

**CS394R**  
**Reinforcement Learning:**  
**Theory and Practice**

**Scott Niekum and Peter Stone**

Department of Computer Science  
The University of Texas at Austin

# Good Morning Colleagues

---

# Good Morning Colleagues

---

- Are there any questions?

# Logistics

---

- Registration

# Logistics

---

- Registration
- Reading responses on edX

# Logistics

---

- Registration
- Reading responses on edX
- Exercises - let us know if not worded clearly
  - Please bear with us

# Logistics

---

- Registration
- Reading responses on edX
- Exercises - let us know if not worded clearly
  - Please bear with us
  - Turn in when you want to

# Logistics

---

- Registration
- Reading responses on edX
- Exercises - let us know if not worded clearly
  - Please bear with us
  - Turn in when you want to
- Still working on programming submission site



# Logistics

---

- Registration
- Reading responses on edX
- Exercises - let us know if not worded clearly
  - Please bear with us
  - Turn in when you want to
- Still working on programming submission site
- Switched to Piazza

# Logistics

---

- Registration
- Reading responses on edX
- Exercises - let us know if not worded clearly
  - Please bear with us
  - Turn in when you want to
- Still working on programming submission site
- Switched to Piazza
- The math is important - use Piazza

# Logistics

---

- Registration
- Reading responses on edX
- Exercises - let us know if not worded clearly
  - Please bear with us
  - Turn in when you want to
- Still working on programming submission site
- Switched to Piazza
- The math is important - use Piazza
- Resources page - and Sutton materials

# Logistics

---

- Registration
- Reading responses on edX
- Exercises - let us know if not worded clearly
  - Please bear with us
  - Turn in when you want to
- Still working on programming submission site
- Switched to Piazza
- The math is important - use Piazza
- Resources page - and Sutton materials
- Lecture videos

# Logistics

---

- Registration
- Reading responses on edX
- Exercises - let us know if not worded clearly
  - Please bear with us
  - Turn in when you want to
- Still working on programming submission site
- Switched to Piazza
- The math is important - use Piazza
- Resources page - and Sutton materials
- Lecture videos
- Next week's readings

# Logistics

---

- Registration
- Reading responses on edX
- Exercises - let us know if not worded clearly
  - Please bear with us
  - Turn in when you want to
- Still working on programming submission site
- Switched to Piazza
- The math is important - use Piazza
- Resources page - and Sutton materials
- Lecture videos
- Next week's readings
- Final exam

# Chapter 3

---

- Defined the problem

# Chapter 3

---

- Defined the problem
- Introduced some important notation and concepts.



# Chapter 3

---

- Defined the problem
- Introduced some important notation and concepts.
  - Returns
  - Markov property
  - State/action value functions
  - Bellman equations

# Chapter 3

---

- Defined the problem
- Introduced some important notation and concepts.
  - Returns
  - Markov property
  - State/action value functions
  - Bellman equations
  - Get comfortable with them!

# Chapter 3

---

- Defined the problem
- Introduced some important notation and concepts.
  - Returns
  - Markov property
  - State/action value functions
  - Bellman equations
  - Get comfortable with them!
    - $q_{\pi}(s, a) =$

# Chapter 3

---

- Defined the problem
- Introduced some important notation and concepts.
  - Returns
  - Markov property
  - State/action value functions
  - Bellman equations
  - Get comfortable with them!
    - $q_{\pi}(s, a) =$
  - Backup diagrams

# Chapter 3

---

- Defined the problem
- Introduced some important notation and concepts.
  - Returns
  - Markov property
  - State/action value functions
  - Bellman equations
  - Get comfortable with them!
    - $q_{\pi}(s, a) =$
    - Backup diagrams
- Solution methods start in Chapter 4

# Chapter 3

---

- Defined the problem
- Introduced some important notation and concepts.
  - Returns
  - Markov property
  - State/action value functions
  - Bellman equations
  - Get comfortable with them!
    - $q_{\pi}(s, a) =$
  - Backup diagrams
- Solution methods start in Chapter 4
  - What does it mean to **solve** an RL problem?

# Formulating the RL problem

---

- Art more than science
- States, actions, rewards
- Rewards: no hints on **how** to solve the problem

# Formulating the RL problem

---

- Art more than science
- States, actions, rewards
- Rewards: no hints on **how** to solve the problem
- Discount factor part of the environment



# Value functions

---

- Consider the week 0 environment

# Value functions

---

- Consider the week 0 environment
- For some  $s$ , what is  $V(s)$ ?

# Value functions

---

- Consider the week 0 environment
- For some  $s$ , what is  $V(s)$ ?
- OK - consider the policy we ended with
- Now, for some  $s$ , what is  $V(s)$ ?

