

# Artificial Intelligence

CS 165A

May 19, 2022

Instructor: Prof. Yu-Xiang Wang

T  
o  
d  
a  
y

→ Reinforcement Learning

# Recap: Tabular MDP

- **Discrete** State, **Discrete** Action, Reward and Observation

$$S_t \in \mathcal{S} \quad A_t \in \mathcal{A} \quad R_t \in \mathbb{R} \quad \text{--- } O_t \in \mathcal{O}$$

- Policy:

– When the state is observable:  $\pi : \mathcal{S} \rightarrow \mathcal{A}$

~~– Or when the state is not observable~~

$$\text{--- } \pi_t : (\mathcal{O} \times \mathcal{A} \times \mathbb{R})^{t-1} \rightarrow \mathcal{A}$$

- Learn the best policy that maximizes the expected reward

– Finite horizon (episodic) RL:  $\pi^* = \arg \max_{\pi \in \Pi} \mathbb{E} \left[ \sum_{t=1}^T R_t \right]$  **T: horizon**

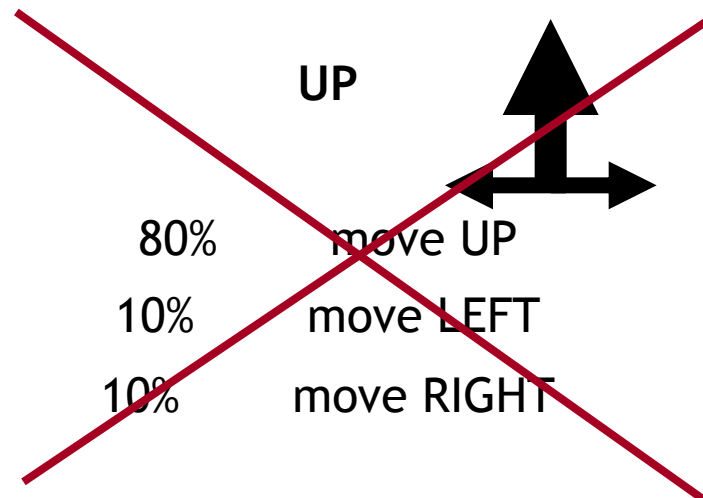
– Infinite horizon RL:  $\pi^* = \arg \max_{\pi \in \Pi} \mathbb{E} \left[ \sum_{t=1}^{\infty} \gamma^{t-1} R_t \right]$

**$\gamma$ : discount factor**

# Recap: Example: Frozen Lake

			<del>+1</del>
			<del>-1</del>
START			

Action 1, Action 2, Action 3, Action 4  
actions. ~~UP, DOWN, LEFT, RIGHT~~



- ~~• reward +1 at [4,3], -1 at [4,2]~~
- ~~• reward -0.04 for each step~~
- what's the strategy to achieve max reward?

Recap: Plug in the estimated MDP parameters and then do offline planning.

$$V_{k+1}^{\pi}(s) \leftarrow \sum_a \pi(a|s) \sum_{s'} \hat{P}(s'|s, a) [\hat{r}(s, a, s') + \gamma V_k^{\pi}(s')]$$

$$\pi' \leftarrow \arg \max_a \sum_{s'} \hat{P}(s'|s, a) [\hat{r}(s, a, s') + \gamma V_k^{\pi}(s')]$$

$$V_{k+1}(s) \leftarrow \max_a \sum_{s'} \hat{P}(s'|s, a) [\hat{r}(s, a, s') + \gamma V_k(s')]$$

Recap: Plug in the estimated MDP parameters and then do offline planning.

$$V_{k+1}^\pi(s) \leftarrow \sum_a \pi(a|s) \sum_{s'} \hat{P}(s'|s, a) [\hat{r}(s, a, s') + \gamma V_k^\pi(s')]$$

$$\pi' \leftarrow \arg \max_a \sum_{s'} \hat{P}(s'|s, a) [\hat{r}(s, a, s') + \gamma V_k^\pi(s')]$$

$$V_{k+1}(s) \leftarrow \max_a \sum_{s'} \hat{P}(s'|s, a) [\hat{r}(s, a, s') + \gamma V_k(s')]$$

Recap: Plug in the estimated MDP parameters and then do offline planning.

$$V_{k+1}^\pi(s) \leftarrow \sum_a \pi(a|s) \sum_{s'} \hat{P}(s'|s, a) [\hat{r}(s, a, s') + \gamma V_k^\pi(s')]$$
$$\pi' \leftarrow \arg \max_a \sum_{s'} \hat{P}(s'|s, a) [\hat{r}(s, a, s') + \gamma V_k^\pi(s')]$$
$$V_{k+1}(s) \leftarrow \max_a \sum_{s'} \hat{P}(s'|s, a) [\hat{r}(s, a, s') + \gamma V_k(s')] \leftarrow U$$

\* Note the “**hat**”. Usually it indicates empirical estimates.

Recap: Plug in the estimated MDP parameters and then do offline planning.

$$V_{k+1}^{\pi}(s) \leftarrow \sum_a \pi(a|s) \sum_{s'} \hat{P}(s'|s, a) [\hat{r}(s, a, s')] + \gamma V_k^{\pi}(s')$$

$$\pi' \leftarrow \arg \max_a \sum_{s'} \hat{P}(s'|s, a) [\hat{r}(s, a, s')] + \gamma V_k^{\pi}(s')$$

$$V_{k+1}(s) \leftarrow \max_a \sum_{s'} \hat{P}(s'|s, a) [\hat{r}(s, a, s')] + \gamma V_k(s')$$

\* Note the “**hat**”. Usually it indicates empirical estimates.

\* These iterations will produce  $\hat{V}^*$  and  $\hat{Q}^*$  functions, and then  $\hat{\pi}^*$


$$\hat{\pi}^*(s) = \arg \max_a Q^*(s, a)$$

# Recap: Three ideas for solving RL

- Idea 1: Model-based approach
  - Estimated the CPTs of MDP by their empirical frequency
  - Plug-in the estimate to Bellman equations for VI / PI
- Idea 2: Model-free approach: Directly estimate  $V$  function and  $Q$  function
  - Monte Carlo: Run many episodes of the MDP
  - First-visit MC: first time you visit State  $s$ , keep all subsequent rewards, then average over many such episodes
- Idea 3: Better model-free approach: Combining MC with VI / PI directly.
  - Temporal difference (TD) learning



# This lecture

- Monte Carlo Methods
  - Temporal Difference Learning
  - Features and linear function approximation
  - Policy Gradient method
- 

# Idea 2: Model-free Reinforcement Learning

- Do we need the model? Can we learn the Q function directly?
  - **How many free parameters are there to represent the Q-function?**

$$Q^* \in \mathbb{R}^{S \times A}$$

$$\hat{P} \rightarrow P$$
$$\mathbb{R}^{S \times S \times A}$$
$$P(s'|sa)$$

- **Ans: SA  $\ll$  O(S<sup>2</sup>A)**

- Recall: Policy iterations

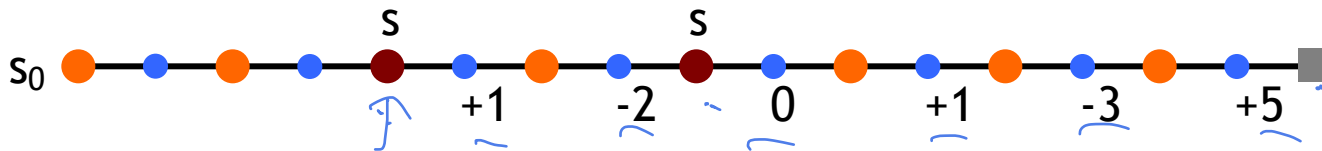
$$\pi_0 \xrightarrow{E} V^{\pi_0} \xrightarrow{I} \pi_1 \xrightarrow{E} V^{\pi_1} \xrightarrow{I} \dots \xrightarrow{I} \pi^* \xrightarrow{E} V^*$$

