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?





• Do programming assignments!





- Do programming assignments!
- Next week's readings





- Do programming assignments!
- Next week's readings
 - Approximate on-policy prediction and control



- Week 0 example
 - (Remember no access to real model)
 - $\ \alpha = .1, \epsilon\text{-greedy} \ \epsilon = .75,$ break ties in favor of \rightarrow



- Week 0 example
 - (Remember no access to real model)
 - $-\alpha = .1, \epsilon$ -greedy $\epsilon = .75$, break ties in favor of \rightarrow
 - Where did policy change?



- Week 0 example
 - (Remember no access to real model)
 - $\ \alpha = .1, \epsilon\text{-greedy} \ \epsilon = .75,$ break ties in favor of \rightarrow
 - Where did policy change?
- How do their convergence guarantees differ?



- Week 0 example
 - (Remember no access to real model)
 - $\alpha = .1$, ϵ -greedy $\epsilon = .75$, break ties in favor of \rightarrow
 - Where did policy change?
- How do their convergence guarantees differ?
 - Sarsa depends on policy's dependence on Q:
 - Policy must converge to greedy



- Week 0 example
 - (Remember no access to real model)
 - $\ \alpha = .1, \epsilon\text{-greedy} \ \epsilon = .75,$ break ties in favor of \rightarrow
 - Where did policy change?
- How do their convergence guarantees differ?
 - Sarsa depends on policy's dependence on Q:
 - Policy must converge to greedy
 - Q-learning value function converges to Q^*
 - As long as all state-action pairs visited infinitely
 - And step-size satisfies stochastic convergence equations



• Why does Q-learning learn to hug the cliff? (p. 139)



- Why does Q-learning learn to hug the cliff? (p. 139)
- Why can expected SARSA use high α , but SARSA can't? (p. 140)

