CSE 573: Artificial Intelligence

Hanna Hajishirzi Reinforcement Learning

slides adapted from Dan Klein, Pieter Abbeel ai.berkeley.edu And Dan Weld, Luke Zettelmoyer



Reinforcement Learning



Double Bandits



Double-Bandit MDP



Offline Planning

• Solving MDPs is offline planning

- o You determine all quantities through computation
- o You need to know the details of the MDP

• You do not actually play the game!





Let's Play!





\$2\$2\$0\$2\$2\$0\$0\$0

Online Planning

• Rules changed! Red's win chance is different.



Let's Play!





\$0
\$0
\$2
\$0
\$2
\$0
\$0

What Just Happened?

• That wasn't planning, it was learning!

- o Specifically, reinforcement learning
- There was an MDP, but you couldn't solve it with just computation
- o You needed to actually act to figure it out

• Important ideas in reinforcement learning that came up

- o Exploration: you have to try unknown actions to get information
- o Exploitation: eventually, you have to use what you know
- o Regret: even if you learn intelligently, you make mistakes
- o Sampling: because of chance, you have to try things repeatedly
- o Difficulty: learning can be much harder than solving a known MDP



Reinforcement Learning

• Still assume a Markov decision process (MDP):

- \circ A set of states $s \in S$
- A set of actions (per state) A
- o A model T(s,a,s')
- A reward function R(s,a,s')
- Still looking for a policy $\pi(s)$





• New twist: don't know T or R

- o I.e. we don't know which states are good or what the actions do
- o Must actually try actions and states out to learn

Reinforcement Learning

• Basic idea:

- Receive feedback in the form of **rewards**
- o Agent's utility is defined by the reward function
- Must (learn to) act so as to maximize expected rewards
- All learning is based on observed samples of outcomes!

Example: Learning to Walk

Initial

A Learning Trial

After Learning [1K Trials]

[Kohl and Stone, ICRA 2004]

Example: Toddler Robot

[Tedrake, Zhang and Seung, 2005]

[Video: TODDLER – 40s]

Robotics Rubik Cub

o <u>https://www.youtube.com/watch?v=x4O8pojMF0w</u>

The Crawler!

[Demo: Crawler Bot (L10D1)] [You, in Project 3]

Video of Demo Crawler Bot

Reinforcement Learning

• Still assume a Markov decision process (MDP):

- \circ A set of states $s \in S$
- A set of actions (per state) A
- o A model T(s,a,s')
- A reward function R(s,a,s')
- Still looking for a policy $\pi(s)$

• New twist: don't know T or R

- o I.e. we don't know which states are good or what the actions do
- o Must actually try actions and states out to learn

Offline (MDPs) vs. Online (RL)

Offline Solution

Online Learning

Model-Based Learning

Model-Based Learning

• Model-Based Idea:

Learn an approximate model based on experiencesSolve for values as if the learned model were correct

• Step 1: Learn empirical MDP model

- Count outcomes s' for each s, a
- Normalize to give an estimate $\hat{T}(s, a, s')$
- Discover each $\hat{R}(s, a, s')$ when we experience (s, a, s')
- Step 2: Solve the learned MDP

• For example, use value iteration, as before

Example: Model-Based Learning

Analogy: Expected Age

Goal: Compute expected age of cse573 students

Without P(A), instead collect samples $[a_1, a_2, ..., a_N]$

Announcements

- o HW1 is due: Feb. 12
- Project proposal is due: Feb 19th
- PS3 is released: Feb 21st
- Wait for paper reports.
- o Mid-quarter review: Feb 12

The Story So Far: MDPs and RL

Unknown MDP: Model-Based

Technique

Goal

Compute V*, Q*, π^* VI/PI on approx. MDP

Evaluate a fixed policy π PE on approx. MDP

Unknown MDP: Model-Free

Analogy: Expected Age

Goal: Compute expected age of cse573 students

Without P(A), instead collect samples $[a_1, a_2, ..., a_N]$

Model-Free Learning

Passive Reinforcement Learning

Passive Reinforcement Learning

• Simplified task: policy evaluation

- ο Input: a fixed policy $\pi(s)$
- You don't know the transitions T(s,a,s')
- o You don't know the rewards R(s,a,s')

o Goal: learn the state values

• In this case:

- o Learner is "along for the ride"
- o No choice about what actions to take
- o Just execute the policy and learn from experience
- o This is NOT offline planning! You actually take actions in the world.

