

Autonomous Robotics

Reinforcement Learning I

Josiah Hanna

University of Wisconsin — Madison

Announcements

Hope you had a good spring break.

Thanks to those who completed the midterm evaluation.

Final project preview

Midterm grades.

Looking ahead:

- We have finished the classical robotics part of the course.
- Final four weeks will cover advanced topics: RL, HRI, robot learning, applications.

Learning Outcomes

After today's lecture, you will:

- Understand the motivation for reinforcement learning (RL) in robotics.
- Understand when RL is (and is not) a good tool in robotics.
- Understand how to define an RL problem.
- Be able to identify key classes of RL methods for robot control problems.

What is Reinforcement Learning?

- Type of machine learning that focuses on learning from rewards and trial and error interaction.
- The learning agent takes actions, receives rewards, and over time learns to take actions that lead to the most reward.
- Think: training a dog to do tricks.



Why (and why not) RL in robotics?

- Opportunities:
 - Well-suited for tasks where success can be defined but the correct actions to achieve success are unclear.
 - Well-suited for addressing unknown environments.
 - Well-suited for changing environments.
- Challenges:
 - May require long training times.
 - In some cases, we already have good existing controllers (e.g., model-predictive control)
 - Success may be difficult to define.

Be an RL Agent*

- You (as a class) are the learning agent.
- Three actions: stand, clap, or wave
- Observations: colors $\in \{\text{red, blue, orange, pink}\}$
- Rewards: depends on color you see and action you take.
- Goal: find the optimal policy.
 - Policy: mapping from colors to actions.
 - Optimal policy: policy that gives you the most reward.