# Value functions

---

- Consider the week 0 environment
- For some  $s$ , what is  $V(s)$ ?
- OK - consider the policy we ended with
- Now, for some  $s$ , what is  $V(s)$ ?
- Construct  $V$  in undiscounted, episodic case

# Value functions

---

- Consider the week 0 environment
- For some  $s$ , what is  $V(s)$ ?
- OK - consider the policy we ended with
- Now, for some  $s$ , what is  $V(s)$ ?
- Construct  $V$  in undiscounted, episodic case
- Construct  $Q$  in undiscounted, episodic case

# Value functions

---

- Consider the week 0 environment
- For some  $s$ , what is  $V(s)$ ?
- OK - consider the policy we ended with
- Now, for some  $s$ , what is  $V(s)$ ?
- Construct  $V$  in undiscounted, episodic case
- Construct  $Q$  in undiscounted, episodic case
- What if it's discounted?

# Value functions

---

- Consider the week 0 environment
- For some  $s$ , what is  $V(s)$ ?
- OK - consider the policy we ended with
- Now, for some  $s$ , what is  $V(s)$ ?
- Construct  $V$  in undiscounted, episodic case
- Construct  $Q$  in undiscounted, episodic case
- What if it's discounted?
- What if it's continuing?

# Value functions

---

- Consider the week 0 environment
- For some  $s$ , what is  $V(s)$ ?
- OK - consider the policy we ended with
- Now, for some  $s$ , what is  $V(s)$ ?
- Construct  $V$  in undiscounted, episodic case
- Construct  $Q$  in undiscounted, episodic case
- What if it's discounted?
- What if it's continuing?
- Continuing tasks without discounting?



# Value functions

---

- Consider the week 0 environment
- For some  $s$ , what is  $V(s)$ ?
- OK - consider the policy we ended with
- Now, for some  $s$ , what is  $V(s)$ ?
- Construct  $V$  in undiscounted, episodic case
- Construct  $Q$  in undiscounted, episodic case
- What if it's discounted?
- What if it's continuing?
- Continuing tasks without discounting?

# Chapter 4

---

- Solution methods **given a model**

# Chapter 4

---

- Solution methods **given a model**
  - So no exploration vs. exploitation

# Chapter 4

---

- Solution methods **given a model**
  - So no exploration vs. exploitation

# Policy Evaluation

---

- $V^\pi$  exists and is unique if  $\gamma < 1$  or termination guaranteed for all states under policy  $\pi$ .

# Policy Evaluation

---

- $V^\pi$  exists and is unique if  $\gamma < 1$  or termination guaranteed for all states under policy  $\pi$ .
- Policy evaluation converges under the same conditions

# Policy Evaluation

---

- $V^\pi$  exists and is unique if  $\gamma < 1$  or termination guaranteed for all states under policy  $\pi$ .
- Policy evaluation converges under the same conditions
- Policy evaluation on the week 0 problem
  - undiscounted, episodic

# Policy Evaluation

---

- $V^\pi$  exists and is unique if  $\gamma < 1$  or termination guaranteed for all states under policy  $\pi$ .
- Policy evaluation converges under the same conditions
- Policy evaluation on the week 0 problem
  - undiscounted, episodic
  - Are the conditions met?



# Policy Evaluation

---

- $V^\pi$  exists and is unique if  $\gamma < 1$  or termination guaranteed for all states under policy  $\pi$ .
- Policy evaluation converges under the same conditions
- Policy evaluation on the week 0 problem
  - undiscounted, episodic
  - Are the conditions met?

# Policy Improvement

---

- Policy improvement theorem:

$$\forall s, q_{\pi}(s, \pi'(s)) \geq v_{\pi}(s) \Rightarrow \forall s, v_{\pi'}(s) \geq v_{\pi}(s)$$

# Policy Improvement

---

- Policy improvement theorem:

$$\forall s, q_{\pi}(s, \pi'(s)) \geq v_{\pi}(s) \Rightarrow \forall s, v_{\pi'}(s) \geq v_{\pi}(s)$$

- Polynomial time convergence (in number of states and actions) even though  $m^n$  policies.
  - Ignoring effect of  $\gamma$  and bits to represent rewards/transitions

# Value Iteration on Week 0 problem

---

- Show the new policy at each step
  - Not actually to compute policy

# Value Iteration on Week 0 problem

---

- Show the new policy at each step
  - Not actually to compute policy
  - Break policy ties with equiprobable actions

# Value Iteration on Week 0 problem

---

- Show the new policy at each step
  - Not actually to compute policy
  - Break policy ties with equiprobable actions
  - No stochastic transitions

# Value Iteration on Week 0 problem

---

- Show the new policy at each step
  - Not actually to compute policy
  - Break policy ties with equiprobable actions
  - No stochastic transitions
- How would policy iteration proceed in comparison?
  - More or fewer policy updates?

# Value Iteration on Week 0 problem

---

- Show the new policy at each step
  - Not actually to compute policy
  - Break policy ties with equiprobable actions
  - No stochastic transitions
- How would policy iteration proceed in comparison?
  - More or fewer policy updates?
  - True in general?



# Chapter 4 Summary

---

- Chapter 4 treats **bootstrapping** with a model

# Chapter 4 Summary

---

- Chapter 4 treats **bootstrapping** with a model
  - Next: no model and no bootstrapping

# Chapter 4 Summary

---

- Chapter 4 treats **bootstrapping** with a model
  - Next: no model and no bootstrapping
  - Then: no model, but bootstrapping