

# Reinforcement Learning Notes

February 18, 2018

## 1 Markov Reward process

A MRP(Markove reward process) does not have actions involved, that concept is MDP(markov decision process)

## 2 return G

The return  $G_t$  is total discounted reward for time-step t. return is defined for a given sample

$$G_t = R_{t+1} + \gamma R_{t+2} + \dots = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \quad (1)$$

The discount  $\gamma \in [0, 1]$

## 3 Bellman equation for MRPs

The main idea is :

The value function can be decomposed into two parts:

- immediate reward  $R_{t+1}$
- discounted value of successor state  $\gamma v(S_{t+1})$

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

## 4 Formal definition of MDP

A Markov decision process is a markov reward process with decisions. It is an environment in which all states are markov.

A Markov decision Process is a tuple  $(S, A, P, R, \gamma)$

- S is a finite set of states
- A is a finite set of actions
- P is a state transition probability matrix,  
 $P_{ss'}^a = P[S_{t+1} = s' | S_t = s, A_t = a]$
- R is a reward function,  
 $R_s^a = E[R_{t+1} | S_t = s, A_t = a]$
- $\gamma$  is a discount factor

## 5 Policy and Stochastic policy

A policy defines behavior/decisions/actions of the agent to look for what actions to take. A policy tells what action to take given a state S i.e.  $\pi : s \mapsto a$

A stochastic policy  $\pi$ , is a distribution over actions given state,

$$\pi(a|s) = P[A_t = a | S_t = s] \quad (3)$$

## 6 Value function

The **state-value function**  $v_\pi(s)$  of an MDP is the expected return starting from state s, and then following policy  $\pi$

$$v_\pi(s) = E_\pi[G_t | S_t = s] \quad (4)$$

The expectation above makes sense for a stochastic policy, whereas for a fixed policy value function is just return defined by policy.

The **action-value function**  $q_\pi(s, a)$  is the expected return starting from state s, taking action a, and then following policy  $\pi$

$$q(s, a) = E_\pi[G_t | S_t = s, A_t = a] \quad (5)$$

Relating state-value  $v_\pi(s)$  and action-value  $q_\pi(s, a)$  (a single step look ahead):

$$v_\pi(s) = \sum_{a \in A} \pi(a|s) q_\pi(s, a) \quad (6)$$

Relating action-value  $q_\pi(s, a)$  and state-value  $v_\pi(s)$  (a single step look ahead):

$$q_\pi(s, a) = R_s^a + \sum_{s' \in S} \gamma P_{ss'}^a v_\pi(s') \quad (7)$$

Using the above two steps look ahead (first over all actions then over all subsequent states) we can specify state-values in terms of itself:

$$v_\pi(s) = \sum_{a \in A} \pi(a|s) \left( R_s^a + \sum_{s' \in S} \gamma P_{ss'}^a v_\pi(s') \right) \quad (8)$$

Similarly using two steps look ahead (first over all states then over all subsequent actions), we can specify action-values recursively in terms of itself TODO