

Lecture 2

1. Recap of definitions
MDPs, T_{μ}^{π} , E_{μ}^{π} , v^{π} , v^*
2. Fundamental theorem:
ML policies & v^* & BOE & DP
3. Value it. Complexity ϵ -opt
4. Policy it. $\tilde{O}\left(\frac{SA}{1-\gamma}\right)$ complexity
5. Lower bounds
6. Conclusions

MDP

$$M = (S, A, P, r, \gamma)$$

$$P = (P_a(s))_{s,a}$$

$$r = (r_a(s))_{s,a}$$

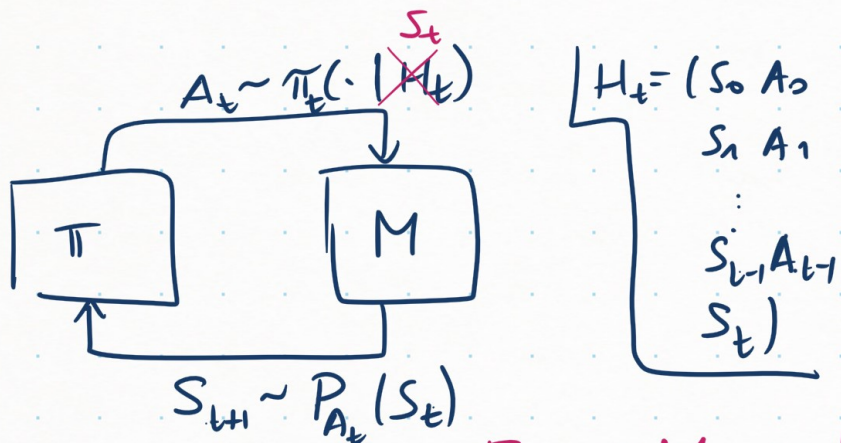
$$0 \leq \gamma < 1$$

Policy $\pi = (\pi_t)_{t \geq 0}$

$$\pi_t: (S \times A)^{t-1} \times S \rightarrow \mathcal{M}_1(A)$$

(discrete S, A)

$$|S|, |A| < \infty$$



Ionescu-Tulcea theorem

$P^\pi \leftarrow$ policy
 $\mu \leftarrow$ initial state distr.

E_μ^π

$$R = \sum_{t \geq 0} \gamma^t r_{A_t}(s_t)$$

$$\mu = \delta_s : P_{\delta_s}^\pi = P_s^\pi / E_s^\pi$$

$$v^\pi(s) = E_s^\pi[R]$$

$$v^*(s) = \sup_{\pi} \max_{\pi} v^\pi(s)$$

optimal value function $[0, 1)$

Prop: $\exists (\Omega, \mathcal{F}, P)$ \rightarrow Fix $\pi, M, \mu \in \mathcal{M}_1(S)$

$\exists s_0, a_0, s_1, a_1, \dots$

random elements over (Ω, \mathcal{F}, P)

s.t. **Markov property**

(1) $P(S_0 = s) = \mu(s) \quad \forall s \in S$

(b) $P(A_t = a | H_t) = \pi_t(a | H_t) \quad \forall a \quad \forall t \geq 0$

(c) $P(\underline{s_{t+1}} = s | \underline{H_t}, \underline{A_t}) = P_{A_t}(s_{t+1}, s) \quad \forall s \in S \quad \forall t \geq 0$

Goal: $\underline{\pi} \ ? \ v^\pi = v^*$

$$v^\pi \leq v^* \quad v^\pi \geq v^* ?$$

$$\Rightarrow v^\pi = v^*$$

Find π s.t. $v^\pi \geq \underline{v^* - \epsilon 1}$

$$1: S \rightarrow \mathbb{R}$$

$$s1 \rightarrow 1$$

ϵ -optimal policy

memoryless policies (ML)

$$\pi: S \rightarrow M_1(A)$$

$$|S|, |A| < \infty \quad \boxed{\text{FT}}$$

Theorem: (a) $\forall \pi$ greedy policy

w.r.t. $v^* \Rightarrow v^\pi = v^*$.

(b) $v^* = \underline{T} v^*$ Bellman optimality equation

Def:

π ML policy is greedy

w.r.t. $v: S \rightarrow \mathbb{R}$ if

the following holds: $\forall s \in S$

$$r^\pi(s) \quad p^\pi$$

$$\rightarrow \sum_a \pi(a|s) \{ r_a(s) + \gamma \langle P_a(s), v \rangle \}$$

$$= \max_a \{ r_a(s) + \gamma \langle P_a(s), v \rangle \}$$

$$r^\pi(s) = \sum \pi(a|s) r_a(s)$$

$$r^\pi \in \mathbb{R}^S \quad S = \{1, \dots, S\}$$

$$P^\pi(s, s') = \sum_a \pi(a|s) P_a(s, s')$$

$$(P_{ss'}^\pi)_{s, s'} \in [0, 1]^{S \times S}$$

$$T^\pi: \mathbb{R}^S \rightarrow \mathbb{R}^S$$

$$T^\pi(v) = r^\pi + \gamma P^\pi v$$

$$T: \mathbb{R}^S \rightarrow \mathbb{R}^S$$

$$(T(v))(s) = \max_a r_a(s) + \gamma \langle P_a(s), v \rangle$$

$$T^\pi(v) = T^\pi v$$

$$T(v) = Tv$$

Operators

$$\boxed{T v}(s) = (T(v))(s)$$

$$(T v)_s$$

$$\boxed{\begin{array}{l} \pi \text{ is greedy wrt } v \\ T^\pi v = Tv \end{array}}$$

