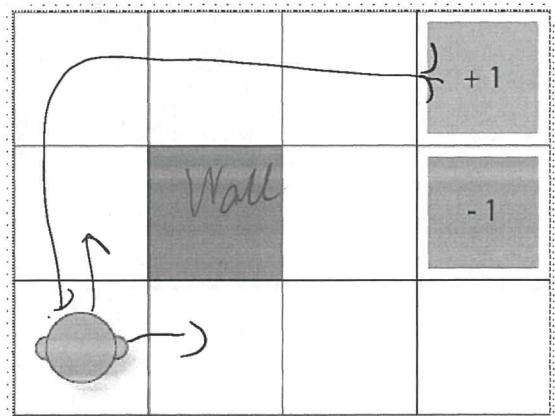


# Markov Decision Processes

Tuesday, 4 September 2018 10:29 AM

# Lecture 8



The grey square is a wall.

The two labelled cells give a reward: 1 and -1 respectively.

But! Things can go wrong:

- If the agent tries to move north 80% of the time, this works as planned (provided the wall is not in the way)
- 10% of the time, trying to move north takes the agent west (provided the wall is not in the way)
- 10% of the time, trying to move north takes the agent east (provided the wall is not in the way)
- If the wall is in the way of the cell that would have been taken, the agent stays put.
- Similar for all other directions

## Classical Planning:

- Set of states S
- Initial state I
- Transition function A
- Goals G

## MDPs:

- Set of states S
- Initial state I
- Transition probabilities:
  - $P_a(s \rightarrow s')$
- Reward function  $r(s, a, s')$  in real numbers
- ← Discount factor (gamma)  $\gamma \in (0, 1)$
- future

## Discounted reward:

$$V = r_1 + \gamma r_2 + \gamma^2 r_3 + \dots$$

$$= r_1 + \gamma(r_2 + \gamma r_3 + \gamma^2 r_4 + \dots)$$

$$V_t = \frac{r_1 + \gamma V_{t+1}}{\gamma(1 - \gamma)}$$

$$V_t = r_t + \gamma V_{t+1}$$

## Probabilistic PDDL:

```
(define (domain bomb-and-toilet)
  (:requirements :conditional-effects :probabilistic-effects)
  (:predicates (bomb-in-package ?pkg) (toilet-clogged)
   (bomb-defused))

  (:action dunk-package
   :parameters (?pkg)
   :effect (and (when (bomb-in-package ?pkg)
                      (bomb-defused))
    (probabilistic 0.05 (toilet-clogged))))
```

deterministic policy  
Solution for MDP is a policy:

function  
at(0,0) => move\_right  
at(0,1) => move\_right  
at(0,2) => move\_right  
at(0,3) => stay  
at(1,0) => move\_up  
at(1,2) => move\_up  
at(1,3) => move\_up  
at(2,0) => move\_up  
at(2,1) => move\_left  
at(2,2) => move\_up  
at(2,3) => move\_left

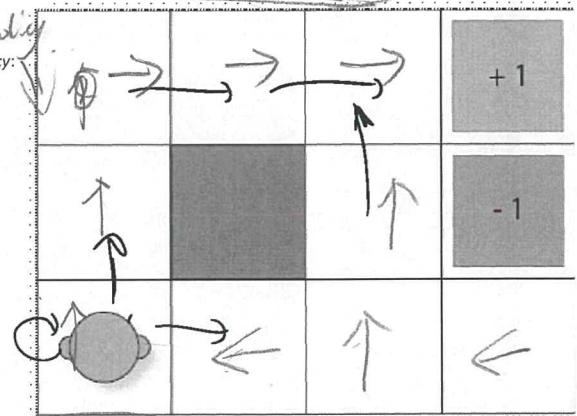
An example  
not best

$\pi(s) \rightarrow a$   
Policy

$\pi(s)$

deterministic

we focus on



Stochastic  $\pi(s, a)$  maximum to get

## Expected return exercise:

You can steal:

- A iPhone, which you think you have a 20% chance of selling for \$500, or an 80% chance of selling for \$250.
- B) A Samsung, which you think you have a 50% chance of selling for \$500, or a 50% chance of selling for \$200.

Which do you steal? <https://pollev.com/timothymille936>

$$A: 0.2 \cdot 500 + 0.8 \cdot 250 = 300$$

$$B: 0.5 \cdot 500 + 0.5 \cdot 200 = 350$$

Expected return of

reward  $\rightarrow s'$

Bellman equations:

$$V(s) = \max_{a \in A} \left( \sum_{s' \in S} P_a(s'|s) \left[ r(s, a, s') + \gamma V(s') \right] \right)$$

immediate reward      reward function      discounted future reward



$$\sum_{a \in A(s)} P_a(s'|s) [r(s, a, s') + \gamma V(s')]$$

prob.

(immediate reward  
+ discounted reward)



# Value iteration and policy iteration

Tuesday, 4 September 2018 9:34 AM

## Lecture 9

$$V(s) = \max_{a \in A(s)} \sum_{s' \in S} P_a(s'|s)[r(s, a, s') + \gamma V(s')]$$

Value iteration:

1) Set  $V_0$  to arbitrary value for each  $s$  in  $S$  (choose 0 as the value)

2) While diff is  $\geq \epsilon$ :

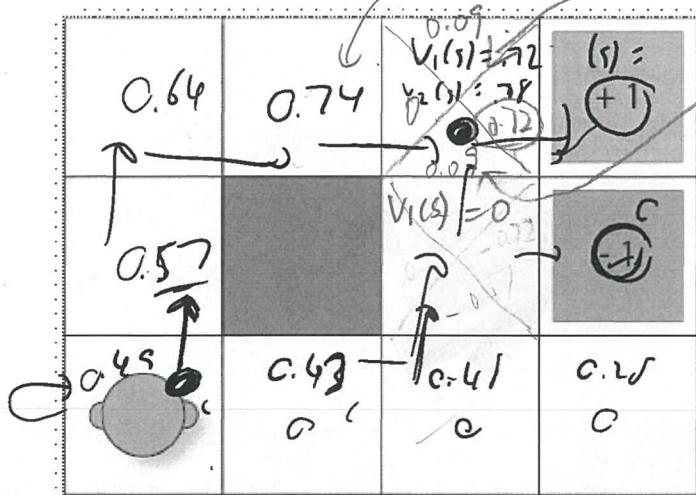
a. For each  $s$  in  $S$  do

$$\text{i. } V_{t+1}(s) := \max_{a \in A(s)} \sum_{s' \in S} P_a(s'|s)[r(s, a, s') + \gamma V_t(s')]$$

3) Select policy

$$\begin{aligned} & 0.8(0 + 0.9 \times 1) = 0.72 \\ & 0.1 \times (0 + 0.9 \times 0) = 0 \cdot \text{right but} \\ & 0.1 \times (0 + 0.9 \times 0) = 0. \end{aligned}$$

$V(s)$



The two labelled cells give a reward: 1 and -1 respectively. (Actually, we will assume  $V(s)=1$  or -1)

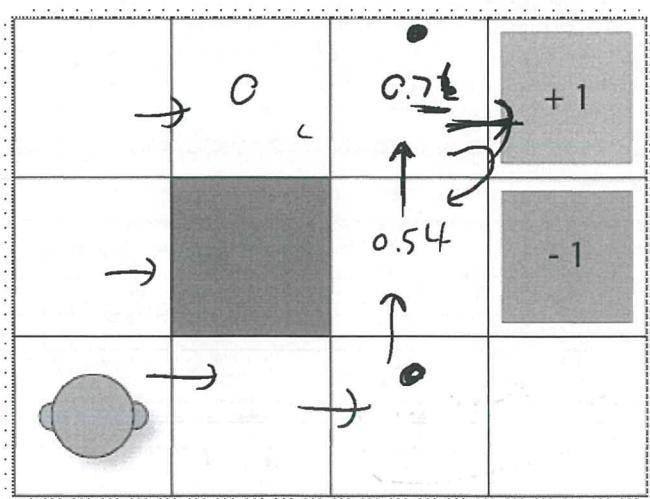
But! Things can go wrong:

- If the agent tries to move north, 80% of the time, this works as planned (provided the wall is not in the way)
- 10% of the time, trying to move north takes the agent west (provided the wall is not in the way)
- 10% of the time, trying to move north takes the agent east (provided the wall is not in the way)
- If the wall is in the way of the cell that would have been taken, the agent stays put.
- Similar for all other directions

$$\begin{aligned} & \text{1st iteration} \\ & \text{if } 0.41 \neq \max \\ & \Rightarrow 0.72 \end{aligned}$$

