

# CS 4803 / 7643

## Deep Learning, Fall 2020

Reinforcement Learning: Module 2/3

Presented by [Nirbhay Modhe](#)

## Previous Lecture

- **RL:** Definitions, interaction API, tasks/challenges
- **MDPs:** Theoretical framework underlying RL, solving MDPs

## Today

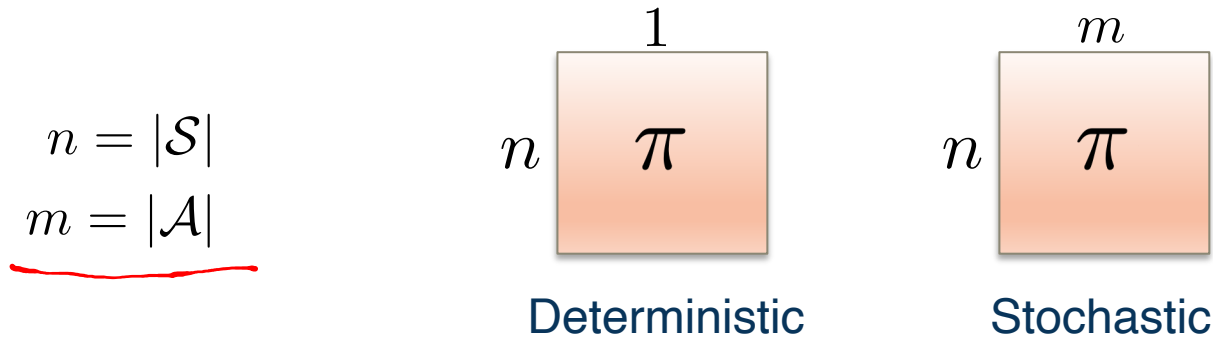
- **Policy** (continued): How an agents acts at states
- **Value function (Utility):** How good is a particular state or state-action pair?
- **Algorithms** for solving MDPs (Value Iteration)
- Departure from known rewards and transitions: Reinforcement Learning (RL), Deep RL

- Markov Decision Processes (MDPs)
  - States, Actions, Reward dist., Transition dist., Discount factor (gamma)

MDP  
 $(\mathcal{S}, \mathcal{A}, \mathcal{R}, \mathbb{T}, \gamma)$

- Policy:
  - Mapping from states to actions (deterministic)
  - Distribution of actions given states (stochastic)

State	Action
A	→ 2
B	→ 1



- Markov Decision Processes (MDPs)
  - States, Actions, Reward dist., Transition dist., Discount factor (gamma)

MDP  
 $(\mathcal{S}, \mathcal{A}, \mathcal{R}, \mathbb{T}, \gamma)$

- Policy:
  - Mapping from states to actions (deterministic)
  - Distribution of actions given states (stochastic)

State	Action
A	→ 2
B	→ 1

- What is a good policy?
  - Maximize discounted sum of future rewards
    - Discount factor:  $\gamma$



1

Worth Now



$\gamma$

Worth Next Step



$\gamma^2$

Worth In Two Steps

- Formally, the **optimal policy** is defined as:

$$\underline{\pi^*} = \arg \max_{\pi} \mathbb{E} \left[ \sum_{t \geq 0} \gamma^t \underline{r_t} \mid \pi \right]$$

discounted sum of future rewards

- Formally, the **optimal policy** is defined as:

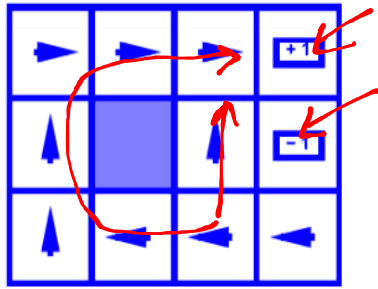
$$\pi^* = \arg \max_{\pi} \mathbb{E} \left[ \sum_{t \geq 0} \gamma^t r_t \mid \pi \right]$$

discounted sum of future rewards

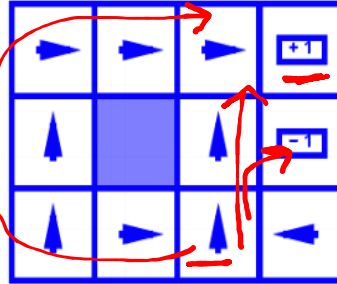
$$\underline{s_0 \sim p(s_0)}, \underline{a_t \sim \pi(\cdot | s_t)}, \underline{s_{t+1} \sim p(\cdot | s_t, a_t)}$$

