

Reinforcement Learning

Part 1

Yingyu Liang

`yliang@cs.wisc.edu`

**Computer Sciences Department
University of Wisconsin, Madison**

Goals for the lecture

you should understand the following concepts

- the reinforcement learning task
- Markov decision process
- value functions
- value iteration

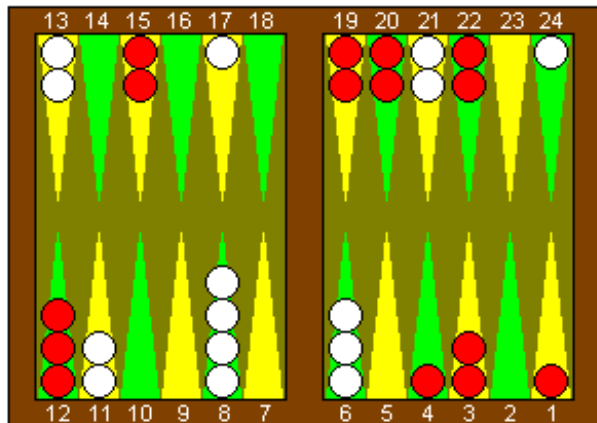
Reinforcement learning (RL)

Task of an agent embedded in an environment

repeat forever

- 1) sense world
- 2) reason
- 3) choose an action to perform
- 4) get feedback (usually reward = 0)
- 5) learn

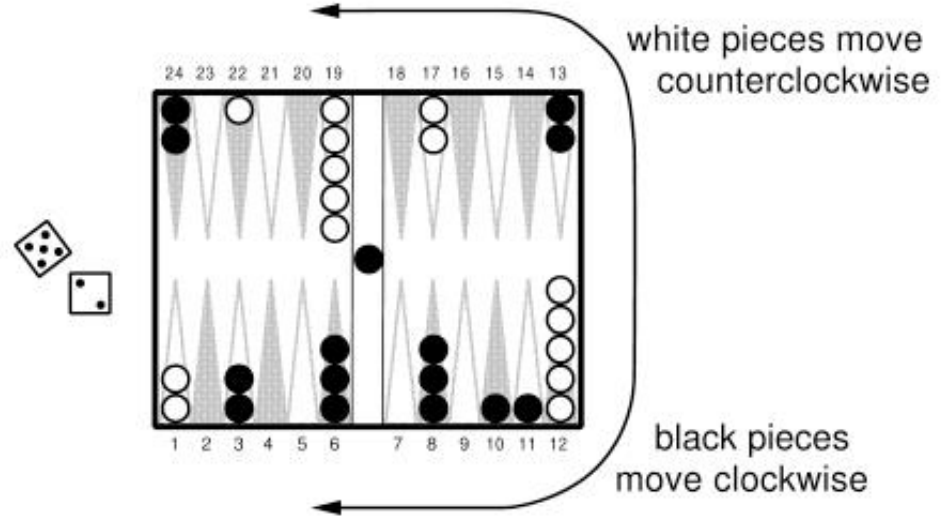
the environment may be the physical world or an artificial one



Example: RL Backgammon Player

[Tesauro, *CACM* 1995]

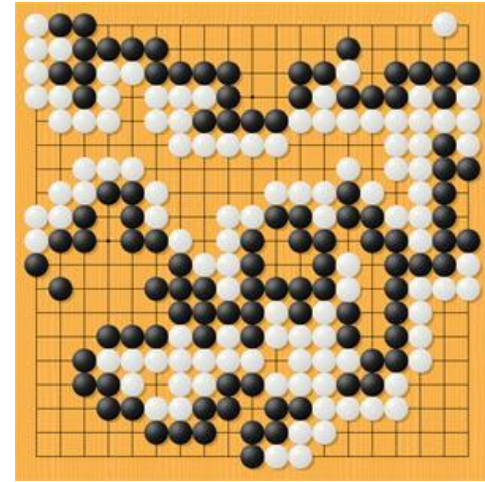
- world
 - 30 pieces, 24 locations
- actions
 - roll dice, e.g. 2, 5
 - move one piece 2
 - move one piece 5
- rewards
 - win, lose
- TD-Gammon 0.0
 - trained against itself (300,000 games)
 - as good as best previous BG computer program (also by Tesauro)
- TD-Gammon 2
 - beat human champion



