

# CS540 Introduction to Artificial Intelligence

## Lecture 13

Young Wu

Based on lecture slides by Jerry Zhu and Yingyu Liang

July 2, 2019

## Secretary Problem

## Motivation

- Interview 10 people, random order, either give an offer or reject immediately after each interview. The goal is to give an offer to the best candidate. Optimal strategy: interview first  $n$  people, give an offer to the first candidate who is better than all previous ones. What is  $n$ ?  
A: 1, B: 2, C: 3, D: 4, E: 5

$$n = \frac{10}{e} \approx 3$$

$\approx 2.7$

Prob { hitkey }  
best  $\leq \frac{1}{e}$

who is better than ]  
always reject n  
filed best so far  
offer n

# Secretary Problem Solution

## Motivation

## Schedule

## Admin

- Thursday, July 4: Post sample midterm and formula sheet.
  - Tuesday, July 9: Tan review session: review + sample midterm.
  - Thursday, July 11: Midterm Version A.
  - Tuesday, July 16: Midterm Version B.

posted

## Midterm

Admin

- 2 hour midterm, 5 : 30 to 7 : 30 +  $\varepsilon$ ,  $\varepsilon > 0$ .
  - Which midterm will you attend?
  - A: Regular: Thursday, July 11.
  - B: Alternative: Tuesday, July 16.
  - C: Cannot make both.

# Reinforcement Learning

## Motivation

- Reinforcement learning is about learning from the outcome of actions.
- ➊ Sense world.
  - ➋ Reason.
  - ➌ Choose an action to perform.
  - ➍ Get feedback.
  - ➎ Learn.

# Applications

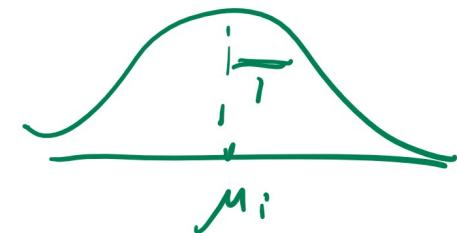
## Motivation

- Actions can be performed in the physical world or artificial ones.
- Board games.
- Robotic control.
- Autonomous helicopter performance.
- Economics models.

# Bandits

## Motivation

- There are  $K$  arms, pulling each arm  $i$  results in reward  $r_i$ .
- The reward  $r_i$  is random and follows Gaussian distribution with mean reward  $\mu_i$  and variance  $\sigma^2 = 1$ .
- Suppose  $\mu_1 \geq \mu_2 \geq \mu_3 \geq \dots \geq \mu_K$ .



# Exploration then Exploitation Algorithm

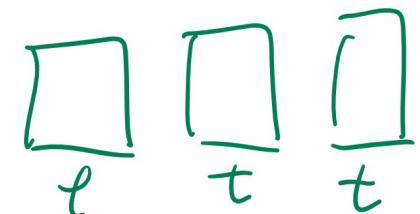
## Motivation

$$\mu_1 > \mu_2 \approx \mu_3$$

- ① Pull each arm  $t$  times to estimate the mean reward.



$$\hat{\mu}_{i,t} = \frac{1}{n} \sum_{t'=1}^t r_{i,t'}$$



$r_{i,t'}$  is the random reward from arm  $i$  and  $t'$ -th pull.

$$\mu_i \approx \hat{\mu}_{i,t}$$

after  
 $tK$   
periods

- ② Pull the arm  $i^*$  with the highest estimated mean reward.

$$i^* = \arg \max_{i=1,2,\dots,K} \hat{\mu}_{i,t}$$



# Upper Confidence Bound Algorithm

## Motivation

- ① Pull the arm  $i^*$  with the highest upper confidence bound.

$$i^* = \arg \max_{i=1,2,\dots,K} \left\{ \begin{array}{l} \text{UCB} \\ \infty \end{array} \right\} = \hat{\mu}_{i,t} + \sqrt{\frac{2 \log \left( \frac{1}{\delta} \right)}{t}} \quad t > 0$$

$\delta$  is the confidence level parameter.