Direct Evaluation

- Goal: Compute values for each state under π
- Idea: Average together observed sample values
 - \circ Act according to π
 - Every time you visit a state, write down what the sum of discounted rewards turned out to be
 - o Average those samples

• This is called direct evaluation

Example: Direct Evaluation

Problems with Direct Evaluation

• What's good about direct evaluation?

- o It's easy to understand
- o It doesn't require any knowledge of T, R
- It eventually computes the correct average values, using just sample transitions

• What bad about it?

- o It wastes information about state connections
- o Each state must be learned separately
- o So, it takes a long time to learn

Output Values

If B and E both go to C under this policy, how can their values be different?

Why Not Use Policy Evaluation?

Key question: how can we do this update to V without knowing T and R?
 In other words, how to we take a weighted average without knowing the weights?

Sample-Based Policy Evaluation?

• We want to improve our estimate of V by computing these averages:

$$V_{k+1}^{\pi}(s) \leftarrow \sum T(s, \pi(s), s') [R(s, \pi(s), s') + \gamma V_k^{\pi}(s')]$$

• Idea: Take samples of outcomes s' (by doing the action!) and average

$$sample_{1} \neq R(s, \pi(s), s_{1}') + \gamma V_{k}^{\pi}(s_{1}')$$

$$sample_{2} = R(s, \pi(s), s_{2}') + \gamma V_{k}^{\pi}(s_{2}')$$

$$\dots$$

$$sample_{n} = R(s, \pi(s), s_{n}') + \gamma V_{k}^{\pi}(s_{n}')$$

$$V_{k+1}^{\pi}(s) \leftarrow \frac{1}{n} \sum_{i} sample_{i}$$

Temporal Difference Learning

 $\pi(s)$

s, $\pi(s)$

 \mathbf{s}'

• Big idea: learn from every experience!

- Update V(s) each time we experience a transition (s, a, s', r)
- Likely outcomes s' will contribute updates more often

Temporal difference learning of values Ο

- Policy still fixed, still doing evaluation!
- Move values toward value of whatever successor occurs: running S, C1, 5', Y average

Sample of V(s): sample = $R(s, \pi(s), s') + \gamma V^{\pi}(s')$

Update to V(s):
$$V^{\pi}(s) \leftarrow (1 - \alpha)V^{\pi}(s) + (\alpha)sample$$

Same update: $V^{\pi}(s) \leftarrow V^{\pi}(s) + \alpha(sample - V^{\pi}(s))$

 $\underbrace{V^{\pi}(s)}_{} \leftarrow \underbrace{V^{\pi}(s)}_{} + \alpha(sample - V^{\pi}(s))$

Exponential Moving Average

• Exponential moving average • The running interpolation update: $\bar{x}_n \neq (1 - \alpha) \cdot \bar{x}_{n-1} + \alpha \cdot x_n$

• Makes recent samples more important:

$$\bar{x}_{n} = \underbrace{x_{n} + (1 - \alpha) \cdot x_{n-1} + (1 - \alpha)^{2} \cdot x_{n-2} + \dots}_{1 + (1 - \alpha) + (1 - \alpha)^{2} + \dots}$$

 $V_{IC}(S')$

o Forgets about the past (distant past values were wrong anyway)

• Decreasing learning rate (alpha) can give converging averages

Example: Temporal Difference Learning

Problems with TD Value Learning

TD value leaning is a model-free way to do policy evaluation, mimicking Bellman updates with running sample averages
However, if we want to turn values into a (new) policy, we're sunk:

$$\pi(s) = \arg\max_{a} Q(s, a)$$

$$Q(s,a) = \sum_{s'} T(s,a,s') \left[R(s,a,s') + \gamma V(s') \right]$$

Idea: learn Q-values, not values
Makes action selection model-free too!