- **Maybe we can do policy evaluation without estimating the model?**

# Monte Carlo Policy Evaluation (Prediction)

- want to estimate  $V^\pi(s) = \mathbb{E}_\pi \left[ \sum_{t=1}^{\infty} \gamma^{t-1} R_t | s \right]$ 
  - = expected return starting from  $s$  and following  $\pi$
  - estimate as average of observed returns in state  $s$
- We can execute the policy  $\pi$
- first-visit MC
  - average returns following the first visit to state  $s$

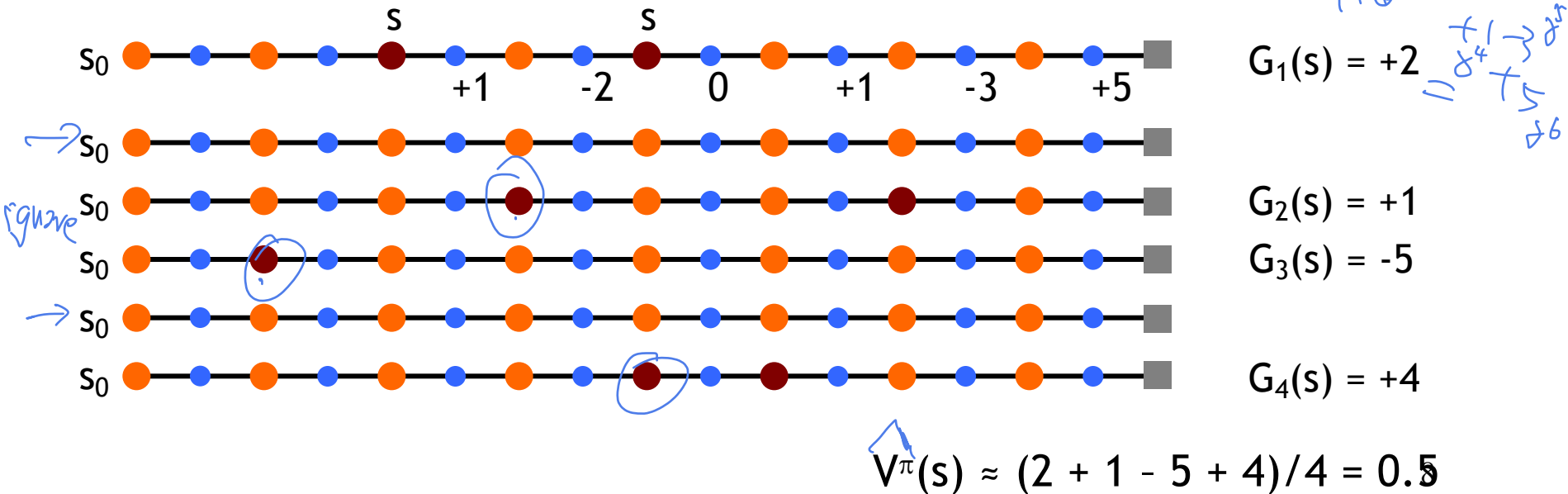
$$\approx \frac{1}{\sum_{i=1}^T \mathbb{1}(S_i = s)} \sum_{i=1}^T \sum_{t=1}^{\infty} \gamma^{t-1} R_t^{(i)} \mathbb{1}(S_t = s)$$



$$\underline{G_1(s) = +2}$$

# Monte Carlo Policy Evaluation (Prediction)

- want to estimate  $V^\pi(s)$ 
  - = expected return starting from  $s$  and following  $\pi$
  - estimate as average of observed returns in state  $s$
- We can execute the policy  $\pi$
- first-visit MC
  - average returns following the first visit to state  $s$



# Monte Carlo Policy Optimization (Control)

$$V^\pi = \sum \pi(a|s) Q^\pi(s,a)$$

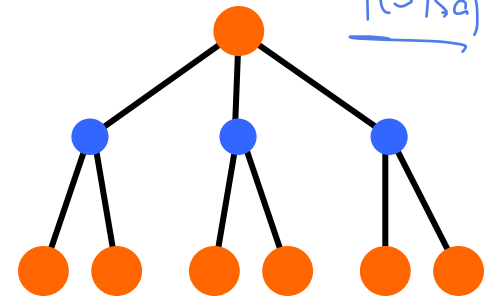
$$Q^\pi(s,a) = \mathbb{E} \left[ \sum_{t=0}^{\infty} \gamma^t r_t \mid s, a \right]$$

$s' \sim P(s'|s,a)$  &  $V^\pi(s)$

- $V^\pi$  not enough for policy improvement
  - need exact model of environment

- estimate  $Q^\pi(s,a)$  *first visit MC for  $Q^\pi$* 

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



- MC control

$$\pi_0 \xrightarrow{E} Q^{\pi_0} \xrightarrow{I} \pi_1 \xrightarrow{E} Q^{\pi_1} \xrightarrow{I} \dots \xrightarrow{I} \pi^* \xrightarrow{E} Q^*$$

- update after each episode

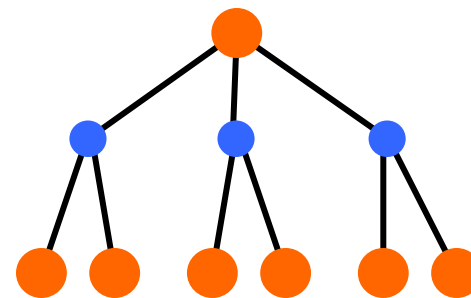
- Two problems

- greedy policy won't explore all actions
- Requires many independent episodes for the estimated value function to be accurate.

# Monte Carlo Policy Optimization (Control)

- $V^\pi$  not enough for policy improvement
  - need exact model of environment

- estimate  $Q^\pi(s,a)$   
 $\pi'(s) = \arg \max_a Q^\pi(s, a)$



- MC control

$$\pi_0 \xrightarrow{E} Q^{\pi_0} \xrightarrow{I} \pi_1 \xrightarrow{E} Q^{\pi_1} \xrightarrow{I} \dots \xrightarrow{I} \pi^* \xrightarrow{E} Q^*$$

- update after each episode

- Two problems

- greedy policy won't explore all actions **eps-greedy!**
- Requires many independent episodes for the estimated value function to be accurate.