Policy iteration:

$V^\pi(s)$



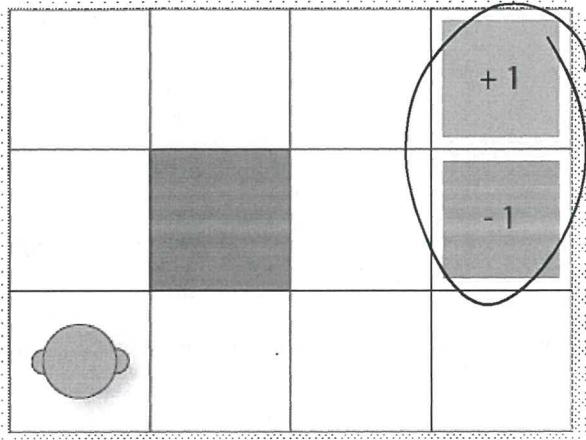
$$Q^\pi(a, s) = \sum_{s' \in S} P_a(s'|s)[r(s, a, s') + \gamma V^\pi(s')]$$

$$V^\pi(s) = \arg \max_{a \in A(s)} Q^\pi(a, s)$$



## Monte-Carlo Tree Search

Tuesday, 11 September 2018 13:58 PM

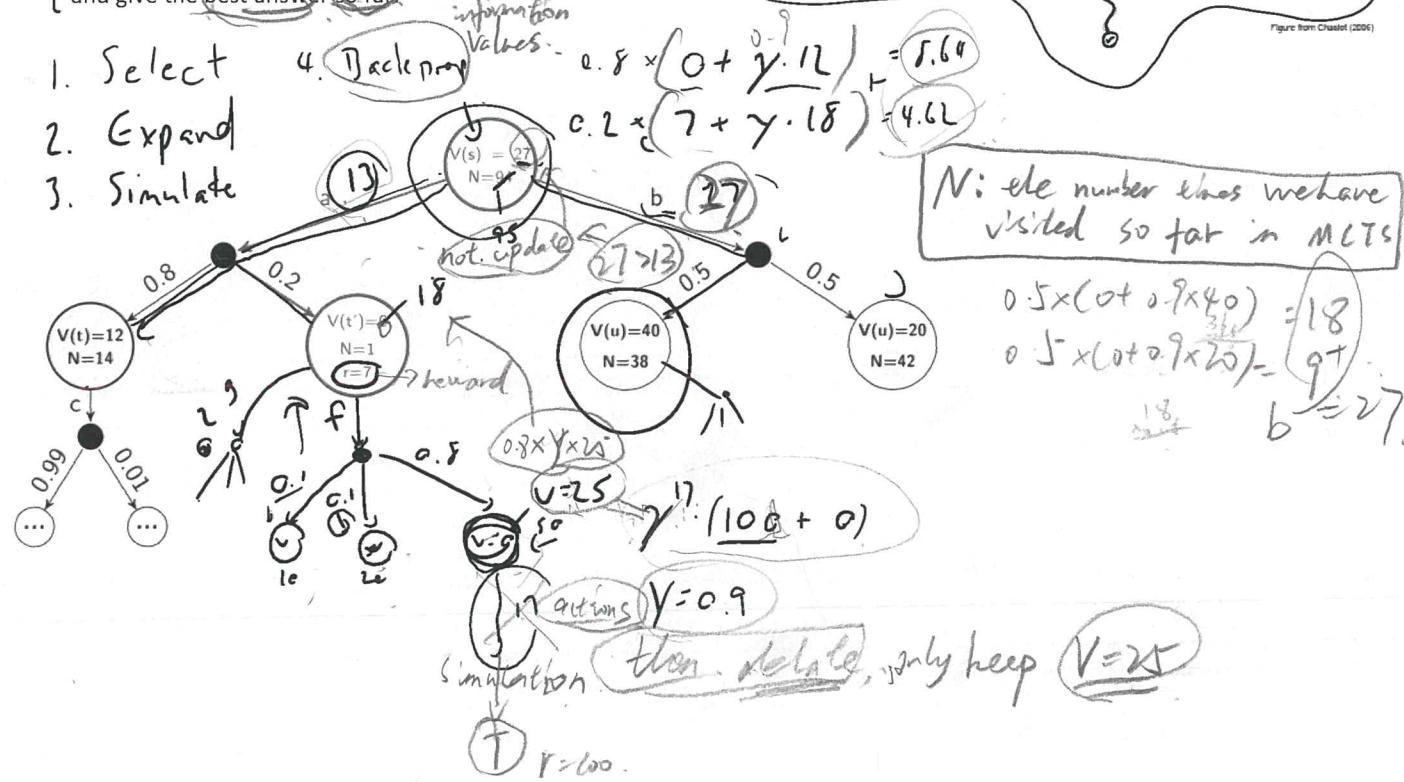
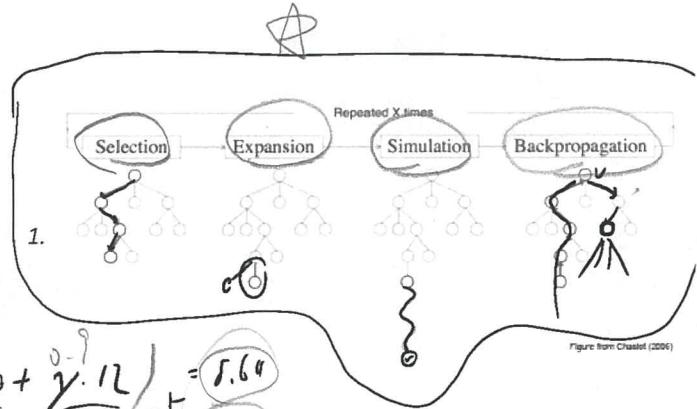


## Rec Lecture 10

$$V(s) = \max_{a \in A(s)} \sum_{s' \text{ children}} [r(s, a, s') + \gamma V(s')]$$

### Monte-Carlo Tree Search (MCTS) overview

1.  $V(s)$  is the estimate of the real value of a state.
2. But we will also use it as an heuristic.
3. The search tree is incrementally built.
4. MCTS is an anytime algorithm: we terminate whenever and give the best answer so far.



### Exercise

Your phone-stealing habits continue, but you are doing well and opening up to a new market of people. Each day, you can sell a bag of 20 iPhones, Samsung, Huawei, or Pixel phones, but the price varies with each buyer, and you don't know the probability that people will buy them. You decide to alternate: iPhones one day, Samsung the next, then Huawei, then Pixel.

After 100 days, you notice the following average return per day:

- Samsung: \$400
- iPhone: \$250
- Pixel: \$200
- Huawei: \$150

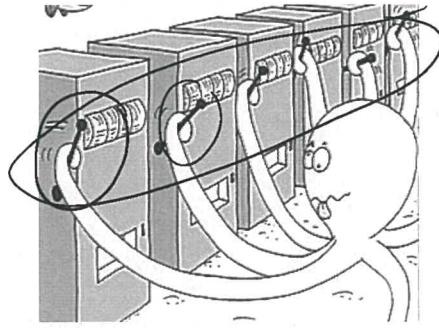


To help with your supply, you need to decide what you want to do for the next 100 days.  
 Assuming the next 100 days will look similar to the previous 100 days, what should you do?

<https://pollev.com/timothymille936>

## Multi-armed bandits

Imagine that you have  $N$  number of slot machines (or poker machines in Australia), which are sometimes called one-armed bandits. Over time, each bandit pays a random reward from an unknown probability distribution. Some bandits pay higher rewards than others. The goal is to maximize the sum of the rewards of a sequence of lever pulls of the machine.



### Exploration vs exploitation

1.  $\epsilon$ -greedy:  $\epsilon$  in  $[0, 1]$  (typically around 0.1), choose the best with probability  $1 - \epsilon$ , and other choose randomly.
2.  $\epsilon$ -decreasing: as above but  $\epsilon$  decreases over time
3. Softmax: choose proportionally

$$\pi(s) = \arg \max_a Q(s, a) + \sqrt{\frac{2 \ln N}{N(s, a)}}$$

which to choose = high reward + less visited.

Upper confidence bounds (UCB)

$$R(b) = \underset{a}{\operatorname{argmax}} Q(s, a) + \sqrt{\frac{2 \ln(N(s))}{N(s, a)}}$$

explore count → high for actions that have been explored less.

