# Reinforcement Learning Notes

February 7, 2018

## 1   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

$$G_t = R_{t+1} + \gamma R_{t+2} + \cdots = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \tag{1}$$

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 $\gamma v(S_{t+1})$

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