# CS395T Reinforcement Learning: Theory and Practice Fall 2004

#### **Peter Stone**

Department or Computer Sciences The University of Texas at Austin

Week4a: Tuesday, September 21st

#### **Good Afternoon Colleagues**

• Are there any questions?



- 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?



 $\bullet$  Relationship between Q and V



- $\bullet$  Relationship between Q and V
- Bellman equations



- $\bullet$  Relationship between Q and V
- Bellman equations unique solution



- $\bullet$  Relationship between Q and V
- Bellman equations unique solution
- Backup diagrams (p. 70, 74, 77)



- $\bullet$  Relationship between Q and V
- Bellman equations unique solution
- Backup diagrams (p. 70, 74, 77)
- Exercise 3.17



#### • Solution methods given a model



- Solution methods **given a model**
- Why is it called dynamic programming?



• Susan on the Gambler's Problem (p. 101)



- Susan on the Gambler's Problem (p. 101)
- Email discussion linked to the book web page



•  $V^{\pi}$  exists and is unique if  $\gamma < 1$  or termination guaranteed for all states under policy  $\pi$ . (p. 90)



- $V^{\pi}$  exists and is unique if  $\gamma < 1$  or termination guaranteed for all states under policy  $\pi$ . (p. 90)
- Policy evaluation converges under the same conditions (p. 91)



- $V^{\pi}$  exists and is unique if  $\gamma < 1$  or termination guaranteed for all states under policy  $\pi$ . (p. 90)
- Policy evaluation converges under the same conditions (p. 91)
- Exercises 4.1, 4.2



- $V^{\pi}$  exists and is unique if  $\gamma < 1$  or termination guaranteed for all states under policy  $\pi$ . (p. 90)
- Policy evaluation converges under the same conditions (p. 91)
- Exercises 4.1, 4.2
- Policy improvement theorem:  $\forall s, Q^{\pi}(s, \pi'(s)) \ge V^{\pi}(s) \Rightarrow \forall s, V^{\pi'}(s) \ge V^{\pi}(s)$



• p. 107: Is LP still inefficient?



- p. 107: Is LP still inefficient?
- p. 109: This chapter treats **bootstrapping** with a model



- p. 107: Is LP still inefficient?
- p. 109: This chapter treats **bootstrapping** with a model
  Next: no model and no bootstrapping



- p. 107: Is LP still inefficient?
- p. 109: This chapter treats **bootstrapping** with a model
  - Next: no model and no bootstrapping
  - Then: no model, but bootstrapping

