CSE 473: Introduction to Artificial Intelligence

Hanna Hajishirzi Reinforcement Learning II

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



The Story So Far: MDPs and RL

Known MDP: Offline Solution

	Goal		Technique		
	Compute V*, Q*, π^*		Value / policy iterati	on	
	Evaluate a fixed policy $ au$	τ	Policy evaluation		
Unknown MDP: Model-Based		Unknown MDP: Model-Free			
Goal	Technique		Goal	Technique	
Compute V*, Q*, π^*	VI/PI on approx. MDP		Compute V*, Q*, π^*	Q-learning	
Evaluate a fixed policy π	PE on approx. MDP		Evaluate a fixed policy π	Value Learning	

Question? Model-Based Learning

Question: Are all episodes observed before learning the model?



Model Free: Direct Evaluation



Model Free: Temporal Difference Learning

Question: Are all episodes observed before learning the model?



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(s,a,s'))
- \circ Consider your old estimatQ(s, a)

• Consider your new sample estimate:

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

o Incorporate the new estimate into a running average

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



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

Q-Learning Demo



)

Video of Demo Q-Learning -- Crawler





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

• Full reinforcement learning: optimal policies

- You don't know the transitions T(s,a,s')
- o You don't know the rewards R(s,a,s')
- o 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



Video of Demo Q-learning – Manual Exploration – Bridge Grid



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)]

Video of Demo Q-learning – Epsilon-Greedy – Crawler



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/nRegular Q-Update: $Q(s,a) \leftarrow_{\alpha} R(s,a,s') + \gamma \max_{a'} Q(s',a') = 1$

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)]

Video of Demo Q-learning – Exploration Function – Crawler



Regret

- Even if you learn the optimal policy, you still make mistakes along the way!
- Regret is a measure of your total mistake cost: the difference between your (expected) rewards, including youthful suboptimality, and optimal (expected) rewards
- Minimizing regret goes beyond learning to be optimal – it requires optimally learning to be optimal
- Example: random exploration and exploration functions both end up optimal, but random exploration has higher regret



Announcements

- o PS3 is due Nov. 15th.
- No class on Nov. 27th! Happy Thanksgiving.
- \circ Although.... HW2 is due on Nov. 27^{th} .
 - You can only use three late day for HW2; we want to release solutions.
- o Minicontest: Due Dec. 2nd

Review

o Q-Learning



o Exploration vs. Exploitation



Approximate Q-Learning



Generalizing Across States

- Basic Q-Learning keeps a table of all q-values
- In realistic situations, we cannot possibly learn about every single state!
 - Too many states to visit them all in training
 - Too many states to hold the q-tables in memory
- Instead, we want to generalize:
 - Learn about some small number of training states from experience
 - o Generalize that experience to new, similar situations
 - This is a fundamental idea in machine learning, and we'll see it over and over again



Video of Demo Q-Learning Pacman – Tiny – Watch All



Video of Demo Q-Learning Pacman – Tiny – Silent Train



Video of Demo Q-Learning Pacman – Tricky – Watch All



Example: Pacman

Let's say we discover through experience that this state is bad:



In naïve q-learning, we know nothing about this state:



Or even this one!



[Demo: Q-learning – pacman – tiny – watch all (L11D5)] [Demo: Q-learning – pacman – tiny – silent train (L11D6)] [Demo: Q-learning – pacman – tricky – watch all (L11D7)]

Feature-Based Representations

- Solution: describe a state using a vector of features (properties)
 - Features are functions from states to real numbers (often 0/1) that capture important properties of the state
 - Example features:
 - Distance to closest ghost
 - Distance to closest dot
 - Number of ghosts
 - \circ 1 / (dist to dot)²
 - \circ Is Pacman in a tunnel? (0/1)
 - o etc.
 - Is it the exact state on this slide?
 - Can also describe a q-state (s, a) with features (e.g. action moves closer to food)



Linear Value Functions

• Using a feature representation, we can write a q function (or value function) for any state using a few weights:

$$V(s) = w_1 f_1(s) + w_2 f_2(s) + \dots + w_n f_n(s)$$

= $w_1 f_1(s, a) + w_2 f_2(s, a) + \dots + w_n f_n(s, a)$

- Advantage: our experience is summed up in a few powerful numbers
- Disadvantage: states may share features but actually be very different in value!

