# CS394R Reinforcement Learning: Theory and Practice

Peter Stone

Department of Computer Science The University of Texas at Austin

#### **Good Afternoon Colleagues**

• Are there any questions?



- State: position, orientation, velocity, angular vels
- Actions: Settings of the 4 or 5 controls
- Goal: Hover



- State: position, orientation, velocity, angular vels
- Actions: Settings of the 4 or 5 controls
- Goal: Hover
- How would you formulate the problem "by the book"?



- State: position, orientation, velocity, angular vels
- Actions: Settings of the 4 or 5 controls
- Goal: Hover
- How would you formulate the problem "by the book"?
- Could you implement that? Why or why not?



- State: position, orientation, velocity, angular vels
- Actions: Settings of the 4 or 5 controls
- Goal: Hover
- How would you formulate the problem "by the book"?
- Could you implement that? Why or why not?
- At a high level, what do they do instead?



- State: position, orientation, velocity, angular vels
- Actions: Settings of the 4 or 5 controls
- Goal: Hover
- How would you formulate the problem "by the book"?
- Could you implement that? Why or why not?
- At a high level, what do they do instead?
  - Collect a small amount of human expert data
  - Use that to train a **1-step** model (simulator)
  - Determine the optimal policy in the simulator
  - Fly it!



- State: position, orientation, velocity, angular vels
- Actions: Settings of the 4 or 5 controls
- Goal: Hover
- How would you formulate the problem "by the book"?
- Could you implement that? Why or why not?
- At a high level, what do they do instead?
  - Collect a small amount of human expert data
  - Use that to train a **1-step** model (simulator)
  - Determine the optimal policy in the simulator
  - Fly it!
- Would this approach work on the Aibo walking task?



• Why quadratic reward (p. 5)?



- Why quadratic reward (p. 5)?
- PEGASUS how does it help policy evaluation?



- Why quadratic reward (p. 5)?
- PEGASUS how does it help policy evaluation?
  - General question: is policy good or lucky?



- Why quadratic reward (p. 5)?
- PEGASUS how does it help policy evaluation?
  - General question: is policy good or lucky?
  - Use same random samples to evaluate each policy



- Why quadratic reward (p. 5)?
- PEGASUS how does it help policy evaluation?
  - General question: is policy good or lucky?
  - Use same random samples to evaluate each policy
- How does he do policy optimization?



- Why quadratic reward (p. 5)?
- PEGASUS how does it help policy evaluation?
  - General question: is policy good or lucky?
  - Use same random samples to evaluate each policy
- How does he do policy optimization?
  - Represent policy as a Neural net (Fig 2c)



- Why quadratic reward (p. 5)?
- PEGASUS how does it help policy evaluation?
  - General question: is policy good or lucky?
  - Use same random samples to evaluate each policy
- How does he do policy optimization?
  - Represent policy as a Neural net (Fig 2c)
  - greedy hillclimbing over few parameters (the weights)!



- Why quadratic reward (p. 5)?
- PEGASUS how does it help policy evaluation?
  - General question: is policy good or lucky?
  - Use same random samples to evaluate each policy
- How does he do policy optimization?
  - Represent policy as a Neural net (Fig 2c)
  - greedy hillclimbing over few parameters (the weights)!
- Shaping rewards



- Why quadratic reward (p. 5)?
- PEGASUS how does it help policy evaluation?
  - General question: is policy good or lucky?
  - Use same random samples to evaluate each policy
- How does he do policy optimization?
  - Represent policy as a Neural net (Fig 2c)
  - greedy hillclimbing over few parameters (the weights)!
- Shaping rewards
- Topics for other courses:
  - Kalman filter (robotics)
  - Cross-validation/hold-out testing (supervised learning)



- Why quadratic reward (p. 5)?
- PEGASUS how does it help policy evaluation?
  - General question: is policy good or lucky?
  - Use same random samples to evaluate each policy
- How does he do policy optimization?
  - Represent policy as a Neural net (Fig 2c)
  - greedy hillclimbing over few parameters (the weights)!
- Shaping rewards
- Topics for other courses:
  - Kalman filter (robotics)
  - Cross-validation/hold-out testing (supervised learning)
- Can it generalize to adverse conditions?



- Why quadratic reward (p. 5)?
- PEGASUS how does it help policy evaluation?
  - General question: is policy good or lucky?
  - Use same random samples to evaluate each policy
- How does he do policy optimization?
  - Represent policy as a Neural net (Fig 2c)
  - greedy hillclimbing over few parameters (the weights)!
- Shaping rewards
- Topics for other courses:
  - Kalman filter (robotics)
  - Cross-validation/hold-out testing (supervised learning)
- Can it generalize to adverse conditions?
- Easy problem or a powerful approach?