# Improved Monte-Carlo Q-function estimate using Bellman equations

- Recall:

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

$$Q^\pi(s, a) = \underbrace{r^\pi(s, a)} + \gamma \mathbb{E}_{s' \sim \underbrace{P(s'|s, a)}} [\underbrace{V^\pi(s')}_{\text{MC estimate}}]$$

# Improved Monte-Carlo Q-function estimate using Bellman equations

- Recall:

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

$$Q^\pi(s, a) = r^\pi(s, a) + \gamma \mathbb{E}_{s' \sim P(s'|s, a)}[V^\pi(s')]$$

- We can use the empirical (Monte Carlo) estimate.

$$\widehat{Q}^\pi(s, a) = \widehat{r}^\pi(s, a) + \gamma \widehat{\mathbb{E}}_{s' \sim P(s'|s, a)}[\widehat{V}^\pi(s')]$$



# Improved Monte-Carlo Q-function estimate using Bellman equations

- Recall:

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

$$Q^\pi(s, a) = r^\pi(s, a) + \gamma \mathbb{E}_{s' \sim P(s'|s, a)} [V^\pi(s')]$$

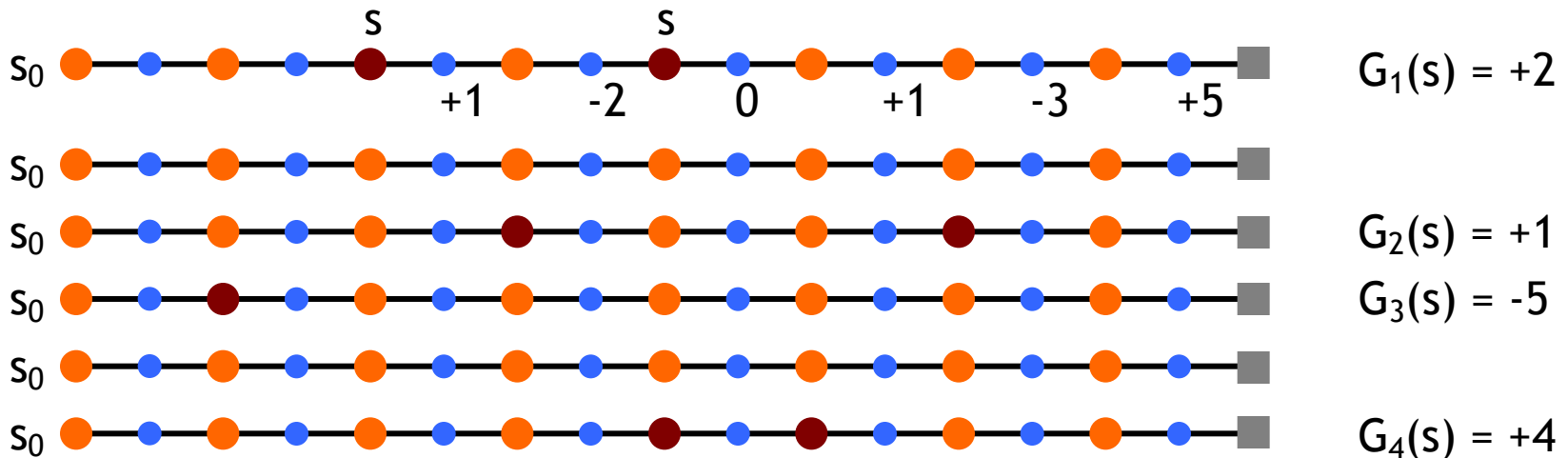
- We can use the empirical (Monte Carlo) estimate.

$$\widehat{Q}^\pi(s, a) = \widehat{r}^\pi(s, a) + \gamma \widehat{\mathbb{E}}_{s' \sim P(s'|s, a)} [\widehat{V}^\pi(s')]$$

\*No need to estimate  $P(s' | s, a)$  or  $r(s, a, s')$  as intermediate steps.

\*Require only  $O(SA)$  space, rather than  $O(S^2A)$

# Online averaging representation of MC



$$V^\pi(s) \approx (2 + 1 - 5 + 4) / 4 = 0.5$$

$$\bar{X}_t = \frac{1}{t} \sum_{i=1}^t X_i \quad \text{receive } X_{t+1}$$

$$\bar{X}_{t+1} = \frac{1}{t+1} \sum_{i=1}^{t+1} X_i = \frac{1}{t+1} (\cancel{t} \bar{X}_t + X_{t+1}) = \bar{X}_t + \frac{1}{t+1} X_{t+1} - \frac{1}{t+1} \bar{X}_t$$

$$\frac{1}{t+1} (X_{t+1} - \bar{X}_t)$$

- Alternative, *online averaging* update

$$V(S_t) \leftarrow V(S_t) + \alpha [G_t - V(S_t)], \quad \text{where } \alpha = 1/N_{S_t}$$

# DP + MC = Temporal Difference Learning

- Monte Carlo

$$V(S_t) \leftarrow V(S_t) + \alpha [G_t - V(S_t)],$$

# DP + MC = Temporal Difference Learning

- Monte Carlo

$$V(S_t) \leftarrow V(S_t) + \alpha [G_t - V(S_t)],$$

Issue:  $G_t$  can only be obtained after the entire episode!

# DP + MC = Temporal Difference Learning

- Monte Carlo

$$V(S_t) \leftarrow V(S_t) + \alpha [G_t - V(S_t)],$$

Issue:  $G_t$  can only be obtained after the entire episode!

- The idea of TD learning:

$$\mathbb{E}_\pi [G_t] = \mathbb{E}_\pi [R_t | S_t] + \gamma V^\pi(S_{t+1})$$

# DP + MC = Temporal Difference Learning

- Monte Carlo

$$V(S_t) \leftarrow V(S_t) + \alpha [G_t - V(S_t)],$$

Issue:  $G_t$  can only be obtained after the entire episode!

- The idea of TD learning:

$$\mathbb{E}_\pi [G_t] = \mathbb{E}_\pi [R_t | S_t] + \gamma V^\pi(S_{t+1})$$

We only need one step before we can plug-in and estimate the RHS!

# DP + MC = Temporal Difference Learning

- Monte Carlo

$$V(S_t) \leftarrow V(S_t) + \alpha [G_t - V(S_t)],$$

Issue:  $G_t$  can only be obtained after the entire episode!

- The idea of TD learning:

$$\mathbb{E}_\pi [G_t] = \mathbb{E}_\pi [\underline{R_t} | S_t] + \gamma \underline{V^\pi(S_{t+1})}$$

We only need one step before we can plug-in and estimate the RHS!

- TD-Policy evaluation

$(S_t, \pi(S_t), S_{t+1}, R_t) \leftarrow$  Data point

$$\underline{V(S_t)} \leftarrow \underline{V(S_t)} + \alpha \left[ \underline{R_t} + \gamma \underline{V(S_{t+1})} - \underline{V(S_t)} \right]$$

Realized at time  $t+1$

your estimate of at time  $t$

# DP + MC = Temporal Difference Learning

- Monte Carlo

$$V(S_t) \leftarrow V(S_t) + \alpha [G_t - V(S_t)],$$

Issue:  $G_t$  can only be obtained after the entire episode!

- The idea of TD learning:

$$\mathbb{E}_\pi [G_t] = \mathbb{E}_\pi [R_t | S_t] + \gamma \underbrace{V^\pi}_{\text{bootstrapping}}(S_{t+1})$$

We only need one step before we can plug-in and estimate the RHS!

- TD-Policy evaluation

$$V(S_t) \leftarrow V(S_t) + \alpha [R_{t+1} + \gamma \underbrace{V(S_{t+1})}_{\text{bootstrapping}} - V(S_t)]$$

**Bootstrapping!**



# Bootstrap's origin

- “The Surprising Adventures of Baron Münchhausen”
  - Rudolf Erich Raspe, 1785



**PULL  
YOURSELF  
UP BY  
THE  
BOOT  
STRAPS!!!**



- In statistics: Brad Efron's resampling methods
- In computing: Booting...
- In RL: It simply means TD learning

# TD policy optimization (TD-control)

- SARSA (On-Policy TD-control)

- Update the Q function by bootstrapping Bellman Equation

$$Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma Q(S', A') - Q(S, A)]$$

- Choose the next  $A'$  using Q, e.g., eps-greedy.