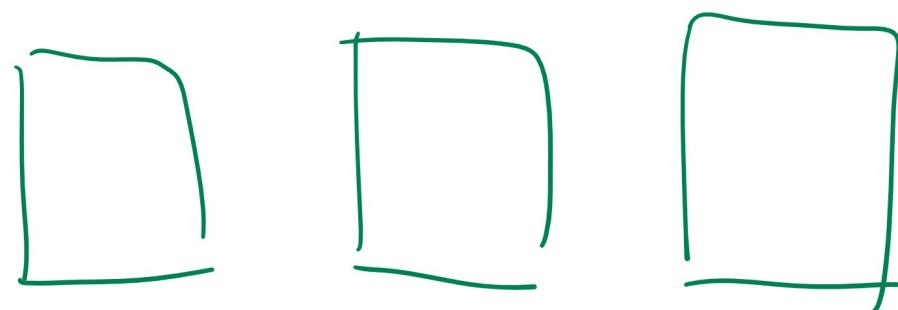
$$\mathbb{P} \left\{ \mu_i \leq \hat{\mu}_{i,t} + \sqrt{\frac{2 \log \left( \frac{1}{\delta} \right)}{t}} \right\} \leq 1 - \delta$$

0.99  
0.01

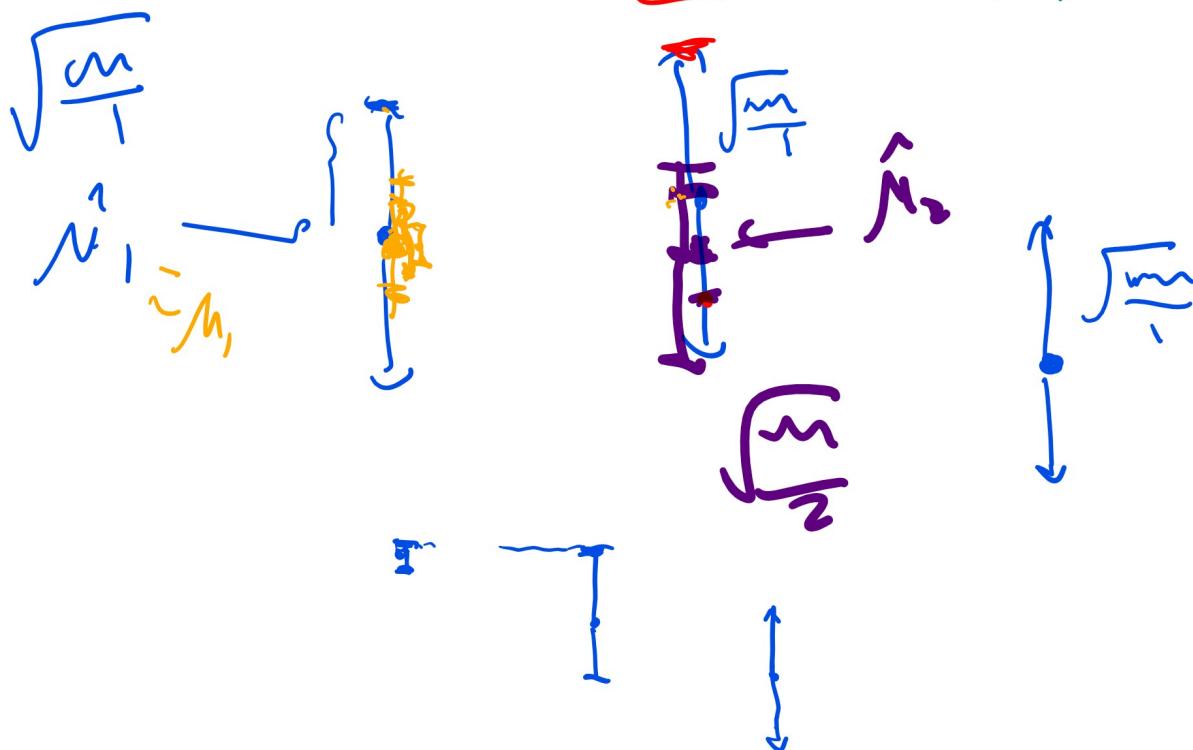
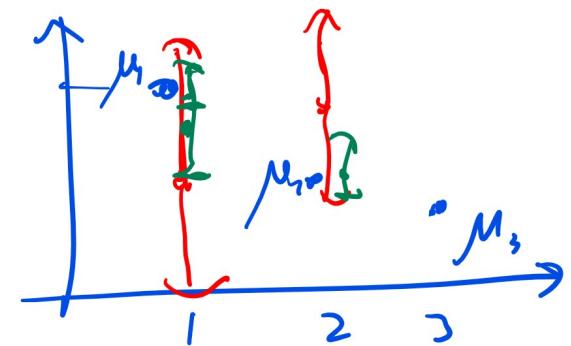
0.99 with high confidence