Example: AlphaGo

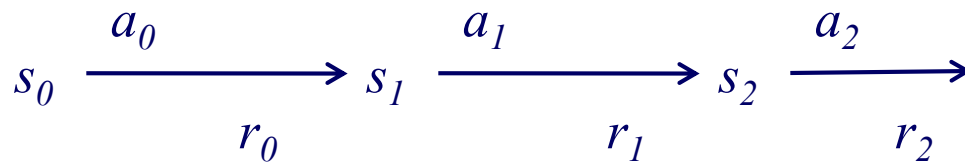
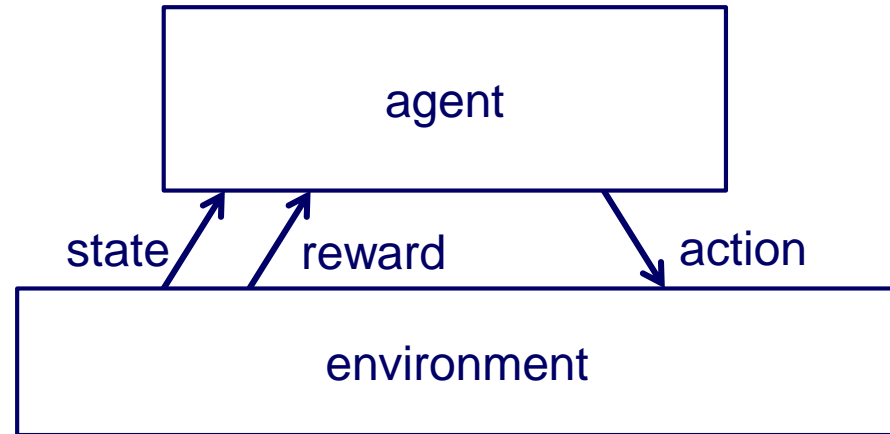
[Nature, 2017]

- world
 - 19x19 locations
- actions
 - Put one stone on some empty location
- rewards
 - win, lose
- 2016 beats World Champion Lee Sedol by 4-1
- Subsequent system (AlphaGo Master/zero) shows superior performance than humans
- Trained by supervised learning + reinforcement learning



Reinforcement learning

- set of states S
- set of actions A
- at each time t , agent observes state $s_t \in S$ then chooses action $a_t \in A$
- then receives reward r_t and changes to state s_{t+1}



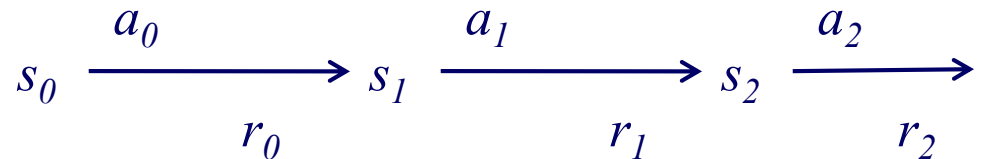
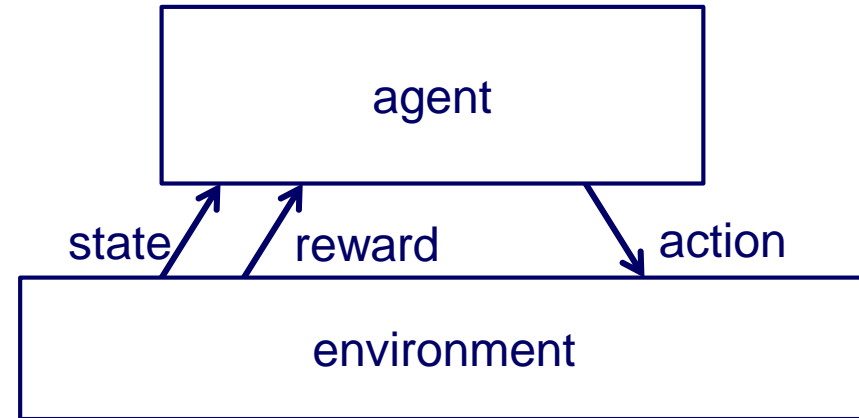
Reinforcement learning as a Markov decision process (MDP)