- Q-Learning (Off-policy TD-control)

- Update the Q function by bootstrapping Bellman Optimality Eq.

$$Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma \max_a Q(S', a) - Q(S, A)]$$

- Choose the next  $A'$  using Q, e.g., eps-greedy, or any other policy.

## Remarks:

- These are **proven to converge** asymptotically.
- Much more data-efficient in practice, than MC.
- Regret analysis is still active area of research.

# Advantage of TD over Monte Carlo

- Given a trajectory, a roll-out, of  $T$  steps.
  - MC updates the  $Q$  function only once
  - TD updates the  $Q$  function (and the policy)  $T$  times!

# Advantage of TD over Monte Carlo

- Given a trajectory, a roll-out, of  $T$  steps.
  - MC updates the  $Q$  function only once
  - TD updates the  $Q$  function (and the policy)  $T$  times!

**Remark:** This is the same kind of improvement from Gradient Descent to Stochastic Gradient Descent (SGD).

# Advantage of TD over Monte Carlo

- Given a trajectory, a roll-out, of  $T$  steps.
  - MC updates the  $Q$  function only once
  - TD updates the  $Q$  function (and the policy)  $T$  times!

**Remark:** This is the same kind of improvement from Gradient Descent to Stochastic Gradient Descent (SGD).

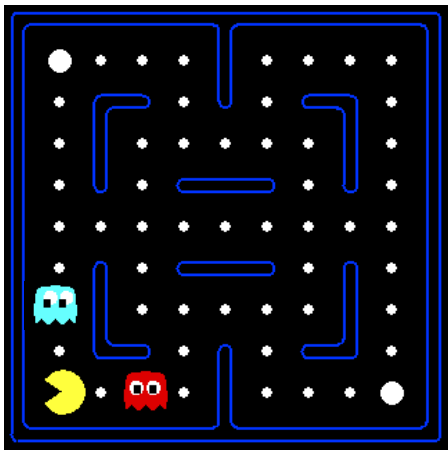
(Optional reading: Sutton and Barton 9.3 Semi-gradient methods)

# The problem of large state-space is still there

- We need to represent and learn SA parameters in Q-learning and SARSA.
- S is often large
  - 9-puzzle, Tic-Tac-Toe:  $9! = 362,800$ ,  $S^2 = 1.3 \cdot 10^{11}$
  - PACMAN with 20 by 20 grid.  $S = O(2^{400})$ ,  $S^2 = O(2^{800})$
- $O(S)$  is not acceptable in some cases.
- Need to think of ways to “generalize”/share information across states.

# Example: Pacman

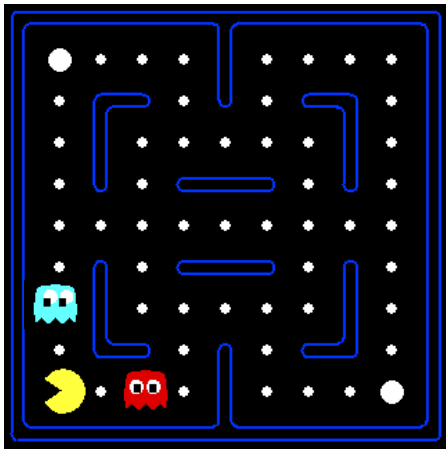
Let's say we discover  
through experience  
that this state is bad:



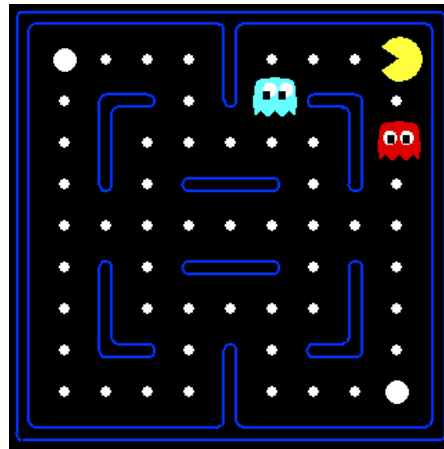
(From Dan Klein and Pieter Abbeel)

# Example: Pacman

Let's say we discover through experience that this state is bad:



In naïve q-learning, we know nothing about this state:

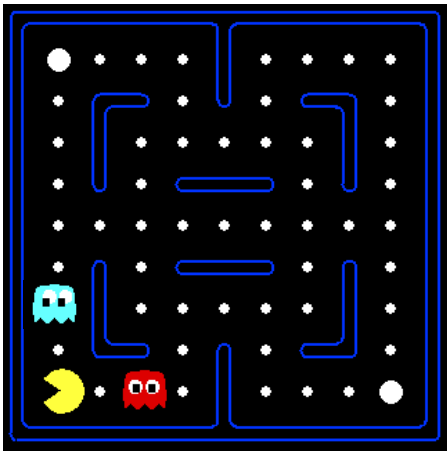


(From Dan Klein and Pieter Abbeel)



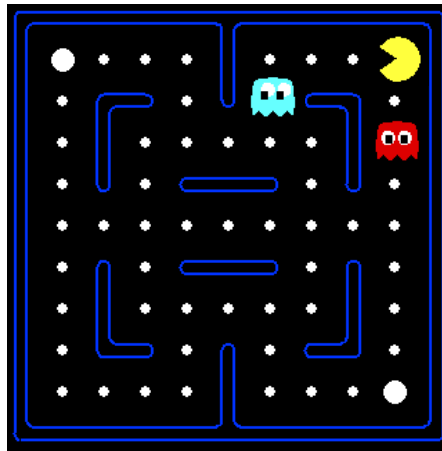
# Example: Pacman

Let's say we discover through experience that this state is bad:

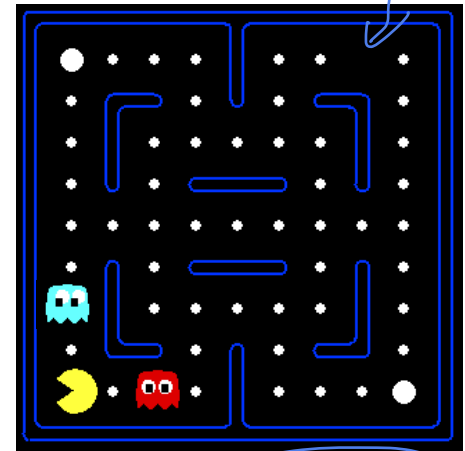


$V(s) = -100$

In naïve q-learning, we know nothing about this state:



Or even this one!



