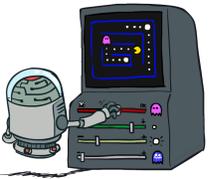


CS 188: Artificial Intelligence Reinforcement Learning II



Instructor: Nathan Lambert, University of California, Berkeley
(These slides were created by Dan Klein, Pieter Abbeel, Anca Dragan, and Nathan Lambert. <http://ai.berkeley.edu>.)

1

Reinforcement Learning

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



- o New twist: **don't know T or R**, so must try out actions
- o Big idea: **Compute all averages over T using sample outcomes**

2

The Story So Far: MDPs and RL

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

3

Approximating Values through Samples

- o Policy Evaluation:

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

- o Value Iteration:

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

- o Q-Value Iteration:

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


4

Q-Learning

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

$$Q_{k+1}(s, a) \leftarrow \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma \max_{a'} Q_k(s', a')]$$
- o Learn $Q(s,a)$ values as you go
 - o Receive a sample (s, a, s', r)
 - o Consider your old estimate: $Q(s, a)$
 - o Consider your new sample estimate:

$$\text{sample} = R(s, a, s') + \gamma \max_{a'} Q(s', a')$$
no longer policy evaluation!
 - o Incorporate the new estimate into a running average:

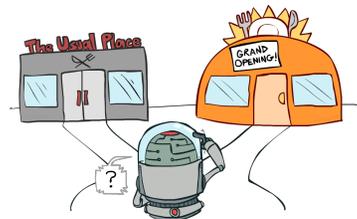
$$Q(s, a) \leftarrow (1 - \alpha)Q(s, a) + (\alpha) [\text{sample}]$$



[Demo: Q-learning - gridworld (L1002)]
[Demo: Q-learning - crawler (L1003)]

5

Exploration vs. Exploitation



6

Exploration Functions

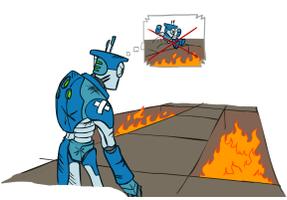
- When to explore?
 - Random actions: explore a fixed amount
 - Better idea: explore areas whose badness is not (yet) established, eventually stop exploring
- Exploration function
 - Takes a value estimate u and a visit count n , and returns an optimistic utility, e.g. $f(u, n) = u + k/n$
 - Regular 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 (L1104)]



7

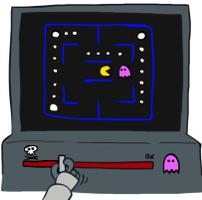
Regret

- 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



8

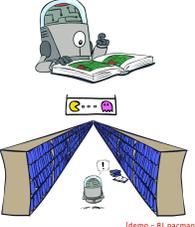
Approximate Q-Learning



9

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
 - Generalize that experience to new, similar situations
 - This is a fundamental idea in machine learning, and we'll see it over and over again



10

Example: Pacman

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



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



Or even this one!



[Demo: Q-learning - pacman - tiny - watch all (L1105)]
[Demo: Q-learning - pacman - tiny - silent train (L1106)]
[Demo: Q-learning - pacman - tricky - watch all (L1107)]

11

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



12

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



13

Feature-Based Representations

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



14

Linear Value Functions

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

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

$$Q(s, a) = w_1 f_1(s, a) + w_2 f_2(s, a) + \dots + w_n f_n(s, a)$$
- o Advantage: our experience is summed up in a few powerful numbers
- o Disadvantage: states may share features but actually be very different in value!

15

Approximate Q-Learning

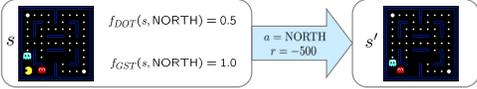
$$Q(s, a) = w_1 f_1(s, a) + w_2 f_2(s, a) + \dots + w_n f_n(s, a)$$

- o Q-learning with linear Q-functions:
 - transition = (s, a, r, s')
 - difference = [r + $\gamma \max_{a'} Q(s', a')$] - Q(s, a)
 - Exact Q's: $Q(s, a) \leftarrow Q(s, a) + \alpha [\text{difference}]$
 - Approximate Q's: $w_i \leftarrow w_i + \alpha [\text{difference}] f_i(s, a)$
- o Intuitive interpretation:
 - o Adjust weights of active features
 - o E.g., if something unexpectedly bad happens, blame the features that were on; disprefer all states with that state's features
- o Formal justification: online least squares