- Markov assumption

$$P(s_{t+1} | s_t, a_t, s_{t-1}, a_{t-1}, \dots) = P(s_{t+1} | s_t, a_t)$$

- also assume reward is Markovian

$$P(r_{t+1} | s_t, a_t, s_{t-1}, a_{t-1}, \dots) = P(r_{t+1} | s_t, a_t)$$



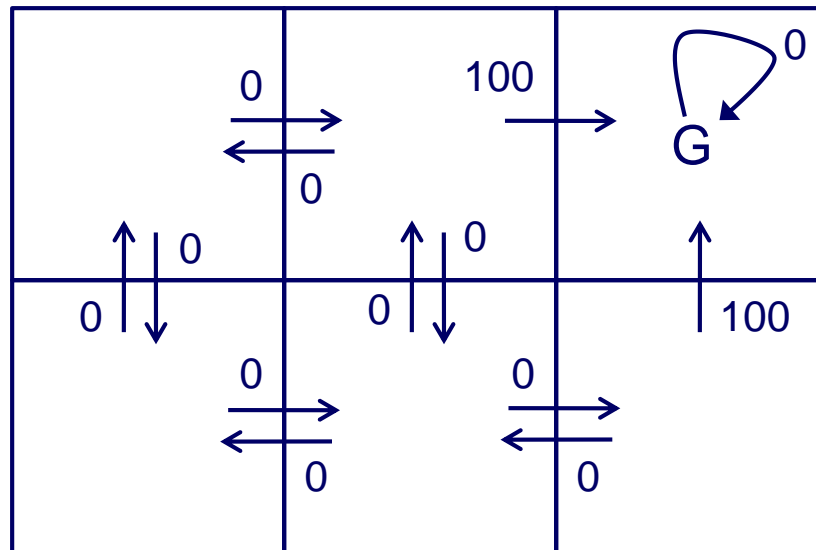
Goal: learn a policy $\pi : S \rightarrow A$ for choosing actions that maximizes

$$E[r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \dots] \quad \text{where } 0 \leq \gamma < 1$$

for every possible starting state s_0

Reinforcement learning task

- Suppose we want to learn a control policy $\pi : S \rightarrow A$ that maximizes $\sum_{t=0}^{\infty} \gamma^t E[r_t]$ from every state $s \in S$



each arrow represents an action a and the associated number represents deterministic reward $r(s, a)$