$V(s') = ?$

(From Dan Klein and Pieter Abbeel)

# Video of Demo Q-Learning Pacman – Tiny – Watch All



# Video of Demo Q-Learning Pacman – Tiny – Silent Train



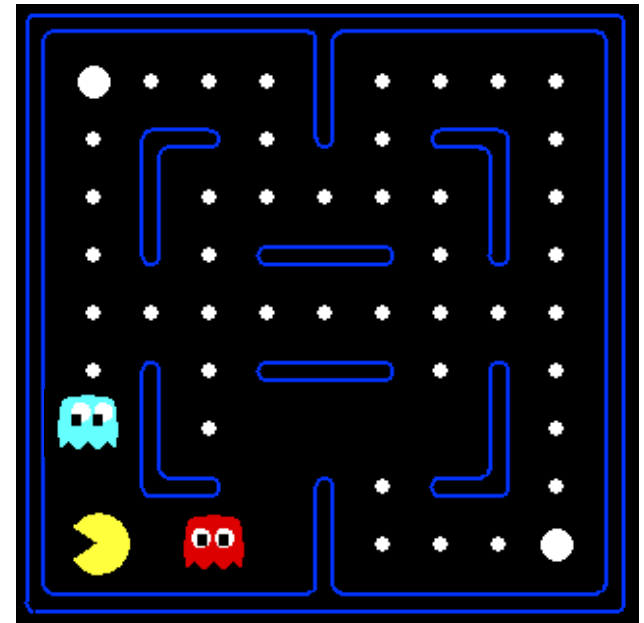
# Video of Demo Q-Learning Pacman – Tricky – Watch All



# Why not use an evaluation function?

## A Feature-Based Representations

- Solution: describe a state using a vector of features (properties)
  - Features are functions from states to real numbers (often 0/1) that capture important properties of the state
  - Example features:
    - Distance to closest ghost
    - Distance to closest dot
    - Number of ghosts
    - $1 / (\text{dist to dot})^2$
    - Is Pacman in a tunnel? (0/1)
    - ..... etc.
    - Is it the exact state on this slide?
  - Can also describe a q-state (s, a) with features (e.g. action moves closer to food)



# Linear Value Functions

- Using a feature representation, we can write a q function (or value function) for any state using a few weights:

$$- V_{\mathbf{w}}(s) = \underline{w_1} \underline{f_1}(s) + \underline{w_2} \underline{f_2}(s) + \dots + \underline{w_n} \underline{f_n}(s)$$

$$- \underline{Q_{\mathbf{w}}(s,a)} = \underline{w_1} \underline{f_1}(s,a) + \underline{w_2} \underline{f_2}(s,a) + \dots + \underline{w_n} \underline{f_n}(s,a)$$

- Advantage: our experience is summed up in a few powerful numbers

- Disadvantage: states may share features but actually be very different in value!

Handwritten notes illustrating the linear value function:

$$Q_{\mathbf{w}}(s,a) = \begin{pmatrix} w_{\phi_1}(s) \\ f_1 \\ \vdots \\ f_n \end{pmatrix} \cdot \begin{pmatrix} w_1 \\ \vdots \\ w_n \end{pmatrix}$$

Additional handwritten notes include:

- $Q_{\mathbf{w}}(s)$  (underlined)
- $\mathbb{R}^n$  (in a box)
- $\mathbf{w} = \begin{pmatrix} w_1 \\ \vdots \\ w_n \end{pmatrix}$
- $w_{\phi_i}(s)$
- $w_i \phi_i^s$  (underlined)

# Updating a linear value function

# Updating a linear value function

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

- Original Q learning rule tries to reduce prediction error at  $s, a$ :

$$Q(s,a) \leftarrow Q(s,a) + \alpha \cdot [R(s,a,s') + \gamma \max_{a'} Q(s',a') - Q(s,a)]$$



# Updating a linear value function

$$\underline{\langle a, b \rangle = a^T b = \sum_{i=1}^n a_i b_i}$$

- Original Q learning rule tries to reduce prediction error at  $s, a$ :

$$Q(s,a) \square Q(s,a) + \alpha \cdot [R(s,a,s') + \gamma \max_{a'} \underline{Q}(s',a') - \underline{Q}(s,a) ]$$

$$\begin{array}{c} \parallel \\ \langle w, f(s,a) \rangle \end{array} \quad \begin{array}{c} \parallel \\ \langle w, f(s,a) \rangle \end{array}$$

- Instead, we update the weights to try to reduce the error at  $s, a$ :

# Updating a linear value function

- Original Q learning rule tries to reduce prediction error at  $s, a$ :

$$Q(s,a) \leftarrow Q(s,a) + \alpha \cdot [R(s,a,s') + \gamma \max_{a'} Q(s',a') - Q(s,a)]$$

- Instead, we update the weights to try to reduce the error at  $s, a$ :

$$\begin{aligned} \underline{w_i} &\leftarrow w_i + \alpha \cdot [R(s,a,s') + \gamma \max_{a'} Q(s',a') - Q(s,a)] \leftarrow \underline{Q_w(s,a)} / \partial w_i \\ &= \underline{w_i} + \alpha \cdot \underline{[R(s,a,s') + \gamma \max_{a'} Q(s',a') - Q(s,a)]} \underline{f_i(s,a)} \end{aligned}$$

# Updating a linear value function

- Original Q learning rule tries to reduce prediction error at  $s, a$ :

$$Q(s,a) \leftarrow Q(s,a) + \alpha \cdot [R(s,a,s') + \gamma \max_{a'} Q(s',a') - Q(s,a)]$$

- Instead, we update the weights to try to reduce the error at  $s, a$ :

$$\begin{aligned} w_i &\leftarrow w_i + \alpha \cdot [R(s,a,s') + \gamma \max_{a'} Q(s',a') - Q(s,a)] \partial Q_w(s,a) / \partial w_i \\ &= w_i + \alpha \cdot [R(s,a,s') + \gamma \max_{a'} Q(s',a') - Q(s,a)] f_i(s,a) \end{aligned}$$

- Qualitative justification:
  - Pleasant surprise: increase weights on positive features, decrease on negative ones
  - Unpleasant surprise: decrease weights on positive features, increase on negative ones

# Q-Learning with function approximation

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

# Q-Learning with function approximation

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

- Q-learning with linear Q-functions:

# Q-Learning with function approximation

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

- Q-learning with linear Q-functions:

$$\text{transition} = (s, a, r, s')$$

# Q-Learning with function approximation

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

- Q-learning with linear Q-functions:

transition =  $(s, a, r, s')$

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

# Q-Learning with function approximation

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

- Q-learning with linear Q-functions:

$$\text{transition} = (s, a, r, s')$$

$$\text{difference} = \left[ r + \gamma \max_{a'} Q(s', a') \right] - Q(s, a)$$

$$Q(s, a) \leftarrow Q(s, a) + \alpha [\text{difference}] \quad \text{Exact Q's}$$



# Q-Learning with function approximation

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

- Q-learning with linear Q-functions:

$$\text{transition} = (s, a, r, s')$$

$$\text{difference} = \left[ r + \gamma \max_{a'} Q(s', a') \right] - Q(s, a)$$

$$Q(s, a) \leftarrow Q(s, a) + \alpha [\text{difference}] \quad \text{Exact Q's}$$

$$\underline{w}_i \leftarrow \underline{w}_i + \alpha [\text{difference}] f_i(s, a) \quad \text{Approximate Q's}$$

Handwritten note:  $f_i(s, a) =$

the number

# Q-Learning with function approximation

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

- Q-learning with linear Q-functions:

$$\text{transition} = (s, a, r, s')$$

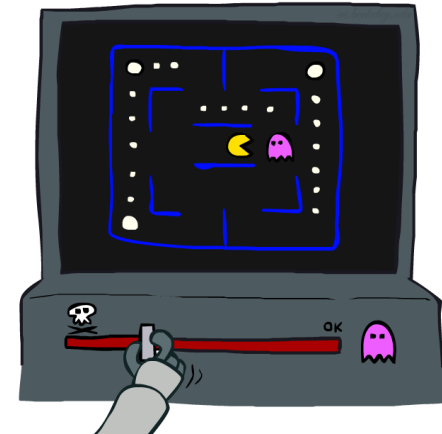
$$\text{difference} = \left[ r + \gamma \max_{a'} Q(s', a') \right] - Q(s, a)$$

$$Q(s, a) \leftarrow Q(s, a) + \alpha [\text{difference}] \quad \text{Exact Q's}$$

$$w_i \leftarrow w_i + \alpha [\text{difference}] f_i(s, a) \quad \text{Approximate Q's}$$

