

## Markov Decision Process

The 1<sup>st</sup> step to apply RL is to formulate the problem as a **MDP**

Starting from a Markov Process:

- if we add  $R$  we get a MARKOV REWARD PROCESS
- if we add  $A$  we get a MARKOV DECISION PROCESS

Why MDP?

1. Actions influence **present and future** rewards
2. MDP involves **delayed** rewards
3. With MDP we introduce the concept of **STATE** (not present in bandit)

MDP describes an environment for RL FULLY OBSERVABLE

STATE: A state  $S_t$  is Markov  $\Leftrightarrow P[S_{t+1} | S_t] = P[S_{t+1} | S_1, \dots, S_t]$

STATE TRANSITION PROBABILITY:  $P_{ss'} = P[S_{t+1} = s' | S_t = s]$   
 $\downarrow$   
the sum of each row is 1

1

MARKOV PROCESS:  $\gamma_t$  is a tuple  $\langle S, P \rangle$

- $S$ : finite set of states
- $P$ : state transition probability matrix

2

MARKOV REWARD PROCESS: tuple  $\langle S, P, R, \gamma \rangle$

- $R$ : reward function  $R_s = E[R_{t+1} | S_t = s]$
- $\gamma$ : discount factor  $[0, 1]$

GOAL: MAXIMIZE TOTAL REWARD ↘ present value of future rewards

RETURN:  $G_t = R_{t+1} + \gamma R_{t+2} + \dots = \sum_{k=0}^{\infty} \gamma^k R_{t+1+k}$

VALUE FUNCTION:  $V_s = E[G_t | S_t = s]$

↳ expected return, starting from  $s$

## BELLMAN EQUATION FOR MRP

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

This value func. can be divided in 2 parts:

- immediate reward
- discounted value of successor state

- Bellman Equation express a relationship between the value of a state and the values of successor states
- BE averages over all the possibilities (SL. 24)
- CC:  $O(n^3)$  → very expensive

3

MARKOV DECISION PROCESS is a tuple  $\langle S, A, P, R, \gamma \rangle$

- $A$ : set of action
- $P_{ss'}^a = P[S_{t+1} = s' | S_t = s, A_t = a]$  → matrix  $\forall$  action
- $R_s^a = E[R_{t+1} | A_t = a, S_t = s]$

How do we make decisions? POLICIES

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

- ↳ distribution over actions given states
- ↳ defines the behaviour of an agent

**VALUE FUNC FOR MDP:**  $V_{\pi}(s) = E_{\pi}[G_t | S_t = s]$

- ↳ Same as value func but it follows the policy  $\pi$

**ACTION-VALUE FUNC:**  $Q_{\pi}(a|s) = E_{\pi}[G_t | S_t = s, A_t = a]$

# BELLMAN EXPECTATION EQUATION

$$V_{\pi}(s) = \sum_{a \in A} \pi(a|s) q_{\pi}(s, a)$$

$$q_{\pi}(s, a) = R_s^a + \gamma \sum_{s' \in S} P_{s, s'}^a V_{\pi}(s')$$

The probability of taking one action or another one depends on the policy.

DIFFERENCE BETWEEN B.E. AND B.E.G.