# UCB Algorithm Diagram

Motivation



$$M_1 > \underline{M}_2 > N_3$$



pull every thing once  
K.

# Bandit Applications

## Motivation

- Managing research projects.
  - Treatment for patients.
  - Search engine ranking.
  - Wireless adaptive routing.
  - Financial portfolio design.

# Q Learning

## Description

- Select an action.
- Receive reward.
- Observe new state.
- Update (learn) the value of the state-action pair.

# State and Actions

## Definition

- The set of possible states is  $s_t \in S$ .
- The set of possible actions is  $a_t \in A$ .
- The set of possible rewards is  $r_t \in R$ .
- At each time  $t$  :
  - ① Observe state  $s_t$ .
  - ② Chooses action  $a_t$ .
  - ③ Receives reward  $r_t$ .
  - ④ Changes to state  $s_{t+1}$ .

# Markov Decision Process

## Definition

- Markov property on states and actions is assumed.

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

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

- The goal is to learn a policy function  $\pi : S \rightarrow A$  for choosing actions that maximize the total expected discounted reward.

$$\mathbb{E}[r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \dots], \gamma \in [0, 1]$$

random

# Expected Reward

## Definition

- The expected reward at a given time  $t$  is the average reward weighted by probabilities.

$$\mathbb{E}[r_t] = \sum_{r_t \in R} r_t \mathbb{P}\{r_t | s_{t-1}, a_{t-1}\}$$

# Discounted Reward

## Definition

- The discounted reward at time 0 is the sum of reward weighted given the time preference, usually described by a constant discount factor.

$$PV(r_t) = \gamma^t r_t, \gamma \in [0, 1]$$

$$PV(r_1, r_2, \dots) = \sum_{t=0}^{\infty} \gamma^t r_t$$

*1 +  $\gamma + \gamma^2 + \gamma^3 + \dots$*

*1 -  $\gamma$*

- $\gamma$  is the value of 1 unit of reward at time 1 perceived at time 0. If  $\gamma = 1$ , the sum over an infinite time period is usually infinity, therefore  $\gamma < 1$  is usually used.

# Value Function

## Definition

- The value function is the expected discounted reward given a policy function  $\pi$ , assuming the action sequence is chosen according to  $\pi$  starting with state  $s$ .