- Intuitive interpretation:

- Adjust weights of active features
- E.g., if something unexpectedly bad happens, blame the features that were on: disprefer all states with that state's features



# Q-Learning with function approximation

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

- Q-learning with linear Q-functions:

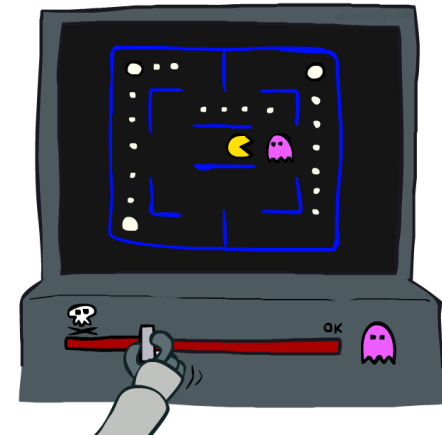
$$\text{transition} = (s, a, r, s')$$

$$\text{difference} = \left[ r + \gamma \max_{a'} Q(s', a') \right] - Q(s, a)$$

$$Q(s, a) \leftarrow Q(s, a) + \alpha [\text{difference}] \quad \text{Exact Q's}$$

$$w_i \leftarrow w_i + \alpha [\text{difference}] f_i(s, a) \quad \text{Approximate Q's}$$

- Intuitive interpretation:
  - Adjust weights of active features
  - E.g., if something unexpectedly bad happens, blame the features that were on: disprefer all states with that state's features
- Formal justification: online least squares (Read the textbook!)



# PACMAN Q-Learning (Linear function approx.)



# So far, in RL algorithms

- Model-based approaches
  - Estimate the MDP parameters.
  - Then use policy-iterations, value iterations.
- Monte Carlo methods:
  - estimating the rewards by empirical averages
- Temporal Difference methods:
  - Combine Monte Carlo methods with Dynamic Programming
- Linear function approximation in Q-learning
  - Similar to SGD
  - Learning heuristic function

# So far, in RL algorithms

- Model-based approaches
  - Estimate the MDP parameters.
  - Then use policy-iterations, value iterations.
- Monte Carlo methods:
  - estimating the rewards by empirical averages
- Temporal Difference methods:
  - Combine Monte Carlo methods with Dynamic Programming
- Linear function approximation in Q-learning
  - Similar to SGD
  - Learning heuristic function

\*Question: What is the policy class  $\Pi$  of interest in these methods?

# Policy gradient

Softmax  $\left( \sum_i \exp(w_i \phi(s, a_i)) \right)$   
 (arg)

$\left. \begin{matrix} = 1 \\ \sum_i w_i \phi(s, a_i) \end{matrix} \right]$

- Let's not worry about states, dynamics, Q function.
  - We might not even observe the true state.
  - Let's specify a class of parametrized policy and hope to compare to the best within this class

- Objective function to maximize:  $J(\theta) \doteq \underline{v}_{\pi_{\theta}}(s_0)$ ,  $\stackrel{=}{=} E_{\pi_{\theta}} [R_1 + \gamma R_2 + \gamma^2 R_3 + \dots | s_0]$

\*Note how this theorem is non-trivial... The first two terms depends on  $\pi$ , but we did not take the gradient w.r.t. them.

# Policy gradient

- Let's not worry about states, dynamics, Q function.
  - We might not even observe the true state.
  - Let's specify a class of parametrized policy and hope to compare to the best within this class

- Objective function to maximize:  $J(\boldsymbol{\theta}) \doteq v_{\pi_{\boldsymbol{\theta}}}(s_0),$

- Do SGD:  $\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t + \alpha \widehat{\nabla J(\boldsymbol{\theta}_t)},$

\*Note how this theorem is non-trivial... The first two terms depends on  $\pi$ , but we did not take the gradient w.r.t. them.



# Policy gradient

- Let's not worry about states, dynamics, Q function.
  - We might not even observe the true state.
  - Let's specify a class of parametrized policy and hope to compare to the best within this class
- Objective function to maximize:  $J(\boldsymbol{\theta}) \doteq v_{\pi_{\boldsymbol{\theta}}}(s_0)$ ,
- Do SGD:  $\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t + \alpha \widehat{\nabla J(\boldsymbol{\theta}_t)}$ ,
- Policy gradient theorem:

\*Note how this theorem is non-trivial... The first two terms depends on  $\pi$ , but we did not take the gradient w.r.t. them.

# Policy gradient

$$V_{\pi_{\theta}}(s_0) = \sum_s d^{\pi_{\theta}}(s) \sum_a \pi_{\theta}(a|s) Q^{\pi_{\theta}}(s,a)$$

$$(xy)' = x'y + xy'$$

- Let's not worry about states, dynamics, Q function.
  - We might not even observe the true state.
  - Let's specify a class of parametrized policy and hope to compare to the best within this class

- Objective function to maximize:  $J(\theta) \doteq v_{\pi_{\theta}}(s_0)$ ,

- Do SGD:  $\theta_{t+1} = \theta_t + \alpha \widehat{\nabla J(\theta_t)}$ ,  
 $R^t + \gamma V^{\pi}(s')$  TD-estimate

- Policy gradient theorem:

$$\nabla J(\theta) = \sum_s d^{\pi}(s) \sum_a Q^{\pi}(s, a) \nabla_{\theta} \pi(a|s, \theta)$$

MC-estimate

\*Note how this theorem is non-trivial... The first two terms depends on  $\pi$ , but we did not take the gradient w.r.t. them.

# Stochastic approximation in policy gradients

$$\nabla J(\boldsymbol{\theta}) = \sum_s d^\pi(s) \sum_a Q^\pi(s, a) \nabla_{\boldsymbol{\theta}} \pi(a|s, \boldsymbol{\theta})$$

- Sample from running policy  $\pi$ 
  - $(S_1, A_1, R_1)$ , ...,  $(S_T, A_T, R_T)$

# Stochastic approximation in policy gradients

$$\nabla J(\boldsymbol{\theta}) = \sum_s \underbrace{d^\pi(s)} \sum_a \underbrace{Q^\pi(s, a)} \nabla_{\boldsymbol{\theta}} \pi(a|s, \boldsymbol{\theta})$$

- Sample from running policy  $\pi$ 
  - $(S_1, A_1, R_1), \dots, (S_T, A_T, R_T)$
- Idea: Sample  $s$ , then the following is an unbiased estimator (finite horizon episodic case)

$$\sum_{t=1}^T \left( \sum_{\ell=t}^T R_\ell \right) \frac{\nabla_{\boldsymbol{\theta}} \pi(A_t|S_t, \boldsymbol{\theta})}{\pi(A_t|S_t, \boldsymbol{\theta})}$$

# Stochastic approximation in policy gradients

$$\nabla J(\boldsymbol{\theta}) = \sum_s d^\pi(s) \sum_a Q^\pi(s, a) \nabla_{\boldsymbol{\theta}} \pi(a|s, \boldsymbol{\theta})$$

