CS394R Reinforcement Learning: Theory and Practice

Peter Stone

Department of Computer Science The University of Texas at Austin

Good Morning Colleagues



• Are there any questions?



Peter Stone



• Reading responses





- Reading responses
- Next week's readings





- Reading responses
- Next week's readings
- The math is important





- Reading responses
- Next week's readings
- The math is important
- Use piazza



• Defines the problem



- Defines the problem
- Introduces some important notation and concepts.



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



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



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

$$-q_{\pi}(s,a) =$$



- Defines the problem
- Introduces some important notation and concepts.
 - Returns
 - Markov property
 - State/action value functions
 - Bellman equations
 - Get comfortable with them!
 - $-q_{\pi}(s,a) =$
 - Backup diagrams (p. 62)



- Defines the problem
- Introduces some important notation and concepts.
 - Returns
 - Markov property
 - State/action value functions
 - Bellman equations
 - Get comfortable with them!
 - $-q_{\pi}(s,a) =$
 - Backup diagrams (p. 62)
- Solution methods start in Chapter 4



- Defines the problem
- Introduces some important notation and concepts.
 - Returns
 - Markov property
 - State/action value functions
 - Bellman equations
 - Get comfortable with them!
 - $-q_{\pi}(s,a) =$
 - Backup diagrams (p. 62)
- 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
- Discounted vs. non-discounted



Formulating the RL problem

- Art more than science
- States, actions, rewards
- Rewards: no hints on **how** to solve the problem
- Discounted vs. non-discounted
- Episodic vs. continuing



• Consider the week 0 environment



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



- 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)?



- 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



- 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



- 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?



- 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?



- 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?



- 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?
- Exercises 3.9, 3.10, 3.16



• What is it?



- What is it?
- Does it hold in the real world?



- What is it?
- Does it hold in the real world?
 - Are any systems "fundamentally" non-Markovian?



- What is it?
- Does it hold in the real world?
 - Are any systems "fundamentally" non-Markovian?
 - What if there's a time horizon?



- What is it?
- Does it hold in the real world?
 - Are any systems "fundamentally" non-Markovian?
 - What if there's a time horizon?
- It's an ideal
 - Will allow us to prove properties of algorithms



- What is it?
- Does it hold in the real world?
 - Are any systems "fundamentally" non-Markovian?
 - What if there's a time horizon?
- It's an ideal
 - Will allow us to prove properties of algorithms
 - Algorithms may still work when not provably correct



- What is it?
- Does it hold in the real world?
 - Are any systems "fundamentally" non-Markovian?
 - What if there's a time horizon?
- It's an ideal
 - Will allow us to prove properties of algorithms
 - Algorithms may still work when not provably correct
 - Could you compensate? Do algorithms change?



- What is it?
- Does it hold in the real world?
 - Are any systems "fundamentally" non-Markovian?
 - What if there's a time horizon?
- It's an ideal
 - Will allow us to prove properties of algorithms
 - Algorithms may still work when not provably correct
 - Could you compensate? Do algorithms change?
 - If not, you may want different algorithms (Monte Carlo)



- What is it?
- Does it hold in the real world?
 - Are any systems "fundamentally" non-Markovian?
 - What if there's a time horizon?
- It's an ideal
 - Will allow us to prove properties of algorithms
 - Algorithms may still work when not provably correct
 - Could you compensate? Do algorithms change?
 - If not, you may want different algorithms (Monte Carlo)
- Exercise 3.6 (broken vision system)



• Solution methods given a model



• Solution methods given a model

- So no exploration vs. exploitation



• Solution methods given a model

- So no exploration vs. exploitation



• V^{π} exists and is unique if $\gamma < 1$ or termination guaranteed for all states under policy π .



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



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



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



- V^{π} exists and is unique if $\gamma < 1$ or termination guaranteed for all states under policy π .
- Policy evaluation converges under the same conditions
- Policy evaluation on the week 0 problem
 - undiscounted, episodic
 - Are the conditions met?
- Exercises 4.1, 4.2



• Policy improvement theorem: $\forall s, q_{\pi}(s, \pi'(s)) \ge v_{\pi}(s) \Rightarrow \forall s, v_{\pi'}(s) \ge v_{\pi}(s)$



- Policy improvement theorem: $\forall s, q_{\pi}(s, \pi'(s)) \ge v_{\pi}(s) \Rightarrow \forall s, v_{\pi'}(s) \ge v_{\pi}(s)$
- Polynomial time convergence (in number of states and actions) even though m^n policies.
 - Ignoring effect of γ and bits to represent rewards/transitions



- Policy improvement theorem: $\forall s, q_{\pi}(s, \pi'(s)) \ge v_{\pi}(s) \Rightarrow \forall s, v_{\pi'}(s) \ge v_{\pi}(s)$
- Polynomial time convergence (in number of states and actions) even though m^n policies.
 - Ignoring effect of γ and bits to represent rewards/transitions



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



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



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



- 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?



- 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 treats bootstrapping with a model



Peter Stone

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



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