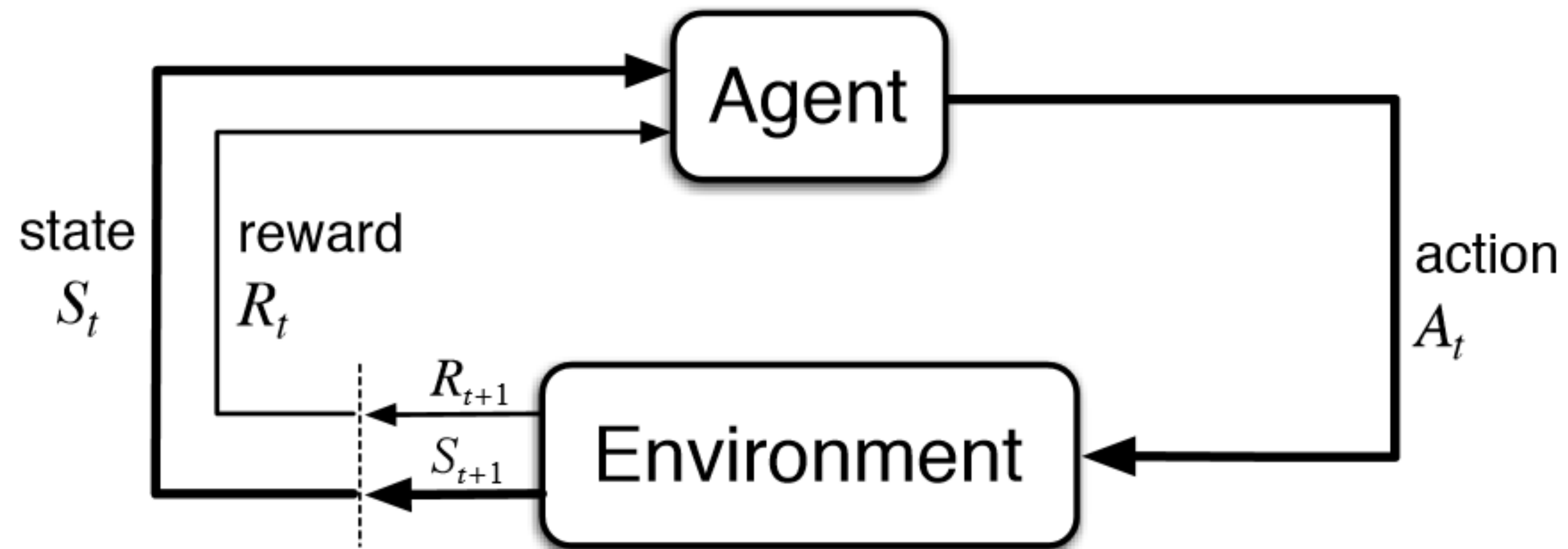
General Reinforcement Learning

- Typical formalism: Markov decision process
- States: $s \in \mathcal{S}$
- Actions: $a \in \mathcal{A}$
- Rewards: $R \sim r(s, a)$
- State transitions: $S \sim p(\cdot | s, a)$ **Markov!**
- Goal: Find a policy, $\pi : \mathcal{S} \rightarrow \mathcal{A}$, that maximizes cumulative reward.

Practice

- Define a Markov decision process (MDP) for:
 - A quadruped robot learning to walk at a desired velocity (forward, lateral, and rotational).
 - A humanoid robot learning to walk forward.
 - A robot arm picking up objects out of a cluttered bin.

General Reinforcement Learning



$\dots S_t, A_t, R_{t+1}, S_{t+1}, A_{t+1}, \dots$

$$S_{t+1}, R_{t+1} \sim p(\cdot | S_t, A_t)$$

$$A_{t+1} \leftarrow \pi(S_{t+1})$$

Returns and Episodes

- The **return** is the discounted sum of future rewards:

$$G_t := R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots$$

- Recursive definition: $G_t = R_{t+1} + \gamma G_{t+1}$.
- **Episodes** are subsequences of interaction that begin in some initial state and end in a special terminal state.
- The initial state of one episode is independent of interaction in the preceding episode.

Value functions

- State transitions and rewards are stochastic so we define the utility of states and actions in terms of **expected return**.
- The expected return from a state, $v_\pi(s)$, is only well-defined with respect to a particular policy, π . (Why?)
- State-value and action-value functions:

$$v_\pi(s) = \mathbb{E}_\pi[G_t | S_t = s] = \mathbb{E}_\pi\left[\sum_{k=0} \gamma^k R_{t+k+1} | S_t = s\right]$$

$$q_\pi(s, a) = \mathbb{E}_\pi[G_t | S_t = s, A_t = a] = \mathbb{E}_\pi\left[\sum_{k=0} \gamma^k R_{t+k+1} | S_t = s, A_t = a\right]$$

Optimality

- Agent's objective: find the policy that maximizes $v_{\pi}(s)$ for all s .
- The optimal policy is the policy that has maximal value in all states. $\pi^{\star} \geq \pi$ if $v_{\pi^{\star}}(s) \geq v_{\pi}(s)$ for all states and possible policies.
- $\pi^{\star}(s) = \arg \max_a q_{\pi^{\star}}(s, a)$
- $v_{\pi^{\star}}(s) = \max_a \mathbf{E}[R_{t+1} + \gamma v_{\pi^{\star}}(S_{t+1}) | S_t = s, A_t = a]$

Value Iteration in RL

- Goal: compute value functions and then use them to compute optimal policies.
- Turn optimality equation into value function updates.
 - $\forall s$, initialize $v_0(s)$, e.g., $v_0(s) \leftarrow 0$.
 - Loop over states and make the update:

$$v_{k+1}(s) \leftarrow \max_a \sum_{s'} \sum_r p(s', r | s, a) [r + \gamma v_k(s')]$$

$$= \max_a \mathbb{E}[R_{t+1} + \gamma v_k(S_{t+1}) | S_t = s, A_t = a]$$

Q-Learning

- Q-learning: value iteration without knowing the transition function.
- At each time-step, take an action and receive data, $(S_t, A_t, R_{t+1}, S_{t+1})$

$$q_{k+1}(S_t, A_t) \leftarrow q_k(S_t, A_t) + \alpha [R_{t+1} + \gamma \max_{a'} q_k(S_{t+1}, a') - q_k(S_t, A_t)]$$

- Converges to q^* for any sufficiently exploratory exploration policy.
- The underlying algorithm for Deep Q-networks, which was a landmark result in the history of RL.