Upper confidence tree (UCT)  
 $Q$  value is high for actions that have a high reward.  
 $UCT = UCB + MCTS$  (almost!)

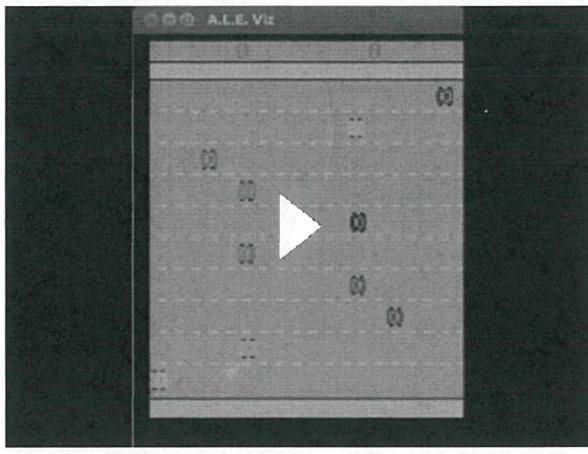
$$\pi(s) = \underset{a \in A(s)}{\operatorname{argmax}} Q(s, a) + 2C_p \sqrt{\frac{2 \ln N(s, a)}{N(s, a)}}$$

UCT playing Mario Brothers: A MCTS-based Mario-playing controller





UCT playing Freeway: [UCT Freeway - atari 2600](#)



Value/policy iteration vs. MCTS



Reinforcement learning: what if we do not know transitions  $P$  and reward function  $r$  of an MDP?

### The Mystery Game:

<https://programmingheroes.blogspot.com/2016/02/udacity-reinforcement-learning-mystery-game.html>

### Q-learning

1. Initialise  $Q(s,a)$  arbitrarily
2. For each episode:
  - a. Initialise  $s$  (go to the initial state)
  - b. Repeat for each step in the episode
    - i. Select the next action  $a$  to apply from  $s$  (using e.g. epsilon greedy, UCT) use  $Q(s,a)$
    - ii. Execute action  $a$  and observe the reward  $r$  and new state  $s'$
    - iii.  $Q(s,a) := Q(s,a) + \alpha[r + \gamma \max_{a'} Q(s',a') - Q(s,a)]$
    - iv.  $s := s'$
  - c. Until  $s$  is terminal

$$\underline{Q(s,a)}$$

learning rate



$$Q(s,a) = \underline{Q(s,a)} + \alpha \cdot [r + \gamma \cdot \max_{a'} Q(s',a')]$$

old value

$$- Q(s,a)$$

↑

do not want extra  $Q(s,a)$

### Q-Tables

State	Action			
	North	South	East	West
(0,0)	0.53	0.36	0.36	0.21
(0,1)	0.61	0.27	0.23	0.23
...	0.792			
(2,2)	0.79	0.72	0.90	0.72
(2,3)	0.90	0.78	0.99	0.81

$$Q(s_{(2,1)}, \text{South}) = 0.72$$

$$Q(s,a) = Q(s,a) + \alpha[r + \gamma \max_{a'} Q(s',a') - Q(s,a)]$$

$$Q(s,a) := Q(s,a) + \alpha[r + \gamma \max_{a'} Q(s',a') - Q(s,a)]$$



Learning rate  $\alpha = 0.1$

Discount reward factor  $\gamma = 0.9$

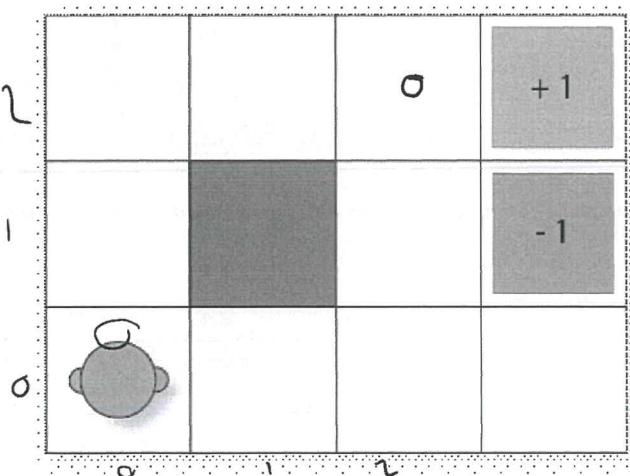
Q-learning:

$$Q((2,2), \text{North}) = 0.79 + 0.1 * (0 + 0.9 * 0.9 - 0.79) = 0.792$$

SARSA, with assumption that  $a'$  is West

$$Q((2,2), \text{North}) = 0.79 + 0.1 * (0 + 0.9 * 0.72 - 0.79) = 0.7828$$

If  $a'$  is East is for SARSA, the update would be just the same as for Q-learning because East is the max action from (2,2)



### Q-learning: Off-policy

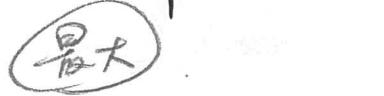
1. Initialise  $Q(s,a)$  arbitrarily
2. For each episode:
  - a. Initialise  $s$  (go to the initial state)
  - b. Repeat for each step in the episode
    - i. Select the next action  $a$  to apply from  $s$  (using e.g. epsilon greedy, UCT) use  $Q(s,a)$

### SARSA: On-policy learning

1. Initialise  $Q(s,a)$  arbitrarily
2. For each episode:
  - a. Initialise  $s$  (go to the initial state)
  - b. Select the next action  $a$  to apply from  $s$  (using e.g. epsilon greedy, UCT)
  - c. Repeat for each step in the episode



- b. Repeat for each step in the episode
- Select the next action  $a$  to apply from  $s$  (using e.g. epsilon greedy, UCT) use  $Q(s, a)$
  - Execute action  $a$  and observe the reward  $r$  and new state  $s'$  *We don't know action*
  - $Q(s, a) := Q(s, a) + \alpha [r + \gamma \max_{a'} Q(s', a') - Q(s, a)]$
  - $s := s'$
- c. Until  $s$  is terminal

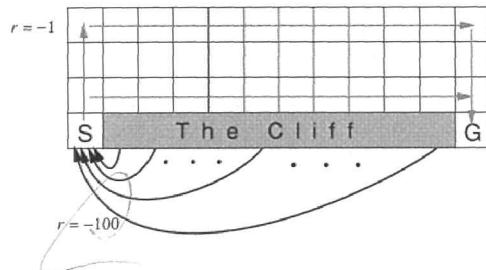


- b. Select the next action  $a$  to apply from  $s$  (using e.g. epsilon greedy, UCT)
- c. Repeat for each step in the episode
- Execute action  $a$  and observe the reward  $r$  and new state  $s'$  *We know action*
  - Select the next action  $a'$  to apply from  $s'$  (using e.g. epsilon greedy, UCT) *we know action*
  - $Q(s, a) := Q(s, a) + \alpha [r + \gamma Q(s', a') - Q(s, a)]$
  - $s := s'; a := a'$
- d. Until  $s$  is terminal

最大

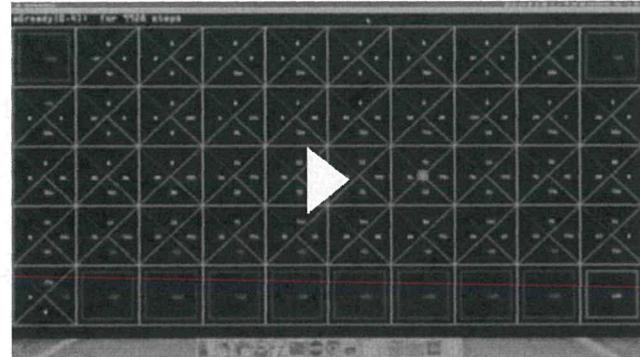
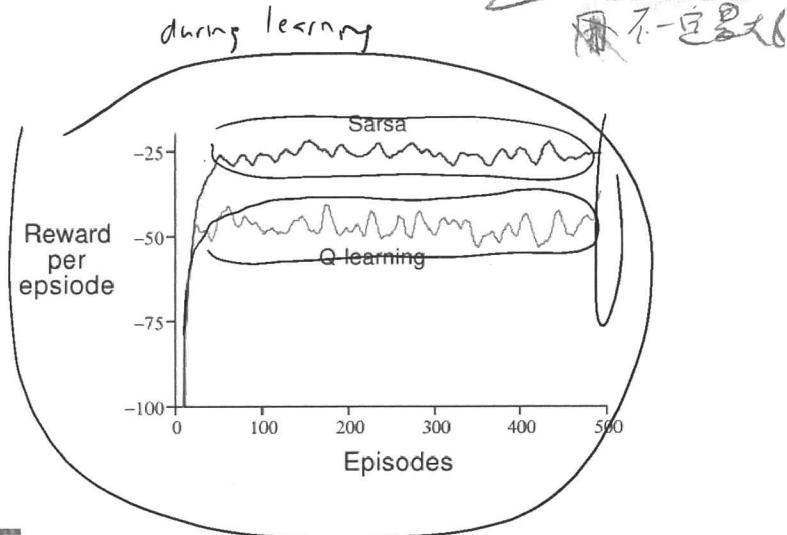
7-8-26