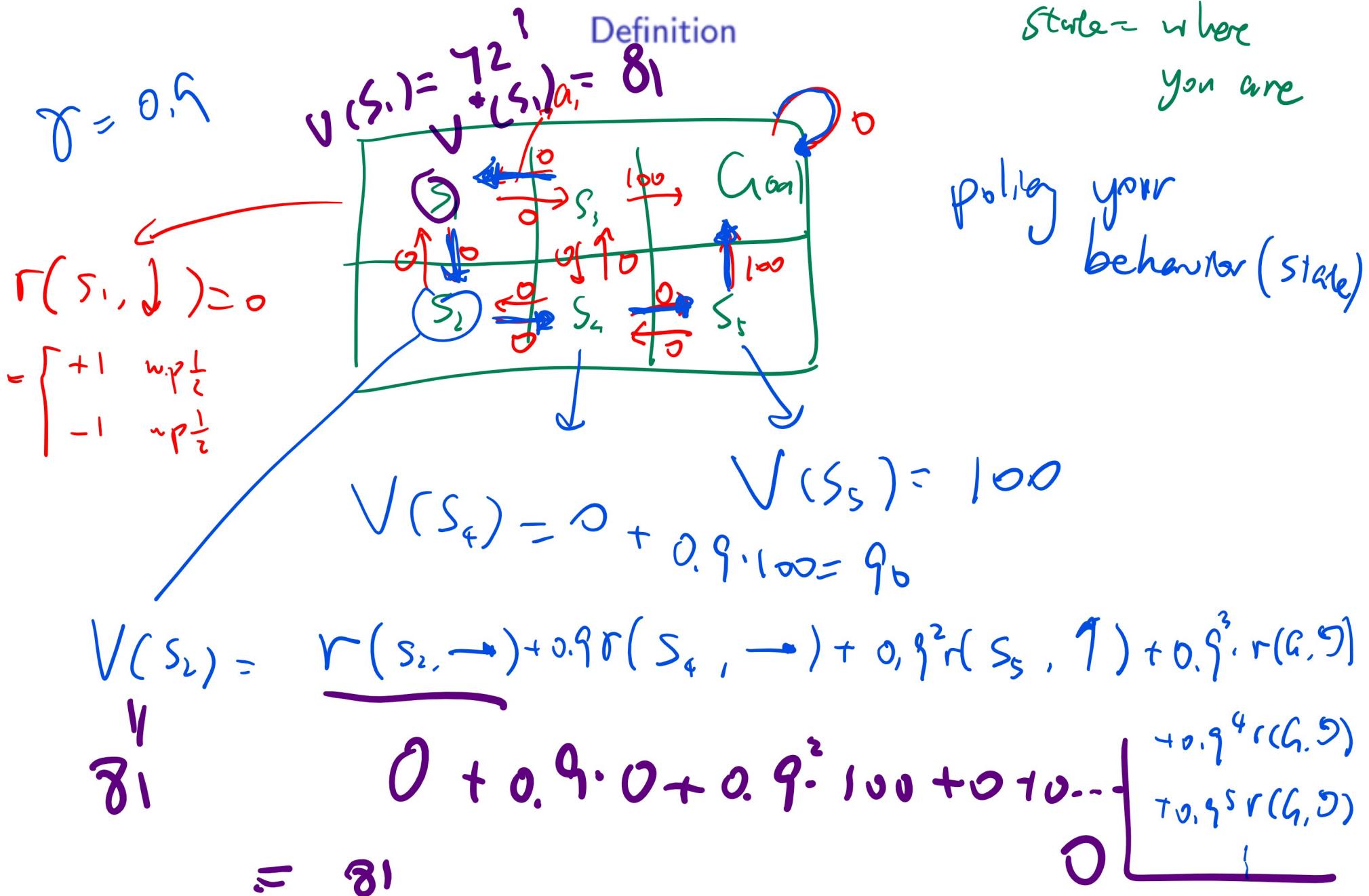
$$V^\pi(s) = \sum_{t=0}^{\infty} \gamma^t \mathbb{E}[r_t]$$

- The optimal policy  $\pi^*$  is the one that maximizes the value function.

