

goal: determine dynamics of finite MP/SDE

refs:

Neuro-Dynamic Programming

Dimitri P. Bertsekas and John N. Tsitsiklis

chapter 2

Reinforcement Learning

An Introduction

Richard S. Sutton and Andrew G. Barto

chapter 3

• consider a finite MP/SDE (X, \mathcal{U}, P) , i.e. $|X|, |\mathcal{U}| < \infty$

$$x^+ \sim P(x^+ | x, u), \quad P: X \times \mathcal{U} \rightarrow \Delta(X)$$

→ starting from state $p \in \Delta(X)$ and applying policy $\pi: X \rightarrow \Delta(\mathcal{U})$

determine next state $p^+ \in \Delta(X)$

(it can help to regard $p: X \rightarrow [0, 1]$ — determine $p^+(x^+)$, all $x^+ \in X$)

— we can compute $p^+(x^+)$ for each $x^+ \in X$ using rules of probability:

$$p^+(x^+) = \sum_{x \in X} \left(p(x) \left[\sum_{u \in \mathcal{U}} \pi(u|x) P(x^+ | x, u) \right] \right)$$

\uparrow conditional probability $= [\pi(x)](u)$

* this determines a deterministic (!) difference equation:

* this determines a deterministic (!) difference equation:

$$p^+ = F(p), \quad \left(\underbrace{p: X \rightarrow [0,1]}_{p \in [0,1]^X} \right) \in \Delta(X) \subset [0,1]^X \subset \mathbb{R}^X = \mathbb{R}^N$$

$$\Rightarrow p \in \mathbb{R}^X = \mathbb{R}^N \quad \text{— more generally we let } B^A = \{f: A \rightarrow B\}$$

$$\text{e.g. } \mathbb{R}^n = \{v: \{1, \dots, n\} \rightarrow \mathbb{R}\}$$

→ show that F is linear (!)

(i.e. find $\Gamma \in \mathbb{R}^{N \times N}$, $N = |X|$, s.t. $F(p) = \Gamma p = p^+$)

$$- [\Gamma]_{x^+, x} = \sum_{u \in \mathcal{U}} \pi(u|x) P(x^+|x, u) \text{ yields } p^+ = \Gamma p$$

* we can use linear systems theory to characterize "solutions" / trajectories including their asymptotic properties !

→ we're studying discrete-time linear time-invariant DE $p^+ = \Gamma p$

• first of all, we know trajectories $p_t = \Gamma^t p_0$, any $t \in \mathbb{N}$
 \uparrow t -fold matrix multiplication

• furthermore, Γ has special properties: $\mathbf{1}^T \Gamma = \mathbf{1}^T \Rightarrow \mathbf{1} \in \text{spec } \Gamma$
 $\uparrow \mathbf{1}^T = (1, \dots, 1) \in \mathbb{R}^N$

i.e. Γ is "left-stochastic" $\Rightarrow \forall \lambda \in \text{spec } \Gamma: |\lambda| \leq 1$

* if $\underbrace{[\Gamma]_{x^+, x}}_{> 0}$ then $\lim_{t \rightarrow \infty} p_t = \bar{p}$ where $\Gamma \bar{p} = \bar{p}$ is unique

\uparrow more generally, if Γ is irreducible and aperiodic)

so all trajectories asymptotically converge to unique \bar{p} !