### Off-policy vs. on policy learning

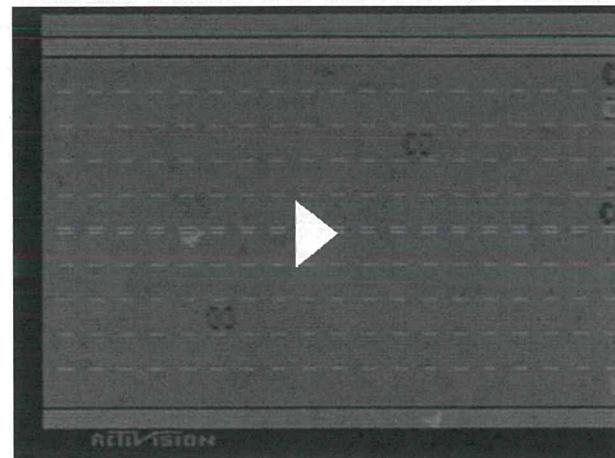


Gridworld Q-Learning - Example 3 - The Cliff

safe path  
optimal path



Learning to Play Freeway, using Reinforcement Learning

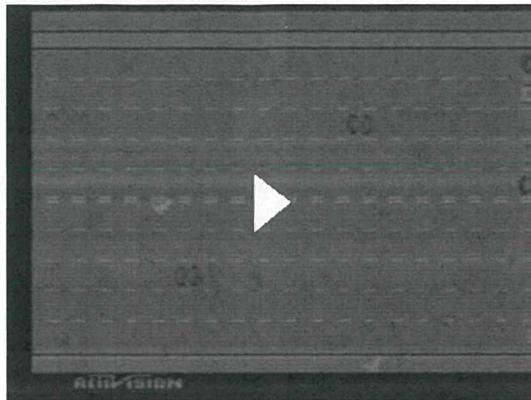




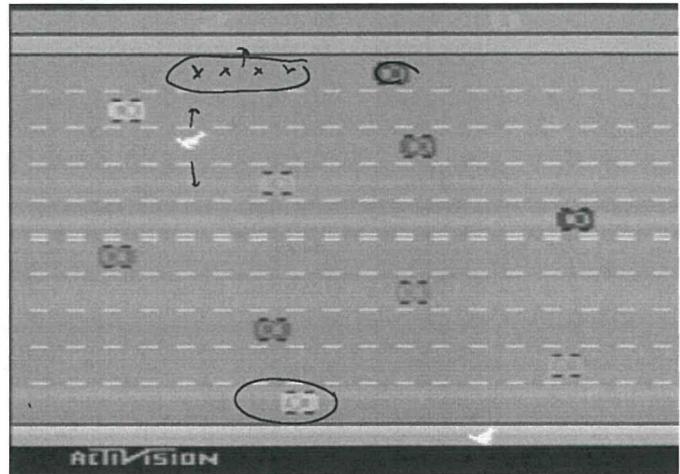
# Q-Function Approximation

## Lecture 12

Learning to Play Freeway, using Reinforcement Learning



$$Q(s, a) \approx f$$



Features and representation:

- Feature function returning a feature vector:

$$f(s, a) = \begin{pmatrix} f_1(s, a) \\ \vdots \\ f_n(s, a) \end{pmatrix} = \begin{pmatrix} 1/1 \\ 1/1 \\ 0 \\ 0 \\ 1 \\ 0 \\ 0 \\ 0 \end{pmatrix}$$

$$f(s, a) = \begin{pmatrix} f_1(s, a) \\ f_2(s, a) \\ \vdots \\ f_n(s, a) \end{pmatrix}$$

- Weight vector:

$$n \times |A|$$

number of actions.

Approximating  $Q(s, a)$

$$Q(s, a) = f_1(s, a) \cdot w_1^a + f_2(s, a) \cdot w_2^a + \dots + f_n(s, a) \cdot w_n^a$$

$$Q(s, a) = f_1(s, a) \cdot w_1^a + f_2(s, a) \cdot w_2^a + \dots + f_n(s, a) \cdot w_n^a = \sum_{i=1}^n f_i(s, a) w_i^a$$

Q-learning Update:

execute a  
→ for each feature  $i$ :

$$w_i^a \leftarrow w_i^a + \alpha [r + \gamma \max_{a'} Q(s', a') - Q(s, a)] f_i(s, a)$$

$Q(s', a')$  for SARSA

(1) Example (Freeway):  
Assume  $Q(s, a) = 0$  for all  $s, a$

(1) initialisation

(2) update

(2) Update:

$$w_i \leftarrow w_i + \alpha[r + \gamma \max_{a'} Q(s', a') - Q(s, a)] f_i(s, a)$$

$$w_{14}^{up} \leftarrow 0 + [10 + 0.9 \times 0]$$

$$w_i \leftarrow 5$$

(last row up)

Deep Q-learning:

$$\theta \leftarrow \theta + \alpha [r + \gamma \max_{a'} Q(s', a') - Q(s, a)] \frac{\partial Q(s, a)}{\partial \theta}$$

gradient

Q-learning:

$$w_i^a \leftarrow w_i^a + \alpha [r + \gamma \max_{a'} Q(s', a') - Q(s, a)] f_i(s, a)$$

SARSA:

$$w_i^a \leftarrow w_i^a + \alpha [r + \gamma Q(s', a') - Q(s, a)] f_i(s, a)$$

Deep Q-learning for robot grasping:

Learning Hand-Eye Coordination for Robotic Grasping

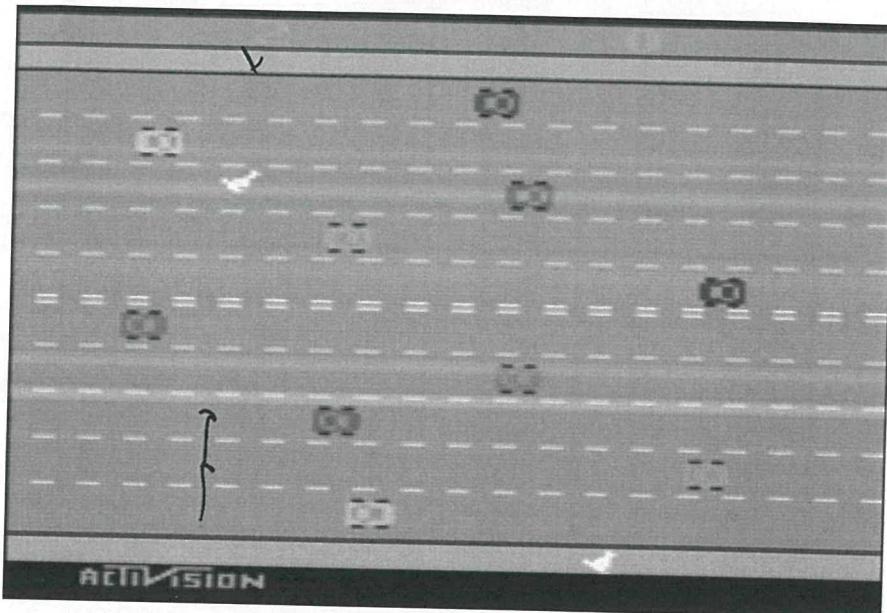


# Reward Shaping

Friday, 21 September 2018

15:22 PM

## Lecture 12



Key idea: add some small additional reward for particular behaviour:

$$Q(s, a) := Q(s, a) + \alpha [r + F(s, s') + \gamma \max a' Q(s', a') - Q(s, a)]$$

shaped reward

Potential-based Reward shaping:

$$F(s, s') = \gamma \Phi(s') - \Phi(s)$$

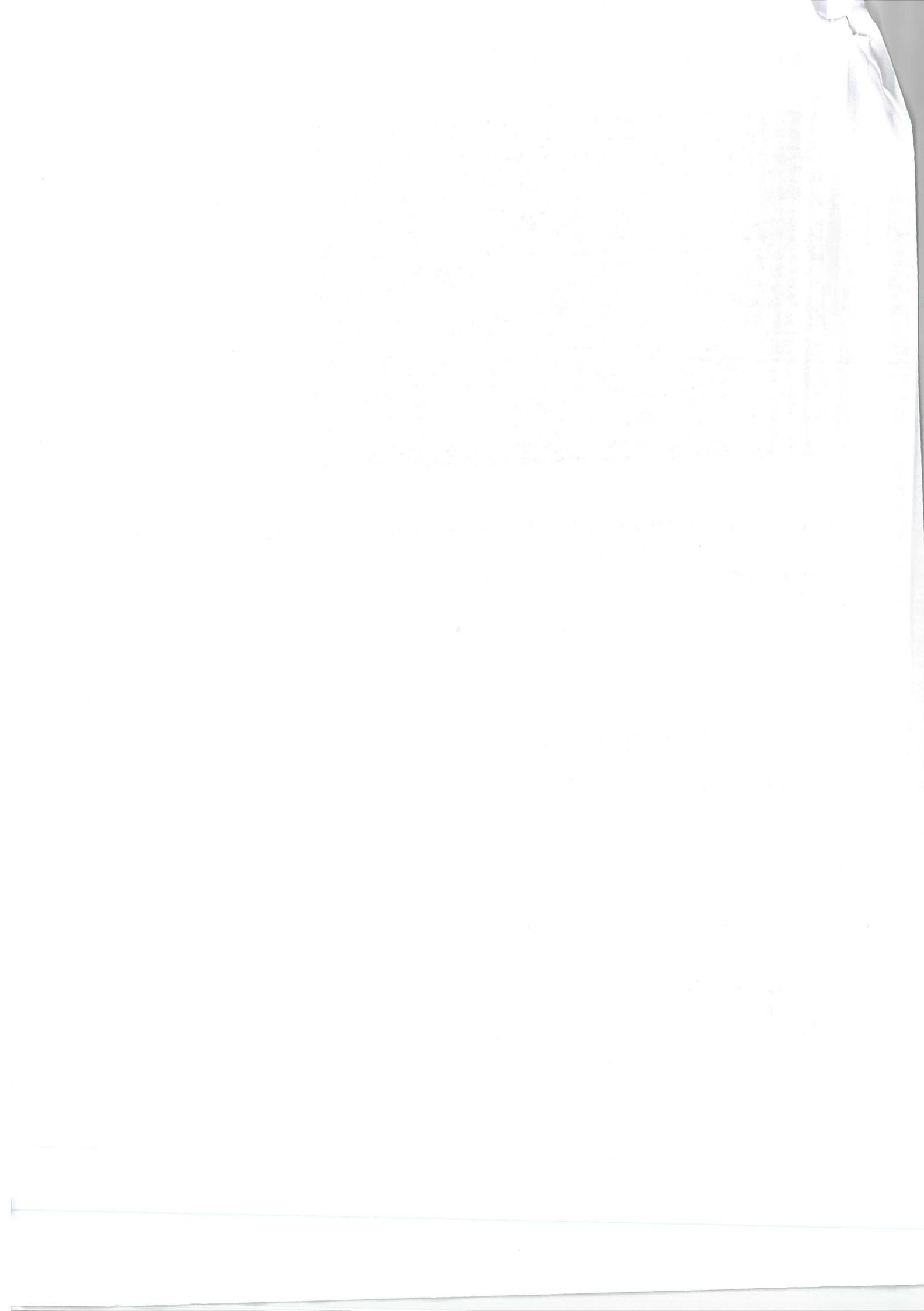
potential function.

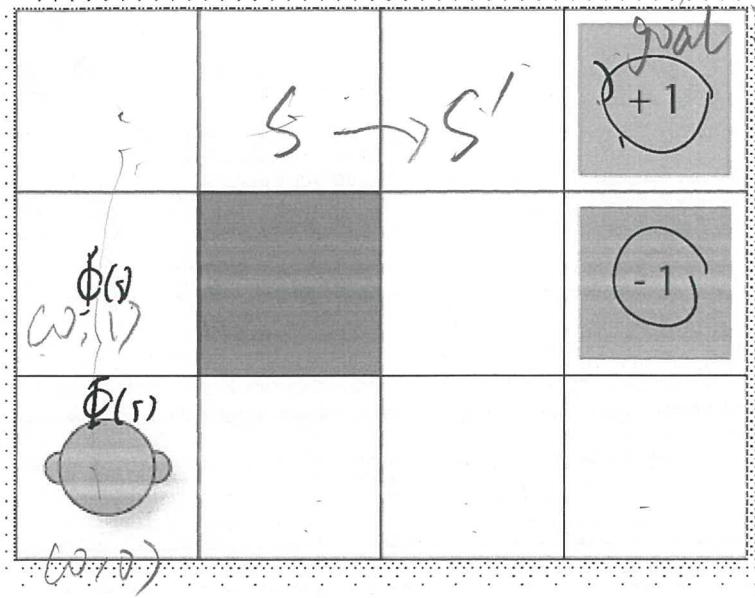
GridWorld example:

$$\Phi(s) = \max - \frac{1}{|x(g) - x(s)| + |y(g) - y(s)|}$$

$$\Phi(s) = \frac{1}{|x(g) - x(s)| + |y(g) - y(s)|}$$

Manhattan distance





Alternative to reward shaping: Q-function initialisation

(3, 2)

$$\begin{aligned}
 F(S, S') &= F((1, 2), (2, 2)) \\
 &= \gamma \Phi(2, 2) - \Phi(1, 2) \\
 &= 0.9 \times \frac{1}{|3-2+1-2-2|} - \frac{1}{|1-1+1-2-2|} \\
 F((1, 2), (2, 2)) &= \cancel{\gamma \cdot \Phi(2, 2)} - \cancel{\Phi(1, 2)} \\
 &= 0.9 \times \frac{1}{5} - \frac{1}{2} = 0.4 \\
 F((0, 0), (0, 1)) &= \cancel{0.9 \times \cancel{\Phi(1, 2)}} \\
 &= 0.9 \times \frac{1}{|3-0+1-2-0|} - \frac{1}{|3-0+1-2|} \\
 &= 0.9 \times \frac{1}{5} - \frac{1}{4}
 \end{aligned}$$



# Lecture 13

n-step temporal difference learning

Discounted future rewards

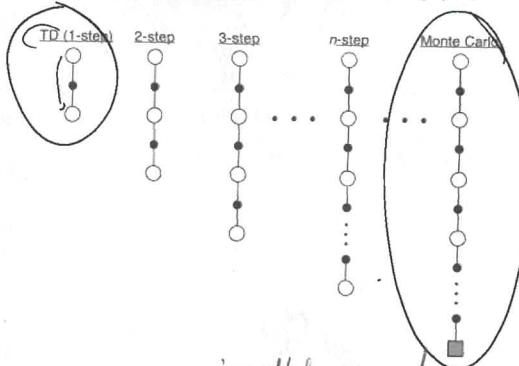
$$G_t = r_t + \gamma V(s')$$

$$= r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \dots + \gamma^{n-1} r_{t+n-1} = R_t + \gamma G_{t+1}$$

$$= r_1 + \gamma r_2 + \gamma^2 r_3 + \gamma^3 r_4 + \dots$$

$$= r_1 + \gamma(r_2 + \gamma(r_3 + \gamma^2 r_4 + \dots))$$

Truncated future rewards



$$\text{SARSA: } Q(s, a) := Q(s, a) + \alpha [r + \gamma Q(s', a') - Q(s, a)]$$

Change the update:

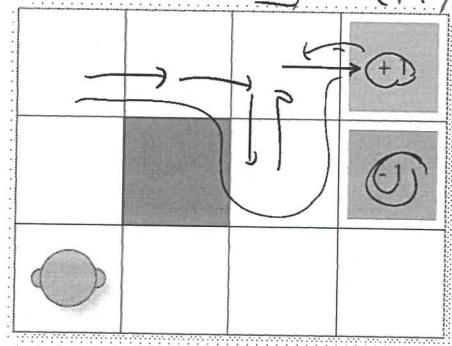
$$G_t = \sum_{i=t+1}^{\text{terminal}} \gamma^{i-t-1} r_i \quad n=1$$

$$Q(s, a) \leftarrow Q(s, a) + \alpha (G_t + \gamma^n Q(s_{t+n}, a_{t+n}) - Q(s, a))$$

With 1-step learning

State	Action			
	North	South	East	West
(0,0)	0	0	0	0
(0,1)	0	0	0	0
(0,2)	0	0	0	0
...	...	...	...	...
(1,2)	0	0	0	0
(2,1)	0	0	0	0
(2,2)	0	0	0.45	0
(2,3)	0	0	0	0
...	...	...	...	...