Q-Learning for Continuous Actions

$$q_{k+1}(S_t, A_t) \leftarrow q_k(S_t, A_t) + \alpha [R_{t+1} + \gamma \max_{a'} q_k(S_{t+1}, a') - q_k(S_t, A_t)]$$

- The max over the action space is difficult to compute when actions take on real-values. Why?
- One idea: use optimization to find the best action.
 - Examples: cross-entropy method, gradient ascent
 - But can be slow.
- Another idea: discretize the action space.
 - Sometimes works but loses precision in control.



QT-Opt: Scalable Deep Reinforcement Learning for Vision-Based Robotic Manipulation

Deterministic Actor-Critic

$$q_{k+1}(S_t, A_t) \leftarrow q_k(S_t, A_t) + \alpha [R_{t+1} + \gamma \max_{a'} q_k(S_{t+1}, a') - q_k(S_t, A_t)]$$

- Final idea: learn a policy (often a neural network) that outputs the maximizing action.
- $\mu_\theta(s) \rightarrow a$.
- Learn θ such that $q_k(s, \mu_\theta(s)) \approx \max_a q_k(s, a)$.
- $\theta_{t+1} \leftarrow \theta_t + \alpha \nabla_\theta q_k(s, \mu_\theta(s))$
- Actor: the policy μ_θ .
- Critic: the action-value function, trained to estimate q_{μ_θ} .

Stochastic Policy Gradient RL

- Policy gradient methods use a differentiable and stochastic policy (e.g., a neural network) and learn policy parameters with gradient ascent.
- $\pi_{\theta}(a | s) = \Pr(A_t = a | S_t = s; \theta)$
- $J(\theta) = v_{\pi_{\theta}}(s_0)$ for some special start state s_0 .
- $\theta_{k+1} \leftarrow \theta_k + \alpha \nabla_{\theta} J(\theta_k)$

REINFORCE

- $\nabla_{\theta} J(\theta)$ can only be estimated.