Value function for a policy

- given a policy $\pi : S \rightarrow A$ define

$$V^\pi(s) = \sum_{t=0}^{\infty} \gamma^t E[r_t]$$

assuming action sequence chosen according to π starting at state s

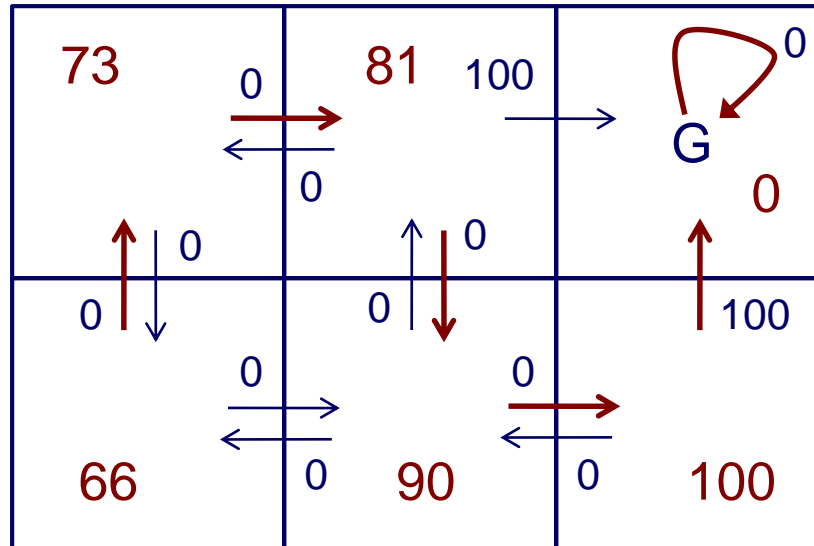
- we want the optimal policy π^* where

$$\rho^* = \arg \max_{\rho} V^\rho(s) \quad \text{for all } s$$

we'll denote the value function for this optimal policy as $V^*(s)$

Value function for a policy π

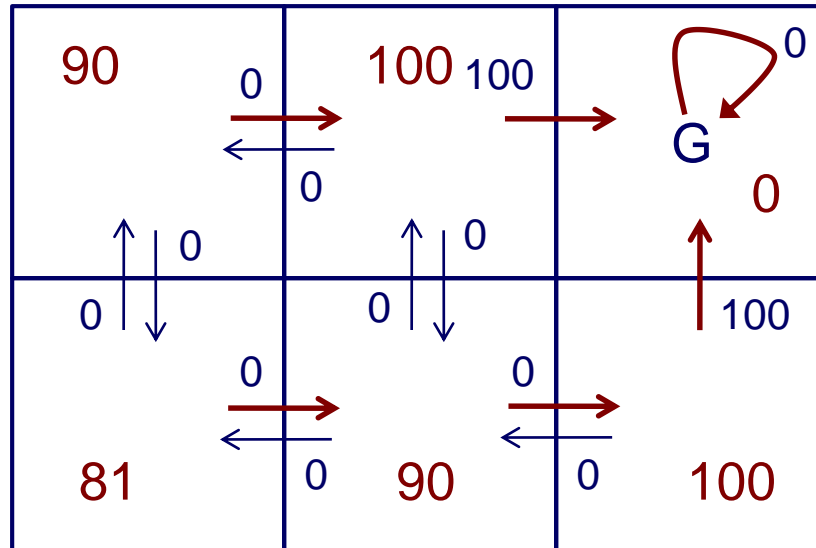
- Suppose π is shown by red arrows, $\gamma = 0.9$



$V^\pi(s)$ values are shown in red

Value function for an optimal policy π^*

- Suppose π^* is shown by red arrows, $\gamma = 0.9$



$V^*(s)$ values are shown in red

Using a value function

If we know $V^*(s)$, $r(s_t, a)$, and $P(s_{t+1} | s_t, a)$
we can compute $\pi^*(s)$

$$\pi^*(s_t) = \arg \max_{a \in A} \left[r(s_t, a) + \gamma \sum_{s \in S} P(s_{t+1} = s | s_t, a) V^*(s) \right]$$

Value iteration for learning $V^*(s)$

initialize $V(s)$ arbitrarily

loop until policy good enough

{

 loop for $s \in S$

 {

 loop for $a \in A$

 {

$$Q(s, a) \leftarrow r(s, a) + \gamma \sum_{s' \in S} P(s' | s, a) V(s')$$

 }

$$V(s) \leftarrow \max_a Q(s, a)$$

 }

}

Value iteration for learning $V^*(s)$

- $V(s)$ converges to $V^*(s)$
- works even if we randomly traverse environment instead of looping through each state and action methodically
 - but we must visit each state infinitely often
- implication: we can do online learning as an agent roams around its environment
- assumes we have a model of the world: i.e. know $P(s_t | s_{t-1}, a_{t-1})$
- What if we don't?