Expectation over initial state, actions from policy,  
next states from transition distribution

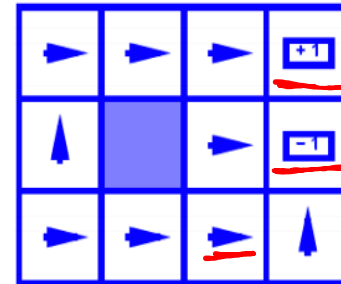
- Some optimal policies for three different grid world MDPs (gamma=0.99)
  - Varying reward for non-absorbing states (states other than +1/-1)



$R(s) = -0.03$



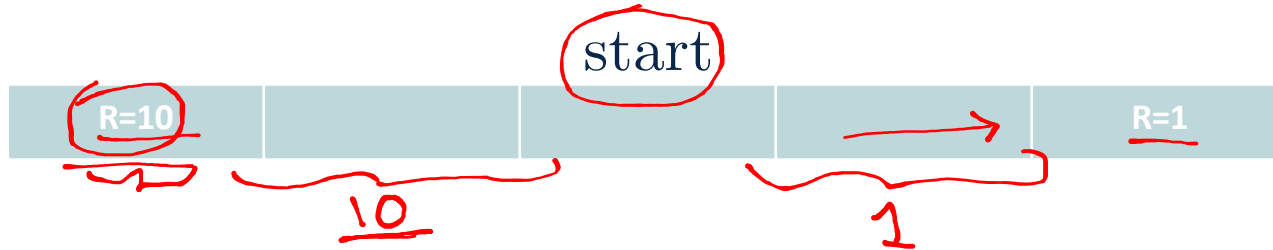
$R(s) = -0.4$



$R(s) = -2.0$

Image Credit: Byron Boots, CS 7641

- For example, with an MDP with 5 states as shown, starting at the middle cell:



- Actions: (Right, Left)
- Deterministic transitions
- What is the optimal policy for:

- $\gamma = 1$

- $\gamma = 0.1$  ← 10

Slides adapted from: Byron Boots, CS 7641



- A value function is a prediction of discounted sum of future reward
- State value function / **V**-function /  $V : \mathcal{S} \rightarrow \mathbb{R}$ 
  - How good is this state?
  - Am I likely to win/lose the game from this state?
- State-Action value function / **Q**-function /  $Q : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ 
  - How good is this state-action pair?
  - In this state, what is the impact of this action on my future?

- For a policy that produces a trajectory sample  $(s_0, a_0, s_1, a_1, s_2 \dots)$
- The **V-function** of the policy at state  $s$ , is the expected cumulative reward from state  $s$ :

$$V^{\pi}(s) = \mathbb{E} \left[ \sum_{t \geq 0} \gamma^t r_t \mid s_0 = s, \pi \right]$$

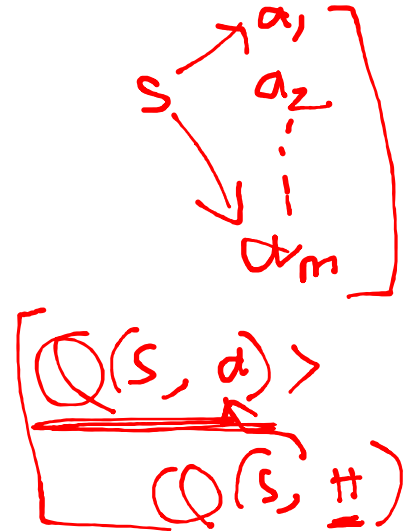
$n \times 1$

$$s_0 \sim p(s_0), a_t \sim \pi(\cdot | s_t), s_{t+1} \sim p(\cdot | s_t, a_t)$$

- For a policy that produces a trajectory sample  $(s_0, a_0, s_1, a_1, s_2 \dots)$
- The **Q-function** of the policy at state  $\mathbf{s}$  and action  $\mathbf{a}$ , is the expected cumulative reward upon taking action  $\mathbf{a}$  in state  $\mathbf{s}$  (and following policy thereafter):

$$Q^\pi(\mathbf{s}, \mathbf{a}) = \mathbb{E} \left[ \sum_{t \geq 0} \gamma^t r_t \mid \underline{s_0 = s, a_0 = a, \pi} \right]$$

$$\mathbf{s}_0 \sim p(\mathbf{s}_0), a_t \sim \pi(\cdot \mid \mathbf{s}_t), \mathbf{s}_{t+1} \sim p(\cdot \mid \mathbf{s}_t, a_t)$$