16

Example: Q-Pacman

$$Q(s, a) = 4.0 f_{DOT}(s, a) - 1.0 f_{GST}(s, a)$$


$f_{DOT}(s, \text{NORTH}) = 0.5$
 $f_{GST}(s, \text{NORTH}) = 1.0$

$Q(s, \text{NORTH}) = +1$
 $r + \gamma \max_{a'} Q(s', a') = -500 + 0$

difference = -501 \rightarrow $w_{DOT} \leftarrow 4.0 + \alpha [-501] 0.5$
 $w_{GST} \leftarrow -1.0 + \alpha [-501] 1.0$

$$Q(s, a) = 3.0 f_{DOT}(s, a) - 3.0 f_{GST}(s, a)$$

(Demo: approximate Q-learning, slides 14, 15, 16)

17

Video of Demo Approximate Q-Learning -- Pacman



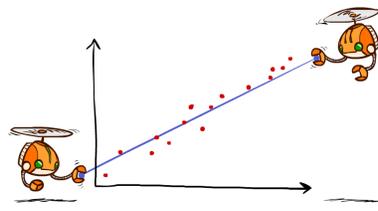
18

DeepMind Atari (©Two Minute Lectures)
approximate Q-learning with neural nets



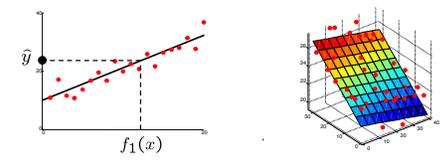
19

Q-Learning and Least Squares



20

Linear Approximation: Regression

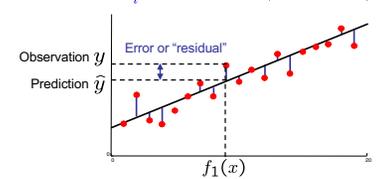


Prediction:
 $\hat{y} = w_0 + w_1 f_1(x)$

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

21

Optimization: Least Squares

$$\text{total error} = \sum_i (y_i - \hat{y}_i)^2 = \sum_i \left(y_i - \sum_k w_k f_k(x_i) \right)^2$$


Observation y

Prediction \hat{y}

Error or "residual"

$f_1(x)$

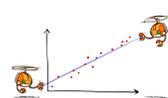
22

Minimizing Error

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

$$\text{error}(w) = \frac{1}{2} \left(y - \sum_k w_k f_k(x) \right)^2$$

$$\frac{\partial \text{error}(w)}{\partial w_m} = - \left(y - \sum_k w_k f_k(x) \right) f_m(x)$$

$$w_m \leftarrow w_m + \alpha \left(y - \sum_k w_k f_k(x) \right) f_m(x)$$


Approximate q update explained:

$$w_m \leftarrow w_m + \alpha [r + \gamma \max_{a'} Q(s', a') - Q(s, a)] f_m(s, a)$$

"target" "prediction"

23

More Powerful Function Approximation

- o Linear:
 $Q(s, a) = w_1 f_1(s, a) + w_2 f_2(s, a) + \dots + w_n f_n(s, a)$
- o Polynomial:
 $Q(s, a) = w_{11} f_1(s, a) + w_{12} f_2(s, a)^2 + w_{13} f_3(s, a)^3 + \dots$
- o Neural Network:
 $Q(s, a) = w_1 f_1(s, a) + w_2 f_2(s, a) + \dots + w_n f_n(s, a)$

Learn These Too

24

More Powerful Function Approximation

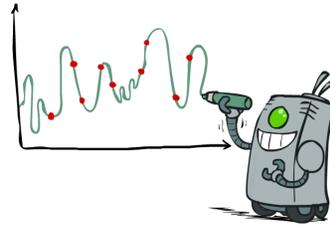
o Summary of iterative feature update equation!

$$w_m \leftarrow w_m + \alpha [r + \gamma \max_a Q(s', a') - Q(s, a)] \frac{dQ}{dw_m}(s, a)$$

