

Value Estimation

- 1 Known Environment
 - Bellman's Equation
 - Solving Techniques

- 2 Unknown Environment
 - Monte-Carlo Method

- 1 Known Environment
 - Bellman's Equation
 - Solving Techniques
- 2 Unknown Environment
 - Monte-Carlo Method

Model of the Environment

- Where does an action take us?

$$\delta(s, a) \mapsto s'$$

- How much reward do we receive?

$$r(s, a) \mapsto \mathfrak{R}$$

The values of different states are interrelated
Bellman's Equation:

$$V^\pi(s) = r(s, a) + \gamma \cdot V^\pi(\delta(s, a)) \quad \text{where } a = \pi(s)$$

Can we solve *Bellman's equation*?

$$V^\pi(s) = r(s, a) + \gamma \cdot V^\pi(\delta(s, a)) \quad \text{where } a = \pi(s)$$

- Direct solution (linear equation system)
- Iteratively (*value iteration*)

$$V_{k+1}^\pi(s) \leftarrow r(s, \pi(s)) + \gamma \cdot V_k^\pi(\delta(s, \pi(s)))$$

- 1 Known Environment
 - Bellman's Equation
 - Solving Techniques
- 2 Unknown Environment
 - Monte-Carlo Method

How do we find an **optimal policy** π^* ?

Easy if the optimal value function V^* was known:

$$\pi^*(s) = \arg \max_a (r(s, a) + \gamma \cdot V^*(\delta(s, a)))$$

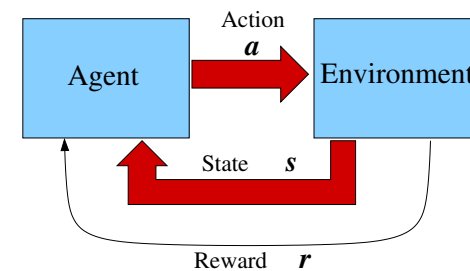
Optimal version of Bellman's equation

$$V^*(s) = \max_a (r(s, a) + \gamma \cdot V^*(\delta(s, a)))$$

Hard to solve

Policy iteration:

Interleaved calculation of policy and values



Normal scenario: *Unknown environment*
 $r(s, a)$ and $\delta(s, a)$ are not known

V^π must be estimated from **experience**

Monte-Carlo Method

- Start at a random s
- Follow π , store the rewards and s_t
- When the goal is reached, update $V^\pi(s)$ -estimation for all visited states with the future reward we actually received

Very slow convergence