- Sample from running policy  $\pi$ 
  - $(S_1, A_1, R_1), \dots, (S_T, A_T, R_T)$
- Idea: Sample  $s$ , then the following is an unbiased estimator (finite horizon episodic case)

$$\begin{aligned} & \sum_{t=1}^T \left( \sum_{\ell=t}^T R_\ell \right) \frac{\nabla_{\boldsymbol{\theta}} \pi(A_t|S_t, \boldsymbol{\theta})}{\pi(A_t|S_t, \boldsymbol{\theta})} \\ &= \sum_{t=1}^T G_t \nabla_{\boldsymbol{\theta}} \log(\pi(A_t|S_t, \boldsymbol{\theta})) \end{aligned}$$

# Stochastic approximation in policy gradients

$$\nabla J(\boldsymbol{\theta}) = \sum_s d^\pi(s) \sum_a Q^\pi(s, a) \nabla_{\boldsymbol{\theta}} \pi(a|s, \boldsymbol{\theta})$$

- Sample from running policy  $\pi$ 
  - $(S_1, A_1, R_1), \dots, (S_T, A_T, R_T)$
- Idea: Sample  $s$ , then the following is an unbiased estimator (finite horizon episodic case)

$$\begin{aligned} & \sum_{t=1}^T \left( \sum_{\ell=t}^T R_\ell \right) \frac{\nabla_{\boldsymbol{\theta}} \pi(A_t|S_t, \boldsymbol{\theta})}{\pi(A_t|S_t, \boldsymbol{\theta})} \\ &= \sum_{t=1}^T G_t \nabla_{\boldsymbol{\theta}} \log(\pi(A_t|S_t, \boldsymbol{\theta})) \end{aligned}$$

**\*Show that this is an unbiased estimator of the gradient.**

# The REINFORCE algorithm (Williams, 1987)

REINFORCE, A Monte-Carlo Policy-Gradient Method (episodic)

Input: a differentiable policy parameterization  $\pi(a|s, \theta)$

Initialize policy parameter  $\theta \in \mathbb{R}^{d'}$

Repeat forever:

Generate an episode  $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T$ , following  $\pi(\cdot|\cdot, \theta)$

For each step of the episode  $t = 0, \dots, T - 1$ :

$G \leftarrow$  return from step  $t$

$\theta \leftarrow \theta + \alpha \gamma^t G \nabla_{\theta} \ln \pi(A_t|S_t, \theta)$

- From Sutton and Barto Ch .13.
- Note the  $\gamma^t$  term. This is for the discounted (episodic) case
- Updating the parameter T times for each episode!
- Easy to implement easy to understand from SGD theory.

# Elements of State-of-the-Art Reinforcement Learning

- Use a deep neural network to parameterize Q-function
- Use a deep neural network to parameterize the policy  $\pi$
- Run a combination of Q-learning and Policy Gradient.
  - Actor-Critics, A3C, etc...
- Heuristic-based exploration: curiosity, reward shaping, etc..
- Experience replay to generate more data from existing data.
- Multi-agent RL: modeling your opponents



# Example of State-of-the-Art RL for Hide-n-Seek

Lowe, Wu et al. (2017)

<https://proceedings.neurips.cc/paper/2017/hash/68a9750337a418a86fe06c1991a1d64c-Abstract.html>

# Example of State-of-the-Art RL for Hide-n-Seek



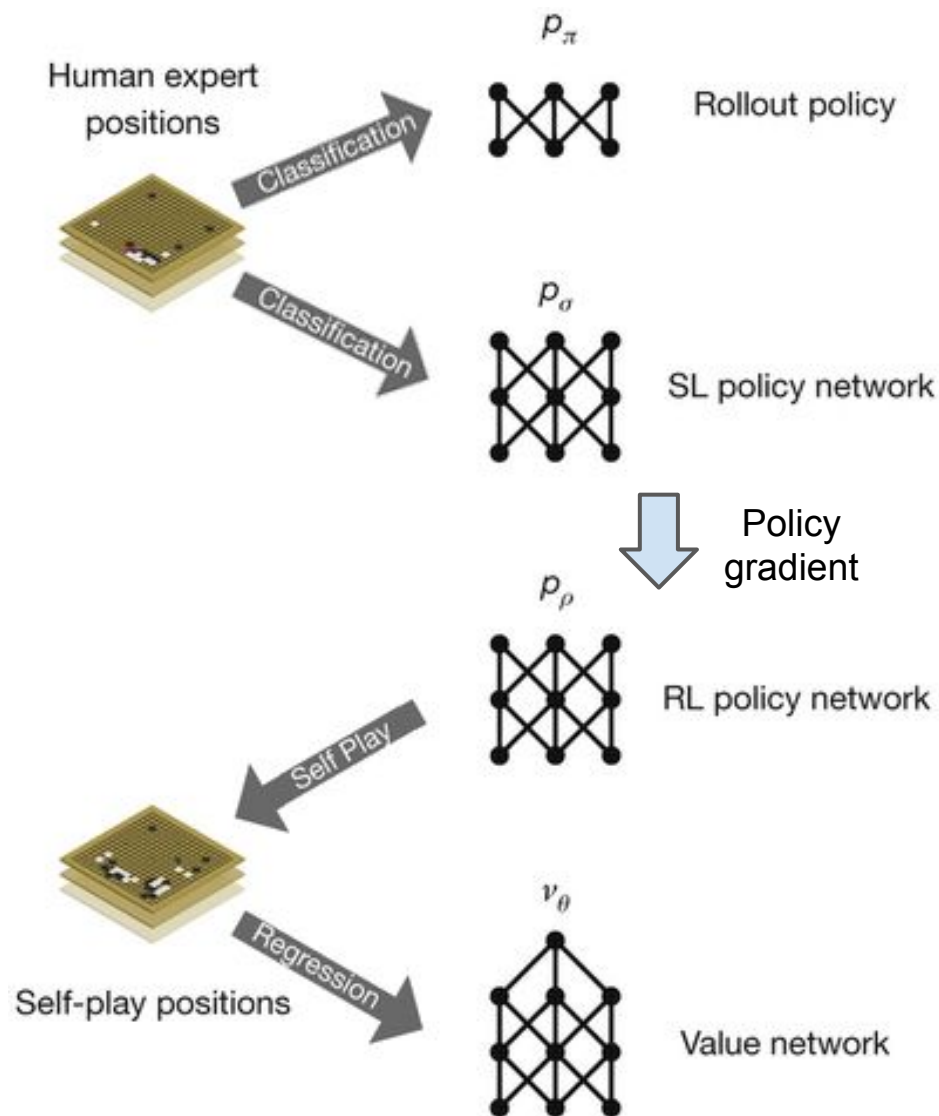
Lowe, Wu et al. (2017)

<https://proceedings.neurips.cc/paper/2017/hash/68a9750337a418a86fe06c1991a1d64c-Abstract.html>

# Alpha-Go and Alpha-Zero

- Parameterize the policy networks with CNN
- Supervised learning initialization
- RL using Policy gradient
- Fit Value Network (This is a heuristic function!)
- Monte-Carlo Tree Search