= $f_m(s, a)$ when linear approximation

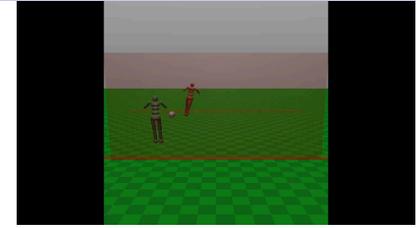
25

Overfitting: Why Limiting Capacity Can Help



26

RL: Learning Soccer



[Bansal et al, 2017]

27

RL: Learning Manipulation



[Levine*, Finn*, Darrell, Abbeel, JMLR 2016]

28

RL: NASA SUPERball



[Guez*, Zhou*, Burch*, Colson*, Mousavi*, Sutton*, Schaal, Abbeel, IJRR 2017]

29

RL: In-Hand Manipulation



Pieter Abbeel - UC Berkeley | CarseyScope | Covariance.AI

30



31



32

Policy Search

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

33

Policy Search

- o Simplest policy search:
 - o Start with an initial linear value function or Q-function
 - o Nudge each feature weight up and down and see if your policy is better than before
- o Problems:
 - o How do we tell the policy got better?
 - o Need to run many sample episodes!
 - o If there are a lot of features, this can be impractical
- o Better methods exploit lookahead structure, sample wisely, change multiple parameters...

34

The Story So Far: MDPs and RL

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

Unknown MDP: Model-Based		Unknown MDP: Model-Free	
Goal	Technique	Goal	Technique
Compute V^*, Q^*, π^*	V/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

35

Advanced Topic: Q-learning for Atari

Bach, Schrittwieser, et al. "Playing atari with deep reinforcement learning." arXiv preprint 1610.02532 (2016).

DeepMind Atari (©Two Minute Lectures)

36

Advanced Topic: Q-learning for Atari

Algorithm 1 Deep Q-learning with Experience Replay

```

Initialize replay memory  $\mathcal{D}$  to capacity  $N$ 
Initialize action-value function  $Q$  with random weights
for episode = 1,  $M$  do
  Initialize sequence  $s_1 = \{x_1\}$  and preprocessed sequence  $\phi_1 = \phi(s_1)$ 
  for  $t = 1, T$  do
    With probability  $\epsilon$  select a random action  $a_t$ 
    otherwise select  $a_t = \arg \max_a Q^*(\phi(s_t), a; \theta)$ 
    Execute action  $a_t$  in emulator and observe reward  $r_t$  and image  $x_{t+1}$ 
    Set  $s_{t+1} = s_t, a_t, x_{t+1}$  and preprocess  $\phi_{t+1} = \phi(s_{t+1})$ 
    Store transition  $(\phi_t, a_t, r_t, \phi_{t+1})$  in  $\mathcal{D}$ 
    Sample random minibatch of transitions  $(\phi_j, a_j, r_j, \phi_{j+1})$  from  $\mathcal{D}$ 
    Set  $y_j = \begin{cases} r_j & \text{for terminal } \phi_{j+1} \\ r_j + \gamma \max_{a'} Q(\phi_{j+1}, a'; \theta) & \text{for non-terminal } \phi_{j+1} \end{cases}$ 
    Perform a gradient descent step on  $(y_j - Q(\phi_j, a_j; \theta))^2$  according to equation 3
  end for
end for

```

Mnih, Volodymyr, et al. "Playing atari with deep reinforcement learning." arXiv preprint arXiv:1312.5602 (2013).

37

Advanced Topic: Q-learning for Atari

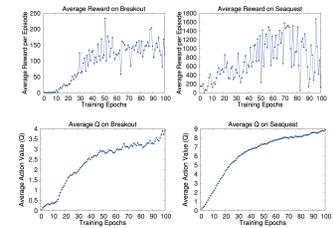
- Key implementations:
 - Neural network approximation of Q-values
 - Handling large state-action space (state compressing function)
 - Experience Replay!
- Experience Replay: idea in RL where you re-select past transitions or episodes to repeatedly train on them

Mnih, Volodymyr, et al. "Playing atari with deep reinforcement learning." arXiv preprint arXiv:1312.5602 (2013).

38

Advanced Topic: Q-learning for Atari

- Convergence of episode rewards vs Q-values



Mnih, Volodymyr, et al. "Playing atari with deep reinforcement learning." arXiv preprint arXiv:1312.5602 (2013).

39

Conclusion

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



40