Approximate Q-Learning

$$Q(s,a) = w_1 f_1(s,a) + w_2 f_2(s,a) + \ldots + w_n f_n(s,a) \checkmark$$

• Q-learning with linear Q-functions:

transition =
$$(s, a, r, s')$$

difference = $[r + \gamma \max_{a'} Q(s', a')] - Q(s, a)$
 $Q(s, a) \leftarrow Q(s, a) + \alpha$ [difference] Exact Q's
 $w_i \leftarrow w_i + \alpha$ [difference] $f_i(s, a)$ Approximate Q's

• Intuitive interpretation:

- Adjust weights of active features
- E.g., if something unexpectedly bad happens, blame the features that were on: disprefer all states with that state's features
- Formal justification: online least squares

Example: Q-Pacman



-- Video of Demo Approximate Q-Learning Pacman



Q-Learning and Least Squares*



Linear Approximation: Regression*



Prediction: $(w_0) + w_1 f_1(x)$ \widehat{y}

Prediction: $\hat{y}_i = w_0 + w_1 f_1(x) + w_2 f_2(x)$

Optimization: Least Squares*



Minimizing Error*

Imagine we had only one point x, with features f(x), target value y, and weights w:

$$\frac{\partial \operatorname{error}(w)}{\partial w_{m}} = -\left(y - \sum_{k} w_{k} f_{k}(x)\right)^{2}$$

$$\frac{\partial \operatorname{error}(w)}{\partial w_{m}} = -\left(y - \sum_{k} w_{k} f_{k}(x)\right) f_{m}(x)$$

$$w_{m} \leftarrow w_{m} + \delta\left(y - \sum_{k} w_{k} f_{k}(x)\right) f_{m}(x)$$
Approximate q update explained:
$$w_{m} \leftarrow w_{m} + \alpha\left[r + \gamma \max_{a} Q(s', a')\right] - Q(s, a) f_{m}(s, a)$$
"target" "prediction"

Overfitting: Why Limiting Capacity Can Help*



Policy Search



Policy Search

- Problem: often the feature-based policies that work well (win games, maximize utilities) aren't the ones that approximate V / Q best
 - E.g. your value functions from project 2 were probably horrible estimates of future rewards, but they still produced good decisions
 - Q-learning's priority: get Q-values close (modeling)
 - Action selection priority: get ordering of Q-values right (prediction)
- Solution: learn policies that maximize rewards, not the values that predict them
- Policy search: start with an ok solution (e.g. Q-learning) then fine-tune by hill climbing on feature weights

Policy Search

- Simplest policy search:
 - o Start with an initial linear value function or Q-function
 - Nudge each feature weight up and down and see if your policy is better than before

• Problems:

- How do we tell the policy got better?
- Need to run many sample episodes!
- o If there are a lot of features, this can be impractical

 Better methods exploit lookahead structure, sample wisely, change multiple parameters...

Iteration 0



Example

o OpenAI Robotic Hand Solving Rubik Cube

https://www.youtube.com/watch?v=M1QD_3Kj7ms

New in Model-Free RL



Summary: MDPs and RL

Known MDP: Of	fline Solution	
Goal	Technique	
Compute V*, Q*, π^*	Value / policy iteration	
Evaluate a fixed policy π	Policy evaluation	

Unknown MDP: Model-Based

Goal	*use features to generalize	Technique	
Compute V*, Q*, π^*		VI/PI on approx. MDP	
Evaluate a fi>	ked policy π	PE on approx. MDP	

Unknown MDP: Model-Free

Goal	*use features to generalize	Technique	
Compute V*, Q*, π^*		Q-learning	
Evaluate a	fixed policy π	Value Learning	

Conclusion

- We're done with Part I: Search and Planning!
- We've seen how AI methods can solve problems in:
 - o Search
 - o Games
 - o Markov Decision Problems
 - o Reinforcement Learning
- Next up: Part II: Uncertainty and Learning!

