

# Reinforcement Learning Markov Decision Process Notes

Lu Hong <sup>1</sup>

<sup>1</sup> Nanjing University of Aeronautics and Astronautics

September 21, 2019

## 1 Markov Process

"The future is independent of the past given the present". A state  $S_t$  is *Markov* iif.

$$P[S_{t+1}|S_t] = P[S_{t+1}|S_1, S_2, \dots, S_t]$$

Markov Process (or Markov Chain) is a *memoryless random process*, represented by a tuple  $\langle S, P \rangle$

## 2 Markov Reward Process

Markov Reward Process is Markov Process with values, represented by a tuple  $\langle S, P, R, \gamma \rangle$ , where  $R = E[R_{t+1}|S_t = s]$  and  $\gamma$  is a discount factor.

Return  $G_t$  is the *total discounted reward* at time-step  $t$  presented by

$$G_t = \sum_{i=0}^{\infty} \gamma^i R_{t+i+1}$$

Also, we define Value Function to indicate the long-term value of the state  $s$ .

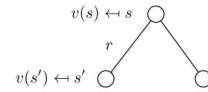
$$V(s) = E[G_t|S = s]$$

\* Bellman Equation for MRPs, it demonstrate that MRPs can be presented in recursive format.

$$\begin{aligned} v(s) &= E[G_t|S_t = s] \\ &= E[R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots | S_t = s] \\ &= E[R_{t+1} + \gamma(R_{t+2} + \gamma R_{t+3} + \dots) | S_t = s] \\ &= E[R_{t+1} + \gamma G_{t+1} | S_t = s] \\ &= E[R_{t+1} + \gamma v(S_{t+1}) | S_t = s] \end{aligned}$$

### Bellman Equation for MRPs (2)

$$v(s) = E[R_{t+1} + \gamma v(S_{t+1}) | S_t = s]$$



$$v(s) = R_s + \gamma \sum_{s' \in S} P_{ss'} v(s')$$

To express Bellman Equation in matrix form as

$$v = R + \gamma P v,$$

we can solve Bellman Equation easily as linear equation.

$$\begin{aligned} v &= R + \gamma P v \\ (I - \gamma P)v &= R \\ v &= (I - \gamma P)^{-1} R \end{aligned}$$

•

## 3 Markov Decision Process

A MDP is a MRP with decisions, represented by a tuple  $\langle S, A, P, R, \gamma \rangle$ , to be specific, some params have been changed after actions added.  $P_{ss'}^a = P[S_{t+1} = s' | S_t = s, A_t = a]$  and  $R_s^a = E[R_{t+1} | S_t = s, A_t = a]$

\*Policies: A policy  $\pi$  is a distribution over actions given states,

$$\pi(a|s) = P[A_t = a | S_t = s]$$

## Policies (2)

- Given an MDP  $\mathcal{M} = \langle \mathcal{S}, \mathcal{A}, \mathcal{P}, \mathcal{R}, \gamma \rangle$  and a policy  $\pi$
- The state sequence  $S_1, S_2, \dots$  is a Markov process  $\langle \mathcal{S}, \mathcal{P}^\pi \rangle$
- The state and reward sequence  $S_1, R_2, S_2, \dots$  is a Markov reward process  $\langle \mathcal{S}, \mathcal{P}^\pi, \mathcal{R}^\pi, \gamma \rangle$
- where

$$\mathcal{P}_{s,s'}^\pi = \sum_{a \in \mathcal{A}} \pi(a|s) \mathcal{P}_{ss'}^a$$
$$\mathcal{R}_s^\pi = \sum_{a \in \mathcal{A}} \pi(a|s) \mathcal{R}_s^a$$