$$\pi^* = \arg \max_{\pi} V^\pi(s) \text{ for all } s \in S$$

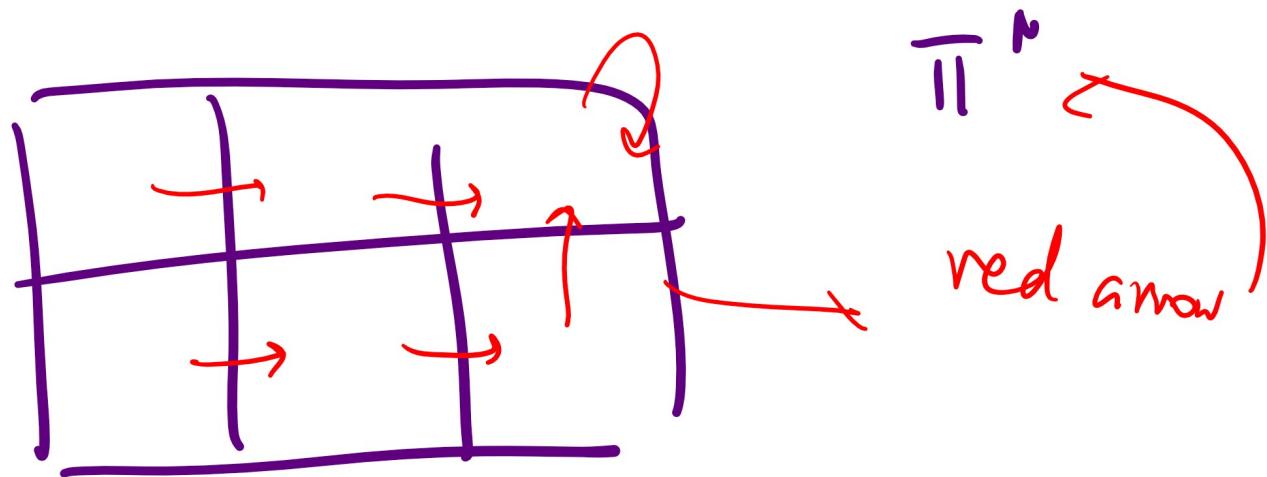
$$V^*(s) = V^{\pi^*}(s)$$

## Goal Learning Example, Part I



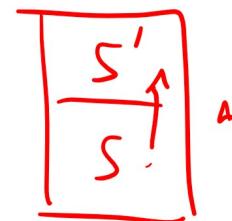
# Goal Learning Example, Part II

Definition



# Optimal Policy Given Value Function

Definition



$$V(s) = r(s \rightarrow s', a) + V^*(s')$$

- Given  $V^*(s)$ ,  $r(s, a)$ ,  $\mathbb{P}(s'|s, a)$ ,  $\pi^*$  can be computed directly.

$$\begin{aligned} \pi^*(s) &= \arg \max_{a \in A} (\underbrace{\mathbb{E}[r|s, a] + \gamma \mathbb{E}[V^*(s')|s, a]}_{\text{Current reward} + \text{discounted all reward}}) = V^*(s) \\ &= \arg \max_{a \in A} \left( \sum_{r \in R} r \mathbb{P}\{r|s, a\} + \gamma \sum_{s' \in S} \mathbb{P}\{s'|s, a\} V^*(s') \right) \\ &\quad \text{in next period starting at state } s' \end{aligned}$$

Annotations in red:

- $\mathbb{E}[r|s, a]$  is circled and labeled "Current reward".
- $\mathbb{E}[V^*(s')|s, a]$  is circled and labeled "discounted all reward".
- $\sum_{r \in R} r \mathbb{P}\{r|s, a\}$  is circled and labeled "Q(s, a)".
- $\sum_{s' \in S} \mathbb{P}\{s'|s, a\} V^*(s')$  is circled and labeled "in next period starting at state  $s'$ ".

- Define the function inside the arg max as the Q function.