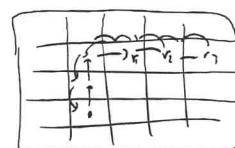
$$G(s_t) \leftarrow Q(s_t, a_t) + \alpha (r_t + \gamma G(s_{t+1}))$$



$$G_t = r_t + V(s_{t+1}) \text{ one step value.}$$

$$G_t = r_t + \gamma r_{t+1} + \gamma^2 V(s_{t+2}) \text{ two step value.}$$

$$G_t = r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \dots + \gamma^n V(s_{t+n}) \text{ n step.}$$



$$Q(s, a) \leftarrow Q(s, a) + \alpha [G_t + \gamma^n Q(s_{t+n}, a_{t+n}) - Q(s, a)]$$

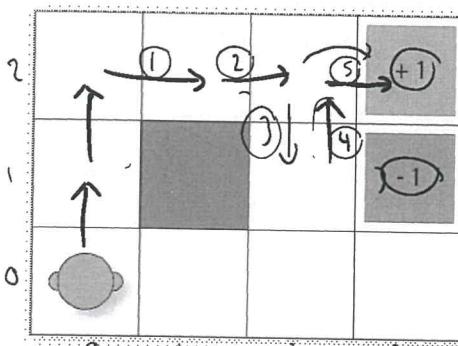
```

n-step SARSA for estimating Q ≈ q_n or Q ≈ q_\infty for a given \alpha
Initialize Q(s, a) arbitrarily, for all s \in S, a \in A
Initialize \pi to be \gamma-greedy with respect to Q, or to a fixed given policy
Parameters: step size \alpha \in (0, 1], small \gamma > 0, a positive integer n
All store and access operations (for S_t, A_t, and R_t) can take their index mod n
Repeat (for each episode):
    Initialize and store an action A_0 \sim \pi(\cdot | S_0)
    T \leftarrow \infty
    For t = 0, 1, 2, ...:
        If t < T, then:
            Take action A_t
            Observe and store the next reward as R_{t+1} and the next state as S_{t+1}
            If S_{t+1} is terminal, then:
                T \leftarrow t+1
            else:
                Select and store an action A_{t+1} \sim \pi(\cdot | S_{t+1})
        If t > 0:
            r \leftarrow t - n + 1 (r is the time whose estimate is being updated)
            G \leftarrow \sum_{i=r}^{\min(t-n+1, T)} \gamma^{i-r-1} R_i
            If r + n < T, then G \leftarrow G + \gamma^n Q(S_{r+n}, A_{r+n})
            Q(S_r, A_r) \leftarrow Q(S_r, A_r) + \alpha [G - Q(S_r, A_r)]
            If \pi is being learned, then ensure that \pi(\cdot | S_r) is \gamma-greedy wrt Q
        Until r = T - 1
    
```

$$G_t = \sum_{i=t+1}^{\text{terminal}} \gamma^{i-t-1} r_i$$

$$Q(s, a) \leftarrow Q(s, a) + \alpha [G_t + \gamma^n Q(s_{t+n}, a_{t+n}) - Q(s, a)]$$

Exercise: Grid World



Compute 5-step SARSA update

$$\alpha = 0.5$$

$$\gamma = 0.9$$

$$\begin{aligned}
 G_{-5} &= \gamma \cdot 1 & \gamma \cdot 1 &= 0.9 \\
 G_{-4} &= \alpha \gamma^2 \cdot 1 & \gamma^2 \cdot 1 &= 0.81 \\
 G_{-3} &= 0 + \alpha + \gamma^3 \cdot 0 + \alpha + \gamma^3 \cdot 1 & 0 + \alpha + \gamma^3 \cdot 1 &= 0.81 \\
 G_{-2} &= ; & 0 + \alpha + \gamma^4 \cdot 1 &= 0.81 \\
 G_{-1} &= ; & 0 + \alpha + \gamma^5 \cdot 1 &= 0.81 \\
 Q(s_{(2,2)}, E) &= 0 + 0.5 [G_5 + 0 - 0] & 0 + 0.5 [G_5 + 0 - 0] &= 0.45 \\
 Q(s_{(1,2)}, N) &= 0 + 0.5 [G_4 + 0 - 0] & 0 + 0.5 [G_4 + 0 - 0] &= 0.405 \\
 \dots & & & \\
 Q(s_{(1,1)}, S) &= ; & &
 \end{aligned}$$

$$\begin{aligned}
 0 + 0.5 [G_5 + 0 - 0] &= 0.45 \\
 0 + 0.5 [G_4 + 0 - 0] &= 0.405
 \end{aligned}$$



Diagram showing a 2x3 grid world with actions North (N), South (S), East (E), and West (W). States are labeled (0,0) through (2,3).

$$Q(s_{(2,2)}, E) = Q + \alpha [0.5(0.5 + 0 - 0)] = 0.45$$

$$Q(s_{(1,2)}, N) = Q + \alpha [0.5(0.5 + 0 - 0)] = 0.405$$

$$\dots$$

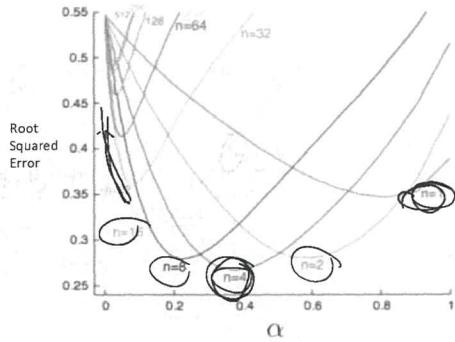
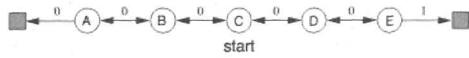
$$Q(s_{(2,2)}, S) = \dots$$

$$Q(s_{(2,1)}, E) = \dots$$

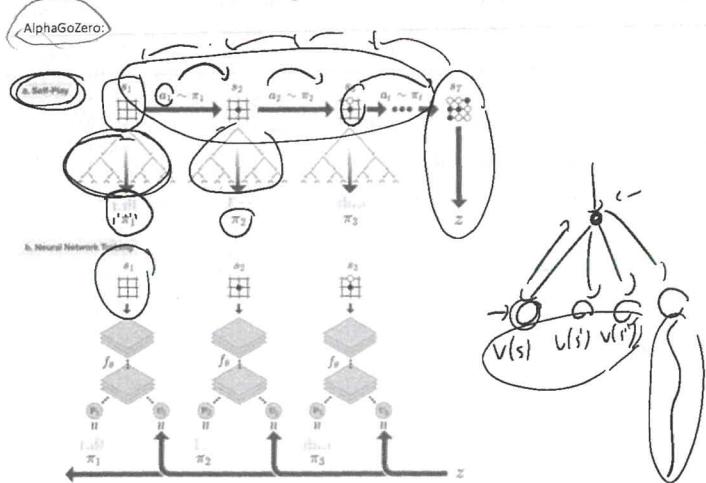
$$Q(s_{(1,1)}) = Q(r_{(1,1)}) + \alpha [Q_{(1,1)} + \gamma^* Q(s_{(1,1}') - Q(s_{(1,1)}))]$$

With 1-step learning					With 5-step learning				
State	Action				State	Action			
	North	South	East	West		North	South	East	West
(0,0)	0	0	0	0	(0,0)	0	0	0	0
(0,1)	0	0	0	0	(0,1)	0	0	0	0
(0,2)	0	0	0	0	(0,2)	0	0	0.2953	0
...					...				
(1,2)	0	0	0	0	(1,2)	0	0	0.3281	0
(2,1)	0	0	0	0	(2,1)	0.405	0	0	0
(2,2)	0	0	0.45	0	(2,2)	0	0.3645	0.45	0
(2,3)	0	0	0	0	(2,3)	0	0	0	0
...					...				

