq \beta \|u - v\| \quad \text{for } \beta \in (0, 1)$$

$$\|u\| = \max_i |u(i)|$$

Contraction mapping Theorem:Say $T: B(S) \rightarrow B(S)$ is a contraction mapping. Then it has a unique fixed point. g i.e. $Tg = g$ Moreover, for any $u \in B(S)$, $T^n u \xrightarrow{n \uparrow \infty} g$

Lemma

$T_\alpha : B(s) \rightarrow B(s)$ is defined as

$$(T_\alpha u)(i) = \min_a [C(i,a) + \alpha \sum_j P_{ij}(a) u(j)]$$

is a contraction mapping

The interpretation of T_α is the minimum cost starting from state i .

Further, V^* is a fixed point of T_α

Our Bellman equations state that $T_\alpha V^* = V^*$

But since T_α is a contraction mapping, V^* is THE unique fixed point.

Algorithm

Initialize $V(1)$ arbitrarily

For $t \geq 1$

$$V(t+1) = T_\alpha V(t)$$

$$V(t) \xrightarrow{t \uparrow \infty} V^*$$

value iteration

Proof to be shared on Moodle

V^* solves the LP :

$$\max \sum_i V(i)$$

$$\text{s.t. } V(i) \leq C(i, a) + \alpha \sum_j P_{ij}(a) V(j) \quad \forall i \in S, a \in A$$

linear objective and linear constraints

The dual of this LP has some nice interpretations

For stationary policy π ,

$$\Phi_{\pi}(i, a) = C(i, a) + \alpha \sum_j P_{ij}(a) V_{\pi}(j)$$

Action value
function

Interpretation: Start at a state, take given action & then play the move as dictated by the policy.

$$V_{\pi}(i) = \Phi_{\pi}(i, \pi(i))$$

Similarly,

$$Q^*(i, a) \neq$$

Given stationary policy f ,

$$\hookrightarrow V_{f^*} = A_f$$

$$f^*(i) = \operatorname{argmin}_a [C(i, a) + \alpha \sum P_{ij}(a) V_f(j)]$$

$$= \operatorname{argmin}_a Q_{\pi}(i, a)$$

If we observe, there is some sort of a mapping between T_f and this quantity.

Lemma

$$V_{f^*}(i) \leq V_f(i) \quad \forall i$$

An iteration would give something which is strictly better

Proof

$$T_{f^*} V_f(i) = C(i, f^*(i)) + \alpha \sum P_{ij}(f^*(i)) V_f(j)$$

... by definition

Proof

$$T_{f^*} V_f(i) = C(i, f^*(i)) + \alpha \sum_j P_{ij}(f^*(i)) V_f(j)$$

... by definition

By construction, since the RHS is the minimum

$$T_{f^*} V_f(i) \leq C(i, f(i)) + \alpha \sum_j P_{ij}(f(i)) V_f(j) \\ = V_f(i)$$

$$\Rightarrow T_{f^*} V_f(i) \leq V_f(i)$$

$$\Rightarrow T_{f^*}^n V_f(i) \leq V_f(i)$$

POLICY ITERATION

At $t \geq 1$

- Compute V_{π_t}

- $\pi_{t+1} = \pi_t^*$

- If $\pi_{t+1} = \pi_t$, STOP

→ Policy Evaluation

→ Policy Improvement

REFERENCE :

Application: Selling an asset

- iid offers taking values $\{0, 1, \dots, N\}$

P_i = probability of offer i

- Daily maintenance cost C

- Discount future cost by α

- Once sold, no cost going forward

$$S = \{0, 1, 2, \dots, N\} \cup \{\infty\}$$

$$A = \{0, 1\} \text{ OR } \{S_{\text{SELL}}, W_{\text{WAIT}}\}$$

$$V^*(i) = \min \left\{ -i, C + \alpha \sum_{j=0}^N P_j V^*(j) \right\}$$

Note that the second term is not dependent on the state we are in.

$$\text{Let } i^* = \min \{ i : -i < C + \alpha \sum P_j V^*(j) \}$$

The optimal policy is then: sell if the bid is greater than i^* , else hold.

Note that we first need to obtain V^* . The way to do this is use the fact that we know the structure of this policy, use it to sweep over the CLASS of all such policies, find their value functions, and then optimize. This is now a 1-D optimization problem.

According to JKN, this is also a method to solve ARRANGED MARRIAGES

$f_i \sim$ policy of accepting offers $\geq i$

one of f_i is optimal.

The bids are i.i.d. Hence, the number of days before we accept the offer is a Geometrically distributed variable.

$T = \#$ of considered offers $\sim \text{Geometric} \left(\sum_i P_i \right)$

Given T , conditional expected discounted cost

$$= C + \alpha C + \dots + \alpha^{T-2} C - \alpha^{T-1} \left(\frac{\sum_{j=i}^N j P_j}{\sum_{j=1}^N P_j} \right)$$

$$= C \left(\frac{1 - \alpha^{T-1}}{1 - \alpha} \right) - \alpha^{T-1} \left(\frac{\sum_{j=1}^N j P_j}{\sum_{j=1}^N P_j} \right)$$

T is geometric, so we need to compute the following

Exercise: $E[\alpha^{T-1}] = \frac{p}{1 - \alpha(p-1)}$ where $p = \sum_{i=1}^N p_i$

The expected discounted cost under f_i is the following

$$\frac{c \sum_{j=0}^{i-1} p_j - \sum_{j=i}^N j p_j}{1 - \alpha \sum_{j=0}^{i-1} p_j}$$