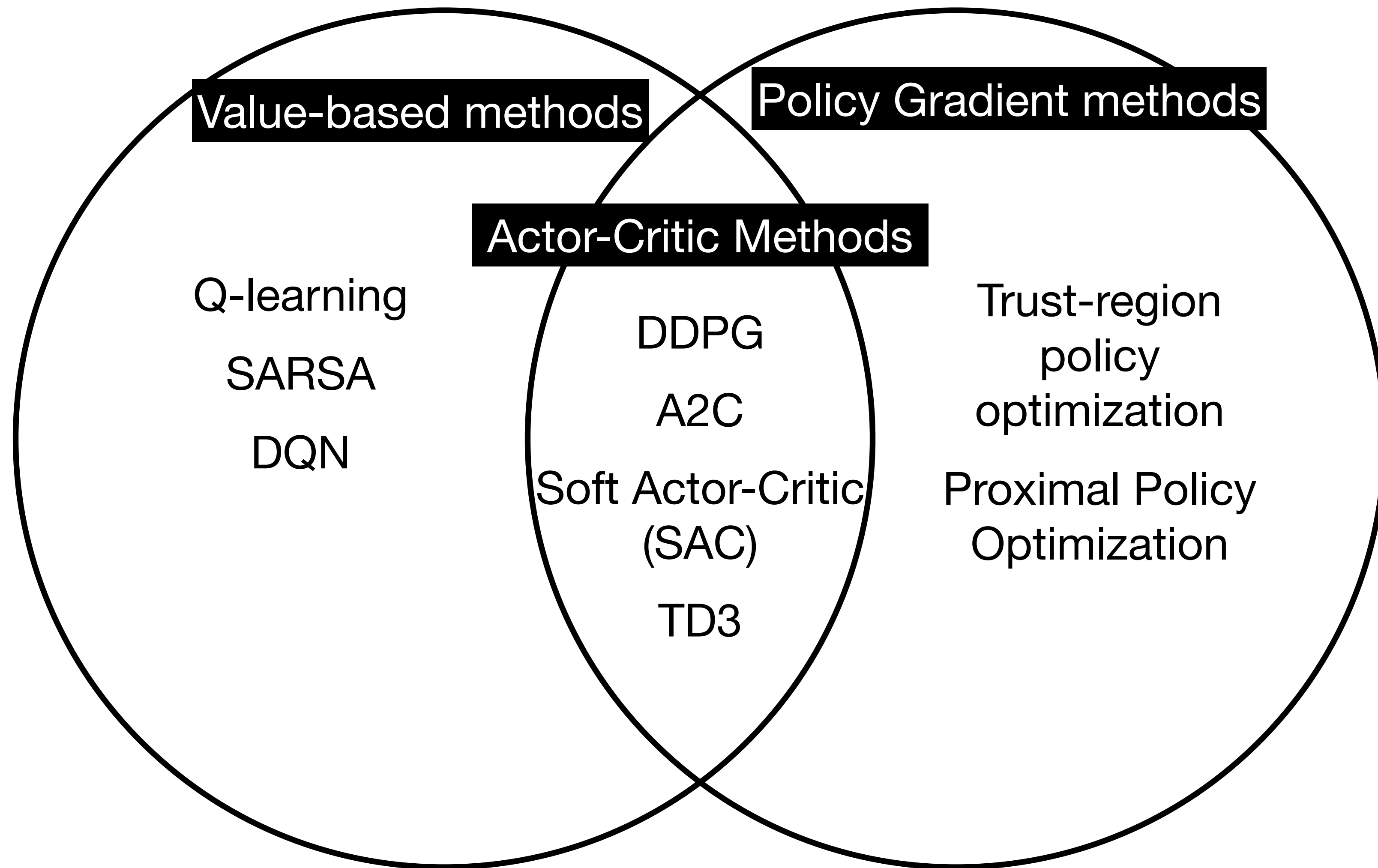
- $$\nabla_{\theta} J(\theta) \propto \mathbf{E} \left[\sum_a \nabla_{\theta} \pi(a | S_t) q_{\pi}(S_t, a) \right] = \mathbf{E} \left[\sum_a \pi_{\theta}(a | S_t) \frac{\nabla_{\theta} \pi_{\theta}(a | S_t)}{\pi_{\theta}(a | S_t)} q_{\pi}(S_t, a) \right]$$

- Finally, replace $q_{\pi}(s, a)$ with G_t .
- $\theta_{k+1} \leftarrow \theta_k + \alpha G_t \nabla_{\theta} \ln \pi_{\theta}(A_t | S_t)$

Actor-Critic Methods

- REINFORCE uses the return following an action to determine which actions are reinforced.
- Actor-critic methods use learned value functions to drive policy changes.
 - Actor: the policy, π .
 - Critic: value function, \hat{v} .
 - $\theta_{t+1} \leftarrow \theta_t + \alpha \delta_t \nabla_{\theta} \ln \pi(A_t | S_t)$ $\delta_t \leftarrow R_{t+1} + \gamma \hat{v}(S_{t+1}) - \hat{v}(S_t)$

Model-Free RL



Model-based RL

- Model-free RL learns value function or policies directly from experience.
- Model-based RL first learns environment transition function and reward functions; then learns/computes policies using estimated transitions.
 - Common approach: simulate experience and apply model-free RL to simulated experience.
- Advantages: often very sample efficient; potentially more generalizable.
- Disadvantages: sometimes dynamics are difficult to learn; in practice, model may generalize poor out-of-distribution of training data.

Summary

Today we covered:

1. The reinforcement learning problem definition.
2. Basic algorithms for RL.
3. Discussion of advantages and challenges with using RL.

Action Items

Complete homework 4

Begin HRI Reading

Start final project!