# Q Function

## Definition

$$\underline{V^*(s)} = \mathbb{E}[r|s, \pi^*(s)] + \gamma \mathbb{E}[V^*(s')|s, \pi^*(s)]$$

$$\underline{Q(s, a)} = \mathbb{E}[r|s, a] + \gamma \mathbb{E}[V^*(s')|s, a]$$

- If the agent knows  $Q$ , then the optimal action can be learned without  $\mathbb{P}\{s'|s, a\}$ .

$$\pi^*(s) = \arg \max_a Q(s, a), V^*(s) = \max_a Q(s, a)$$

# Deterministic Q Learning

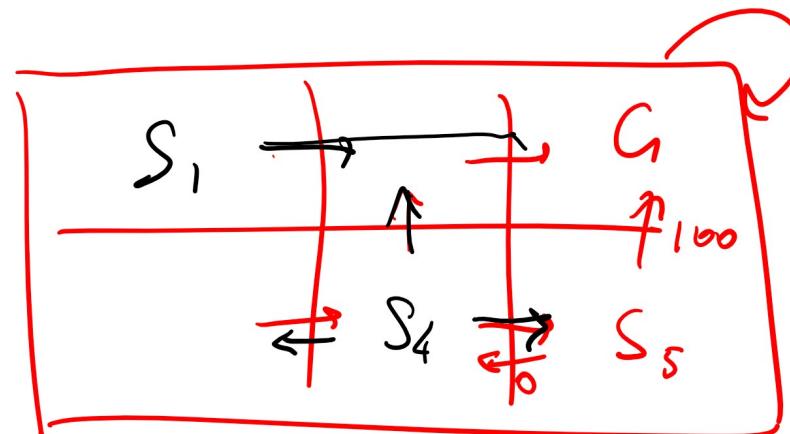
## Definition

- In the deterministic case,  $\mathbb{P}\{s'|s, a\}$  is either 0 or 1, the update formula for the  $Q$  function is the following.

$$\hat{Q}(s, a) = r + \gamma \max_{a'} \hat{Q}(s', a')$$

# Q Learning Example, Part I

## Definition



$$Q(S_1, \uparrow) = 100 + 0.9 \max_{a'} \{ Q(G, a') \} = 0$$

$$Q(S_1, \leftarrow) = 0 + 0.9 \max_{a'} \{ Q(S_4, a') \} = 0$$

$$Q(S_4, \rightarrow) = 0 + 0.9 \max_{a'} \{ Q(S_5, a') \}$$

$$\approx 90$$

initialised at 0

$$\max \{ Q(S_5, \leftarrow), Q(S_5, \uparrow) \}_{100}$$

Reinforcement Learning  
○○○○○○

Multi Armed Bandits  
○○○○○

Q-Learning  
○○○○○○○○○○●○○○○○○○

# Q Learning Example, Part II

## Definition

# Non-Deterministic Q Learning

## Definition

- In the nondeterministic case, the update formula for the  $Q$  function is the following.

$$\hat{Q}(s, a) = (1 - \alpha) \hat{Q}(s, a) + \alpha \left( r + \gamma \max_{a'} \hat{Q}(s', a') \right)$$

$$\alpha = \frac{1}{1 + \text{visits}(s, a)}$$

- Q learning will converge to the correct  $Q$  function in both deterministic and non-deterministic cases. In practice, it takes a very large number of iterations.

# Q Learning, Part I

## Algorithm

- Input: the state and reward processes.
- Output: optimal policy function  $\pi^*(s)$
- Initialize the Q table.

$$\hat{Q}(s, a) = 0, \text{ for each } s \in S, a \in A$$

# Q Learning, Part II

## Algorithm

- Observe current state  $s$ .
- Select an action  $a$  and execute it.
- Receive immediate reward  $r$ .
- Observe the new state  $s'$ .
- Update the table entry.

$$\hat{Q}(s, a) = (1 - \alpha) \hat{Q}(s, a) + \alpha \left( r + \gamma \max_{a'} \hat{Q}(s', a') \right)$$

$$\alpha = \frac{1}{1 + \text{visits } (s, a)}$$

- Update the state and repeat forever.

$$s = s'$$

# Exploration vs Exploitation

## Discussion

- There is a trade-off between learning about possibly better alternatives and following the current policy. Sometimes, random actions should be selected.