T: Bellman operator

T^π : Policy eval. op. of π

Proof of the FT

$$\tilde{v}^*(s) = \sup_{\pi \in \Pi} v^\pi(s)$$

$$\left[\begin{array}{l} \pi \text{ is greedy wrt } \tilde{v}^* \\ \Rightarrow v^\pi = \tilde{v}^* \\ \Rightarrow \tilde{v}^* = T \tilde{v}^* \end{array} \right. \text{fixed-point}$$

$$\left[\tilde{v}^* = v^* \right]$$

$$\tilde{v}^* \leq v^* \quad \checkmark$$

$$v^* \leq \tilde{v}^* \quad ??$$

$$\mathbb{I}(\text{Log}) = \begin{cases} 1, & \text{Log} = \text{True} \\ 0, & \text{Log} = \text{False} \end{cases}$$

(Discounted) occupancy

measures

$$\mu \in \mathcal{M}_\eta(S) \quad , \quad \pi \text{ policy}$$

$$\Rightarrow \nu_\mu^\pi \in \mathcal{M}_1(S \times A) :$$

$$\nu_\mu^\pi(s, a) = \sum_{t \geq 0} \gamma^t \mathbb{P}_\mu^\pi(S_t = s, A_t = a)$$

$$\begin{aligned} v^\pi(s) &= \sum_{s, a} \nu_\mu^\pi(s, a) r_a(s) \\ &= \langle \nu_\mu^\pi, r \rangle \end{aligned}$$

$$\begin{aligned} &= \mathbb{E}_\mu^\pi \left[\sum_{t \geq 0} \gamma^t r_{A_t}(S_t) \right] = \sum_{s, a} \sum_{t \geq 0} \gamma^t \mathbb{E}_\mu^\pi \left[r_{A_t}(S_t) \mathbb{I}(S_t = s, A_t = a) \right] \end{aligned}$$

$\nearrow r_a(s)$

$$r_{A_t}(S_t) \mathbb{I}(S_t = s, A_t = a)$$

$$= r_a(s) \mathbb{I}(S_t = s, A_t = a)$$

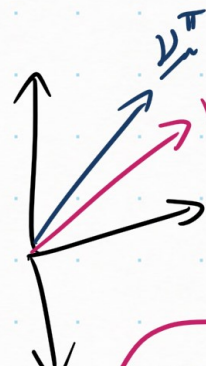
$$v^\pi(s) = \sum_{s, a} \sum_{t \geq 0} \gamma^t r_a(s) P_\mu^\pi(S_t = s, A_t = a)$$

$$= \sum_{s, a} r_a(s) \underbrace{\sum_{t \geq 0} \gamma^t P_\mu^\pi(S_t = s, A_t = a)}_{V_\mu^\pi(s, a)}$$

$$= \langle V_\mu^\pi, r \rangle$$

$V_\mu^\pi(s, a) \xrightarrow{\delta_s \rightarrow V_s^\pi}$

$$\tilde{V}_\mu^\pi(s) = \sum_a V_\mu^\pi(s, a)$$



$$\pi \text{ arb.} \rightarrow ML(\pi)$$

$$v^\pi(s) = \langle v_s^\pi, r \rangle$$

$$= \langle v_s^{ML(\pi)}, r \rangle / \sup_{\pi} \langle v_s^\pi, r \rangle$$

$$v^*(s) = \sup_{\pi} \langle v_s^\pi, r \rangle \leq \sup_{\pi \in ML} \langle v_s^\pi, r \rangle = \tilde{v}^*(s)$$

Theorem: $\forall \pi, \forall \mu \in M_n(S)$

$\exists \pi' \in ML$ s.t.

$$V_\mu^\pi = V_\mu^{\pi'}$$

Proof (hint):

$$\pi'(a|s) = \begin{cases} \frac{V_\mu^\pi(s, a)}{\tilde{V}_\mu^\pi(s)}, & \tilde{V}_\mu^\pi(s) \neq 0 \\ \pi_0(a), & \text{otherwise} \end{cases}$$

$\pi_0 \in M_n(A)$ arbitrary.



$\Theta(A^S)$

FT : v^* known

$$T_{\pi} v^* = T v^* \quad O(1)$$

$$\text{argmax}_a \{ r_a(s) + \gamma \langle P_a(s), v^* \rangle \}$$

$$\downarrow$$

$$O(A)$$

$$\underbrace{\hspace{10em}}_{O(S)}$$

① Compute v^*

② Find greedy w.r.t v^*
 $O(S^2 A)$

value iteration

$$T v^* = v^*$$

$$v_0 = 0$$

$$v_{k+1} = T v_k$$

$$\downarrow$$

$$v^* \quad k \rightarrow \infty$$

Prop:

$$\|T v - T u\|_{\infty} \leq \gamma \|v - u\|_{\infty}$$

(γ -contraction)

Banach's fixed point thm

Thm: $\|v_k - v^*\|_{\infty} \leq \gamma^k \|v^*\|_{\infty} \leq \epsilon$

$$\|v^*\|_{\infty} \leq \frac{1}{1-\gamma}$$

$\gamma \in [0, 1)$, $\epsilon \rightarrow 0$
 $r_a(s) \in [0, 1]$