- The V and Q functions corresponding to the optimal policy  $\pi^*$

$$\underline{V^*(s)} = \mathbb{E} \left[ \sum_{t \geq 0} \gamma^t r_t \mid s_0 = s, \underline{\pi^*} \right]$$

$$\underline{Q^*(s, a)} = \mathbb{E} \left[ \sum_{t \geq 0} \gamma^t r_t \mid s_0 = s, a_0 = a, \underline{\pi^*} \right]$$

◆ Recursive Bellman expansion (from definition of Q)

$$Q^*(s, a) = \mathbb{E}_{\substack{a_t \sim \pi^*(\cdot | s_t) \\ s_{t+1} \sim p(\cdot | s_t, a_t)}} \left[ \sum_{t \geq 0} \gamma^t r(s_t, a_t) \mid s_0 = s, a_0 = a \right]$$

(Expected) return from t = 0

(Reward at t = 0) + gamma \* (Return from expected state at t=1)

$$\begin{aligned} &= \gamma^0 r(s, a) + \mathbb{E}_{s' \sim p(\cdot | s, a)} \left[ \gamma \mathbb{E}_{\substack{a_t \sim \pi^*(\cdot | s_t) \\ s_{t+1} \sim p(\cdot | s_t, a_t)}} \left[ \sum_{t \geq 1} \gamma^{t-1} r(s_t, a_t) \mid s_1 = s' \right] \right] \\ &= r(s, a) + \gamma \mathbb{E}_{s' \sim p(s' | s, a)} [V^*(s')] \\ &= \mathbb{E}_{s' \sim p(s' | s, a)} [r(s, a) + \gamma V^*(s')] \end{aligned}$$

- Equations relating optimal quantities

$$\rightarrow \underline{V^*(s)} = \max_a \underline{Q^*(s, a)}$$

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

- Recursive Bellman optimality equation

$$\begin{aligned} \underline{Q^*(s, a)} &= \mathbb{E}_{s' \sim p(s'|s, a)} [r(s, a) + \gamma \underline{V^*(s')}] \\ &= \sum_{s'} p(s' | s, a) [r(s, a) + \gamma \underline{V^*(s')}] \\ &= \sum_{s'} p(s' | s, a) [r(s, a) + \gamma \max_a \underline{Q^*(s', a)}] \end{aligned}$$

Diagram illustrating the Bellman optimality equation with handwritten annotations. A tree structure shows a state  $s$  at the root, branching into actions  $a^*$ . The tree is labeled with  $t=1$  and  $t=2$ , indicating time steps. A red arrow points from the  $V^*(s')$  term in the equation to the state  $s'$  in the tree. Another red arrow points from the  $\max_a Q^*(s', a)$  term to the action  $a^*$  in the tree. A red asterisk is placed next to the action  $a^*$  in the tree.

$$\underline{V^*(s)} = \max_a \sum_{s'} p(s' | s, a) [r(s, a) + \gamma \underline{V^*(s')}]$$

Diagram illustrating the Bellman optimality equation with handwritten annotations. The equation is written with red circles around  $V^*(s)$ ,  $\max_a$ , and  $r(s, a)$ . A red arrow points from the  $V^*(s')$  term in the equation to the state  $s'$  in the tree. Another red arrow points from the  $\max_a Q^*(s', a)$  term to the action  $a^*$  in the tree. A red asterisk is placed next to the action  $a^*$  in the tree.

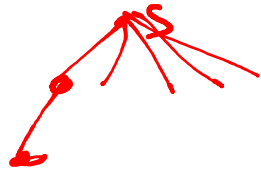
- Based on the bellman optimality equation

$$V^*(s) = \max_a \sum_{s'} p(s'|s, a) [r(s, a) + \gamma V^*(s')]$$

- Algorithm

- Initialize values of all states

$n \times 1$   $V(s)$   $V_0$



- While not converged:

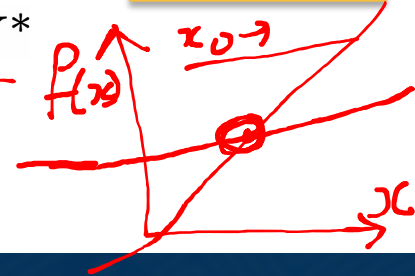
For each state:  $V^{i+1}(s) \leftarrow \max_a \sum_{s'} p(s'|s, a) [r(s, a) + \gamma V^i(s')]$