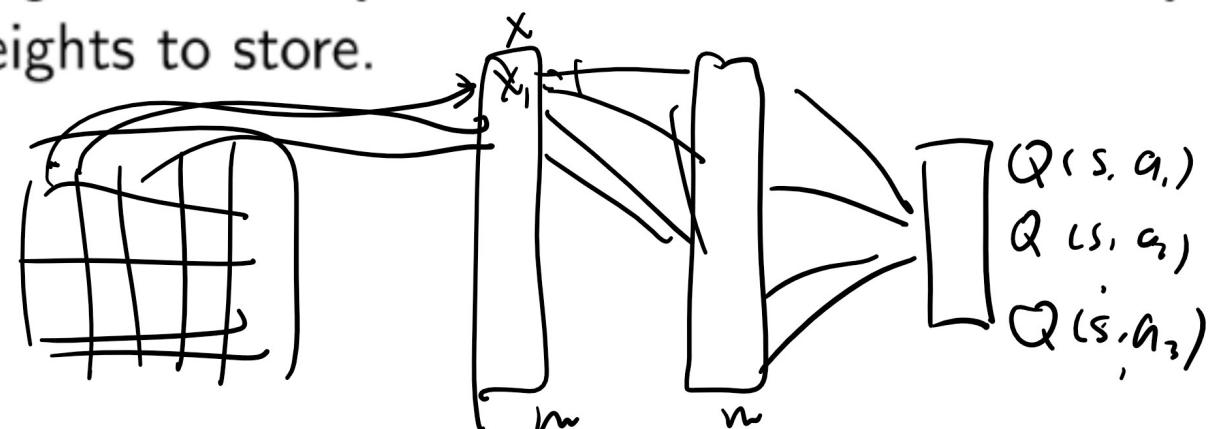
$$\mathbb{P}\{a|s\} = \frac{c^{\hat{Q}(s,a)}}{\sum_{a' \in A} c^{\hat{Q}(s,a')}}$$

- $c > 0$  is a constant that determines how strongly selection favors actions with higher Q values.

# Q Table vs Q Net

## Discussion

- In practice, Q table is too large to store since the number of possible states is very large.
- If there are  $m$  binary features that represent the state, the Q table contains  $\underline{2^m |A|}$ .
- However, it can be stored in a neural network called Q net.
- If there is a single hidden layer with  $m$  units, there are only  $m^2 + m |A|$  weights to store.



Reinforcement Learning  
○○○○○○

Multi Armed Bandits  
○○○○○

Q-Learning  
○○○○○○○○○○○○○○○○○○●○○

# Q Net Diagram

## Discussion

# Q Net Training

## Discussion

- Observe the features  $x$  given a state  $s$ .
- Apply action  $a$  and observe new state  $s'$  with features  $x'$  and reward  $r$ .
- Train the network with new instance  $(x, y)$

$$y = (1 - \alpha) \hat{y}(x, a) + \alpha \left( r + \gamma \max_{a'} \hat{y}(x', a') \right)$$

$\hat{y}(x, a)$                        $\hat{Q}(s, a')$   
 $\curvearrowleft$                        $\curvearrowright$   
 $\underbrace{\hspace{10em}}$                $\text{activation in } (x, a)$

- $\hat{y}(x, a)$  is the activation of output unit  $a$  given the input  $x$  in the current neural network.
- $\hat{y}(x', a')$  is the activation output unit  $a'$  given the input  $x'$  in the current neural network.

# Multi-Agent Learning

## Discussion

- Value function and policy function iteration methods can be applied to solve dynamic games with multiple agents.
- It will be used again in game theory in Week 11.