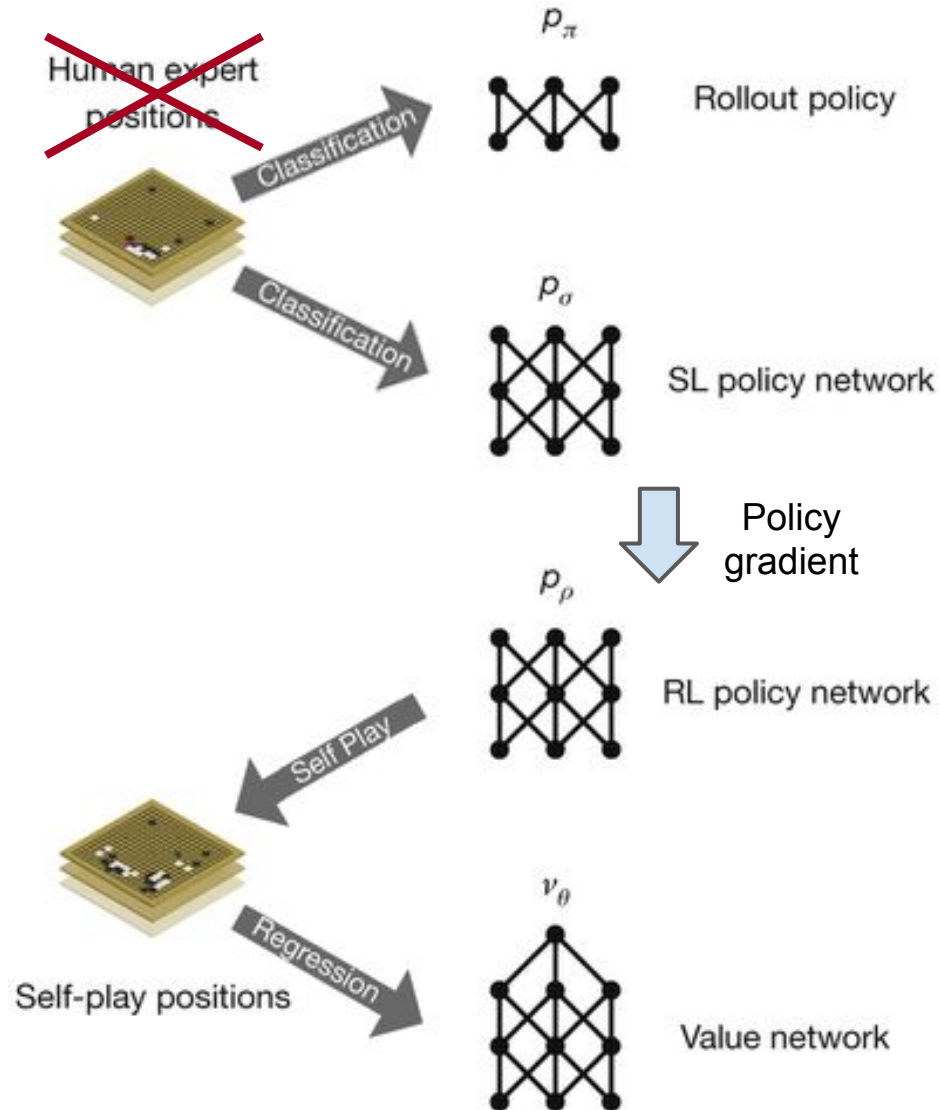
<https://www.youtube.com/watch?v=4D5yGiYe8p4>



# Alpha-Go and Alpha-Zero

- Parameterize the policy networks with CNN
- ~~• Supervised learning initialization~~
- RL using Policy gradient
- Fit Value Network (This is a heuristic function!)
- Monte-Carlo Tree Search

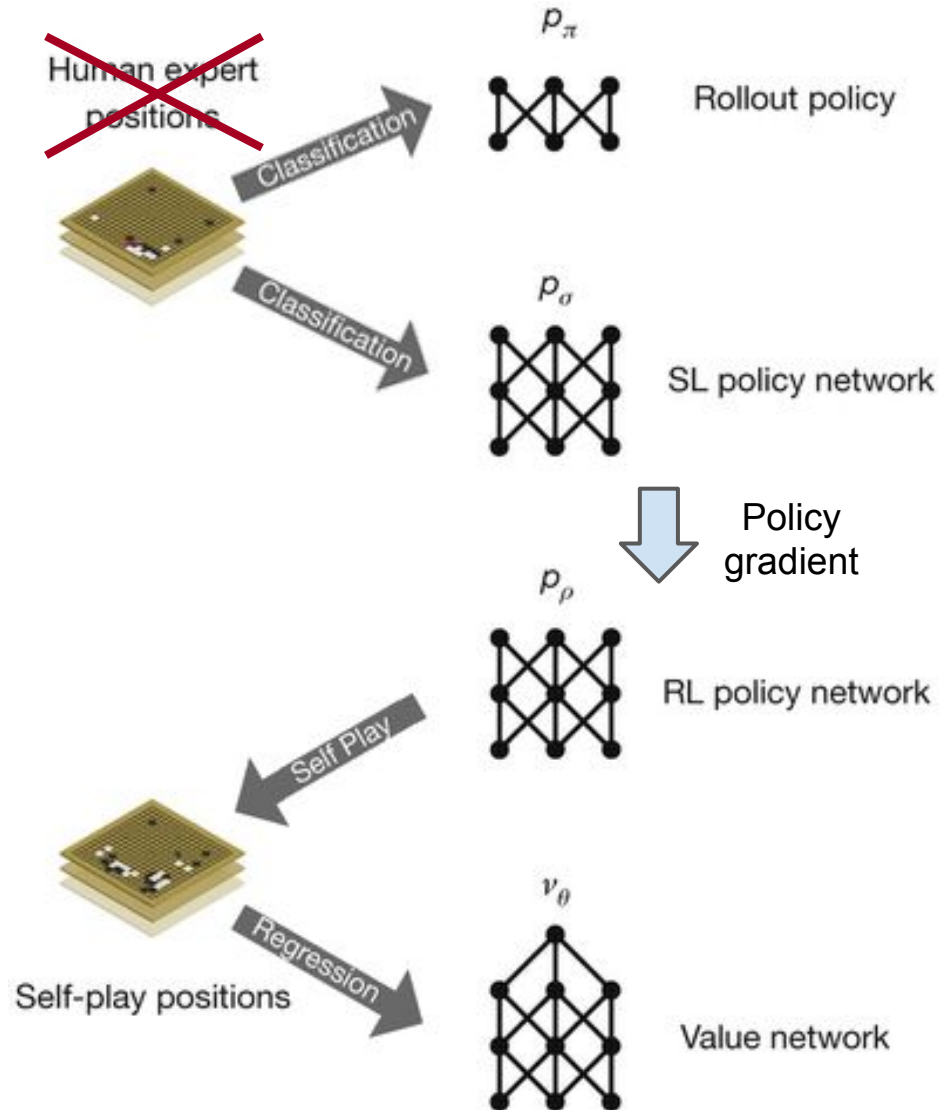
<https://www.youtube.com/watch?v=4D5yGiYe8p4>



# Alpha-Go and Alpha-Zero

- Parameterize the policy networks with CNN
- ~~• Supervised learning initialization~~
- RL using Policy gradient
- Fit Value Network (This is a heuristic function!)
- Monte-Carlo Tree Search

<https://www.youtube.com/watch?v=4D5yGiYe8p4>



# Summary of RL algorithms

- Model-based:
  - Policy iteration / Value iteration
  - Need to estimate the dynamics (MDP parameters)
- Model-free: (no need to “explicitly” estimate dynamics)
  - TD learning: SARSA, Q-learning
  - Function approximation (Share information across states)
- Absolutely model-free (do not even need an MDP model)
  - Policy gradient

# Next lecture

- Start Logical Agent
- Logical inference for propositional logic
- What you should do:
  - Read Chapter 7 of AIMA textbook.
  - Start working on Project 3 if you haven't yet.