- Repeat until convergence (no change in values)

$$V^0 \rightarrow V^1 \rightarrow V^2 \rightarrow \dots \rightarrow V^i \rightarrow \dots \rightarrow V^*$$

Homework

Time complexity per iteration  $O(|S|^2 |A|)$



- Value Iteration Update:

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

$$\underline{V^{i+1}(s)} \leftarrow \max_a \sum_{s'} p(s'|s, a) [r(s, a) + \gamma \underline{V^i(s')}]$$

- Q-Iteration Update:

$$\underline{Q^{i+1}(s, a)} \leftarrow$$



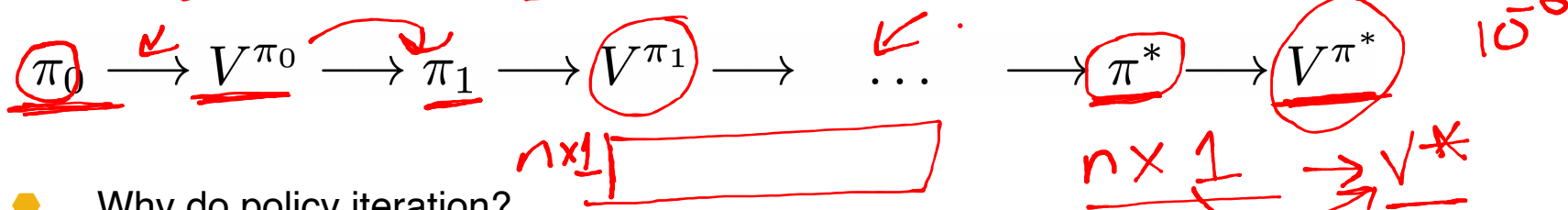
- Policy iteration: Start with arbitrary  $\pi_0$  and refine it.

$$\pi_0 \rightarrow \pi_1 \rightarrow \pi_2 \rightarrow \dots \rightarrow \pi^*$$

- Involves repeating two steps:

- Policy Evaluation: Compute  $V^\pi$  (similar to Value Iteration)

- Policy Refinement: Greedily change actions as per  $V^\pi(s)$

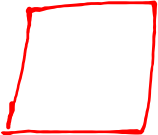


- Why do policy iteration?

- $\pi_i$  often converges to  $\pi^*$  much sooner than  $V^{\pi_i}$  to  $V^{\pi^*}$

For Value Iteration:

Time complexity per iteration  $O(|\mathcal{S}|^2|\mathcal{A}|)$

- 3x4 Grid world? 12 4
- Chess/Go? 10 120
- Atari Games with integer image pixel values [0, 255] of size 16x16 as state?  
256 16x16  


## ◆ Value Iteration

- ◆ Bellman update to state value estimates

## ◆ Q-Value Iteration

- ◆ Bellman update to (state, action) value estimates

## ◆ Policy Iteration

- ◆ Policy evaluation + refinement

# Reinforcement Learning, Deep RL

- Recall RL assumptions:
  - $\mathbb{T}(s, a, s')$  unknown, how actions affect the environment.
  - $\mathcal{R}(s, a, s')$  unknown, what/when are the good actions?
- But, we can learn by trial and error.
  - Gather experience (data) by performing actions.
  - Approximate unknown quantities from data.

# Reinforcement Learning

- Old Dynamic Programming Demo

- [https://cs.stanford.edu/people/karpathy/reinforcejs/gridworld\\_dp.html](https://cs.stanford.edu/people/karpathy/reinforcejs/gridworld_dp.html)

- RL Demo

- [https://cs.stanford.edu/people/karpathy/reinforcejs/gridworld\\_td.html](https://cs.stanford.edu/people/karpathy/reinforcejs/gridworld_td.html)

Slide credit: Dhruv Batra

- In addition to not knowing the environment, sometimes the state space is too large.
- Recall: Value iteration not scalable (chess, RGB images as state space, etc)
- Solution: Deep Learning! ... more precisely, function approximation.
  - Use deep neural networks to learn state representations
  - Useful for continuous action spaces as well

## Deep Reinforcement Learning

In today's class, we looked at

- ◆ **Dynamic Programming** for solving MDPs
  - ◆ Value, Q-Value Iteration
  - ◆ Policy Iteration
  
- ◆ **Reinforcement Learning (RL)**
  - ◆ The challenges of (deep) learning based methods

Next class:

- ◆ **Value-based RL algorithms**
  - ◆ Deep Q-Learning
- ◆ **Policy-based RL algorithms** (policy gradients)