Example: Random walk



MCTS + Reinforcement learning





# Normal form games

## Prisoner's dilemma

Consider the following (hopefully fictional) scenario. You and your partner in crime have been arrested for armed robbery. To avoid a lengthy trial, the police put you and your partner in crime in separate cells, *unable to communicate with each other*. The police offer each of you the following deal:

- If neither of you admit your guilt, you will be charged with carrying illegal weapons, and receive 1 year each in prison.
- If one of you admits both are guilty, while the other does not, the prisoner that admitted will get off free in exchange for the confession, while the other prisoner will receive 4 years in prison.
- If both of you admit you are guilty, you will each receive 2 years in prison — the length is reduced from 4 years in exchange for your confession.

Assuming that you care little for your partner (you are an armed robber, after all), what choice will you make? Will you cooperate with your partner in crime and deny your guilt, or will you betray your partner and admit your guilt? Remember, you have no idea which choice your partner will make, and you cannot communicate.

A: Admit

B: Deny

<https://pollev.com/timothymille936>

## Normal Form game

Players are + players

→ moves (but has)

→ utility / payoffs

$u_i(a_j)$

Nash equilibrium → admit B deny

(A)	admit		deny	
	admit	-2, -2	0, -4	-4, 0
admit	-2, -2	0, -4	-4, 0	-1, -1

① dominant strategy

] agents

② Nash equilibrium

] games

→ best response

(admit, admit)

if B admit  $\Rightarrow$  A admit

A B deny  $\Rightarrow$  A admit

Nash equilibrium

is there any action make better

## The Advertising Game

Firm 2 best response

		Not	Adv
		Not	16, 12
Firm 1	Not	16, 12	7, 13
	Adv	16, 12	7, 13

F1: dominant strategy. (Not)

F2: no dominant strategy. equilibrium

equilibrium

dominated strategy

If Firm 2:



A: Advertise  
B: Not Advertise

<https://pollev.com/timothymille936>

### Split or Steal?

		Player 2	
		Split	Steal
Player 1	Split	0.5, 0.5	0, 1
	Steal	1, 0	0, 0

Three equilibria

- A: Split, Split
- B: Split, Steal
- C: Steal, Split
- D: Steal, Steal

<https://pollev.com/timothymille936>

golden balls. the weirdest split or steal ever!





## Mixed strategies

Pure Strategy. ~~2 ways~~

### Matching pennies

		Player Odd	
		$y$	$1-y$
		Heads	Tails
Player Even	X	(Heads)	$1, -1$
	$1-x$	(Tails)	$-1, 1$

A: Heads

B: Tails

<https://pollev.com/timothymille936>

$$\begin{aligned} E(H) &= E(T) \\ 1-2x &= 2x-1 \\ 2 &= 4x \\ x &= \frac{1}{2} \\ 1-x &= \frac{1}{2} \end{aligned}$$

$$E(t) = E(H)$$

Mixed strategy

Indifference:

$$\begin{aligned} \text{Odd: } & \text{Even play Heads.} \\ E(H) &= -dx + 1-x \\ &= 1-2x \\ E(T) &= x + -1(1-x) \\ &= 2x-1 \end{aligned}$$

Even:

$$\begin{aligned} E(H) &= 2y-1 \\ E(T) &= 1-2y \end{aligned}$$

$$\begin{aligned} 2y-1 &= 1-2y \\ 4y &= 2 \\ y &= \frac{1}{2} \end{aligned}$$

Adversary:

$$\begin{aligned} E(T_1) &= -3x + 5(1-x) \\ &= -8x + 5. \end{aligned}$$

$$E(T_2) = x - (1-x) = 2x-1$$

$$E(T_1) = E(T_2)$$

$$\begin{aligned} -8x+5 &= 2x-1 \\ 6 &= 10x \\ x &= \frac{3}{5}. \end{aligned}$$

Defender

$$E(T_1) = 5y - (1-y) = 6y-1$$

$$\begin{aligned} E(T_2) &= -5y + 2(1-y) \\ \text{Defender: } & -7y+2. \end{aligned}$$

$$\begin{aligned} E(T_1) &= 6y-1 \\ E(T_2) &= -7y+2 \end{aligned}$$

$$\begin{aligned} 6y-1 &= -7y+2 \\ 13y &= 3 \\ y &= \frac{3}{13} \end{aligned}$$

### Security games



		Terminal 1	Terminal 2
		$\frac{3}{13}, -3$	$-1, 1$
		$5, -5$	$2, -1$
Defender	X	$5, -3$	$-1, 1$
	$1-x$	$-5, 5$	$2, -1$

Adversary:

$$\begin{cases} E(T_1) = -3x + 5(1-x) \\ = 5-8x \\ E(T_2) = x + -1(1-x) \\ = 2x-1 \end{cases}$$



$$\begin{aligned}
 E(T_2) &= x + -1(1-x) \\
 &= 2x - 1 \\
 E(T_1) &= E(T_2) \\
 5 - 8x &= 2x - 1 \\
 6 &= 10x \\
 \{ \quad x &= \frac{6}{10} \\
 1 - x &= \frac{4}{10} = \frac{2}{5}
 \}
 \end{aligned}$$

Defender strategy

Ferry escort



$$E(T_1) = 1 - x$$

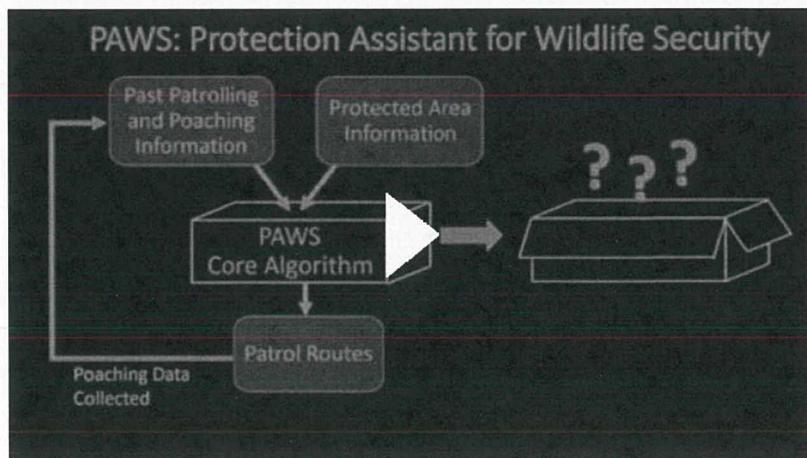
$$E(T_1) = E(T_2)$$

$$6y - 1 = 2 - 7y$$

$$\begin{aligned} y &= \frac{3}{13} \\ 1 - y &= \frac{10}{13} \end{aligned}$$

Adversary strategy

Save the Wildlife, Save the Planet: Protection Assistant for Wildlife Security (PAWS)



Player 2

$$E(\text{Up}) = x + 0.$$

$$\begin{aligned}
 E(\text{Right}) &= 0 \cdot x + 3(1-x) \\
 &= -3x + 3
 \end{aligned}$$

$$E(\text{Left}) = E(\text{Right}).$$

$$x = -3x + 3$$

$$4x = 3$$

$$x = \frac{3}{4}$$

$$\begin{aligned}
 \text{Player 1: } E(\text{Up}) &= 3y & 3y - 1y & 4y = 1 \\
 E(\text{Down}) &= 1y & y = \frac{1}{4}
 \end{aligned}$$

Pure and mixed equilibria

		Player 2	
		Left	Right
Player 1	Up	(3, 1)	0, 0
	Down	0, 0	1, 3

$$E(L) = 1x + 0$$

$$E(R) = 0 + 3(1-x)$$

$$E(U) = E(R)$$

$$1x = 3 - 3x$$



	$x$	Up	<span style="border: 1px solid black; padding: 2px;">3, 1</span>	0, 0	$\frac{1}{4}x - \frac{1}{4}(1-x)$
Player 1	$1-x$	Down	0, 0	<span style="border: 1px solid black; padding: 2px;">1, 3</span>	$\frac{3}{4}x = 3 - 3x$ $x = \frac{3}{4}$

One mixed equilibria

$$\begin{cases} x = \frac{3}{4}, 1-x = \frac{1}{4} \\ y = \frac{1}{4}, 1-y = \frac{3}{4} \end{cases}$$

Two pure equilibria



## Extended form games

### Buyer/seller game

Seller:

- 1) Sell low quality = 3
- 2) Sell high quality = 2
- 3) Don't sell low quality = 1
- 4) Don't sell high quality = 0

Buyer:

- 1) Buy high quality, payoff = 2
- 2) Buy nothing, payoff = 1
- 3) Buy low quality, payoff = 0

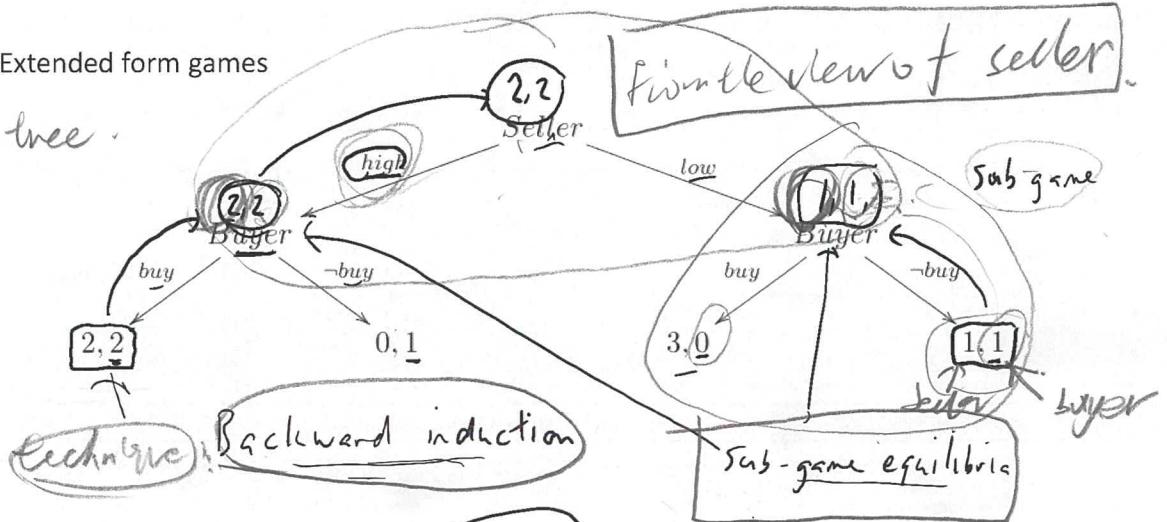
		Buyer	
		buy	don't buy
Seller	high	2, 2	0, 1
	low	3, 0	1, 1

dominated strategy

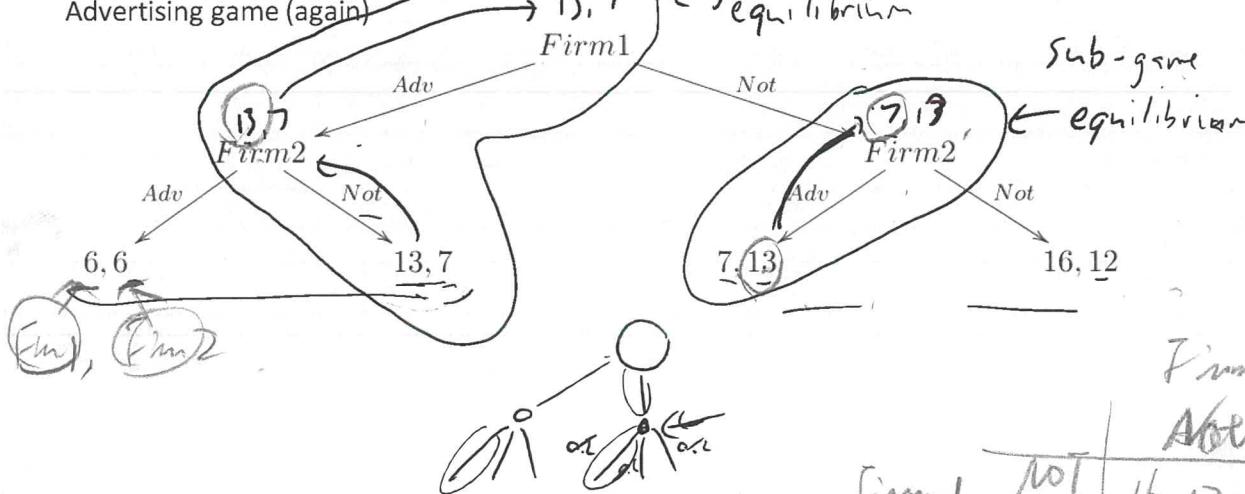
one pure equilibrium

### Extended form games

Game tree



### Advertising game (again)

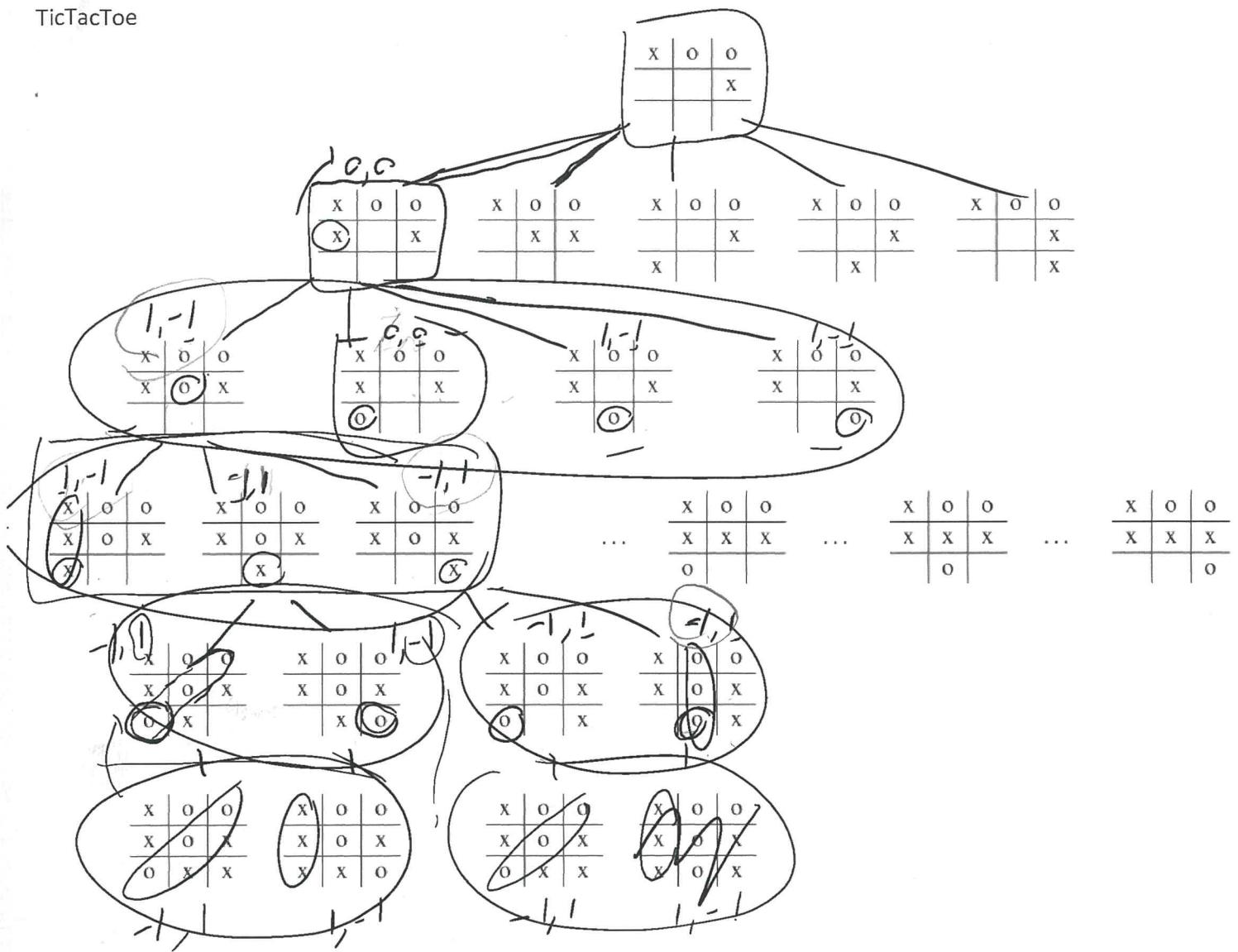


different outputs

		Not	Adv
Firm 1	Not	16, 12	7, 13
	Adv	13, 7	6, 6



# TicTacToe





## Revision (MCTS, Q-learning & SARSA)

$$MCTS \quad V(s) = \max_{a \in A(s)} \sum P(a|s) [r(s,a,s') + \gamma V(s')]$$