Active Reinforcement Learning

Active Reinforcement Learning

- Full reinforcement learning: optimal policies (like value iteration)
 - o You don't know the transitions T(s,a,s')
 - o You don't know the rewards R(s,a,s')
 - You choose the actions now
 - o Goal: learn the optimal policy / values

• In this case:

- o Learner makes choices!
- o Fundamental tradeoff: exploration vs. exploitation
- This is NOT offline planning! You actually take actions in the world and find out what happens...

Detour: Q-Value Iteration

• Value iteration: find successive (depth-limited) values

- Start with $V_0(s) = 0$, which we know is right
- Given $V_{k'}$ calculate the depth k+1 values for all states:

$$V_{k+1}(s) \leftarrow \max_{a} \sum_{s'} T(s, a, s') \left[R(s, a, s') + \gamma V_k(s') \right]$$

• But Q-values are more useful, so compute them instead

- Start with $Q_0(s,a) = 0$, which we know is right
- o Given $Q_{k'}$ calculate the depth k+1 q-values for all q-states:

$$Q_{k+1}(s,a) \leftarrow \sum_{s'} T(s,a,s') \left[R(s,a,s') + \gamma \max_{a'} Q_k(s',a') \right]$$

Q-Learning

• Q-Learning: sample-based Q-value iteration

$$Q_{k+1}(s,a) \leftarrow \sum_{s'} T(s,a,s') \left[R(s,a,s') + \gamma \max_{a'} Q_k(s',a') \right]$$

- Learn Q(s,a) values as you go
 - o Receive a sample (s,a,s',r)
 - \circ Consider your old estimatQ(s, a)
 - o Consider your new sample estimate:

 $sample = R(s, a, s') + \gamma \max_{a'} Q(s', a')$ no longer policy evaluation!

o Incorporate the new estimate into a running average

$$Q(s,a) \leftarrow (1-\alpha)Q(s,a) + (\alpha) [sample]$$

[Demo: Q-learning – gridworld (L10D2)] [Demo: Q-learning – crawler (L10D3)]

Q-Learning Demo

Video of Demo Q-Learning -- Gridworld

Video of Demo Q-Learning -- Crawler

Q-Learning: act according to current optimal (and also explore...)

- Full reinforcement learning: optimal policies (like value iteration)
 - o You don't know the transitions T(s,a,s')
 - o You don't know the rewards R(s,a,s')
 - You choose the actions now
 - o Goal: learn the optimal policy / values

• In this case:

- o Learner makes choices!
- o Fundamental tradeoff: exploration vs. exploitation
- This is NOT offline planning! You actually take actions in the world and find out what happens...

Q-Learning Properties

- Amazing result: Q-learning converges to optimal policy -even if you're acting suboptimally!
- This is called off-policy learning
- Caveats:
 - You have to explore enough
 - You have to eventually make the learning rate small enough
 - o ... but not decrease it too quickly
 - o Basically, in the limit, it doesn't matter how you select actions

Exploration vs. Exploitation

How to Explore?

Several schemes for forcing exploration
 Simplest: random actions (ε-greedy)
 Every time step, flip a coin
 With (small) probability ε, act randomly
 With (large) probability 1-ε, act on current policy

O Problems with random actions?
O You do eventually explore the space, but keep thrashing around once learning is done
One solution: lower ε over time
O Another solution: exploration functions

[Demo: Q-learning – manual exploration – bridge grid (L11D2)] [Demo: Q-learning – epsilon-greedy -- crawler (L11D3)]

Exploration Functions

• When to explore?

- o Random actions: explore a fixed amount
- Better idea: explore areas whose badness is not (yet) established, eventually stop exploring

• Exploration function

• Takes a value estimate u and a visit count n, and returns an optimistic utility, e.g. f(u, n) = u + k/n

Regular Q-Update: $Q(s, a) \leftarrow_{\alpha} R(s, a, s') + \gamma \max_{a'} Q(s', a')$

Modified Q-Update: $Q(s,a) \leftarrow_{\alpha} R(s,a,s') + \gamma \max_{a'} f(Q(s',a'), N(s',a'))$

 Note: this propagates the "bonus" back to states that lead to unknown states as well!
 [Demo: exploration – Q-learning – crawler – exploration function (L11D4)]

Discussion: Model-Based vs Model-Free RL