

# Reinforcement Learning

Final exam : Friday May 15th

- 7:15 PM - 10:15 PM
  - Online
  - Multiple choice
  - Calculator needed
  - Practice Gradiance Quiz will be posted
  - One more Gradiance Quiz left (RL)
- 

$A$  - set of all actions

$$\underline{A_t \in \underline{A(s_t)}} \subset A$$
$$R_t \in \mathbb{R}$$

---

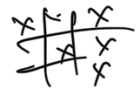
Policy  $\pi$

$$\boxed{s_1 \rightarrow A_{s_1}}$$

$$S_2 \rightarrow \underline{\underline{\cancel{AS_2}}}$$

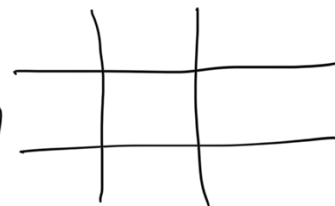
## Tic Tac Toe

Playing against an imperfect opponent.



Should  
be  
 $3^9$

$9^3$  "possible" states



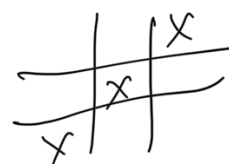
$S$

Value Prob. of winning from that State.

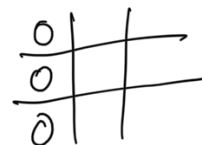
$S$  is small fact'a.



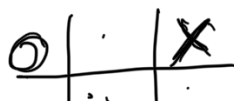
0.5



1



0



$A(S_0)$

~~x to 0~~

$s_1$

$$0.5 + \alpha(0 - 0.5)$$

$$0.5 - 0.5\alpha$$

Markov Decision processes (MDP)

Finite MDP



$$p(s', r | s, a) \triangleq P(S_t = s', R_t = r | S_{t-1} = s, A_{t-1} = a)$$

$$|S| = 2 = \text{on, off}$$

$$|A| = 4 = u, d, l, r$$

$$|R| = 2 = +1, -1$$

$$p(\text{on}, +1 | \text{on}, u)$$

$$\rightarrow p(\underline{s'}, \underline{r} | s, a)$$

$$\underline{p(s' | s, a)} \doteq \sum_{r \in R} p(\underline{s'}, \underline{r} | s, a)$$


---

$$G_t = \text{Expected return} \\ \text{from } t \rightarrow \underline{\underline{T}}$$

# Reinforcement learning - part II

## Announcements

- ① Gradience 10 due on Sunday night
  - ② For finals, a practice quiz and list of topics will be posted soon
  - ③ Assignment 3 due this week on Friday
- 

RL

↳ Markov Decision Processes (MDP)

$$p(s', r | s, a)$$

$$G_t = R_{t+1} + \dots$$

Agent

Policy ( $\pi$ )

$$\pi(a|s) \rightarrow A_t = a \mid S_t = s$$

value function at state  $s$

$$V_{\pi}(s) \doteq \mathbb{E}_{\pi}[\underline{G_t} | S_t = s] \quad \textcircled{\circ}$$

$$Q_{\pi}(s, a) \doteq \mathbb{E}_{\pi}[G_T | S_t = s, A_t = a] \quad \textcircled{\circ}$$

complete model of the environment.

$$\underline{p(s', r | s, a)}$$

Dynamic Programming <sup>(DP)</sup> based methods

init: Start with a policy

Step 1. Find  $\underline{V_{\pi}}(s) \forall s \in S$

Policy prediction

Step 2  $\pi \rightarrow V_{\pi}(s) \rightarrow \pi'$

Policy Improvement

$$\underline{p(s', r | s, a)}$$

DP  $\rightarrow$  works if we ~~have~~ a perfect  
 know the environment model  
 completely  $\underbrace{p(s', r | s, a)}$   
 $\rightarrow$  Data-free  
 $\rightarrow$  not very useful

Monte-carlo methods.

Policy prediction  $\swarrow$   
 Policy Control  $\nwarrow$

first visit MC prediction

$s_1$        $V(s_1) = 1$       Return<sub>1</sub>( $s_1$ ) = [.]  
 $s_2$        $V(s_2) = 1$       Return<sub>1</sub>( $s_2$ ) = [.]

$(s_1) A_0 R_1 s_1 A_1 R_2 s_2 A_2 R_3 s_3 A \dots$

|               |
|---------------|
| $V(s_t)$      |
| $Q(s_t, a_t)$ |

PlayerTotal , Ace usable , House Card

|       |               |   |   |
|-------|---------------|---|---|
| $s_1$ | <del>11</del> | 0 | 6 |
| $s_2$ | 13            | 0 | 6 |
| $s_3$ | 17            | 0 | 6 |
| $s_4$ | 18            | 0 | 6 |
| $s_5$ | 22            | 0 | 6 |

-1

$$\text{Returns}(22, 0, 6) = \underline{\underline{[-1]}}$$

✓

$$\text{Returns}(18, 0, 6) = \underline{\underline{[-1]}}$$

⋮

$$V(\_) = -1,$$

$$\left[ \begin{array}{ccc} 11 & 0 & 6 \end{array} \right]$$

{