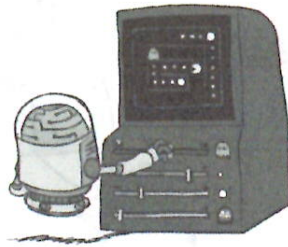


CS 188: Artificial Intelligence

Reinforcement Learning II



Instructors: Dan Klein and Pieter Abbeel --- University of California, Berkeley

[These slides were created by Dan Klein and Pieter Abbeel for CS188 Intro to AI at UC Berkeley. All CS188 materials are available at <http://ai.berkeley.edu>.]

Reinforcement Learning

We still assume an MDP:

- A set of states $s \in S$
- A set of actions (per state) A
- 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 , so must try out actions

Big idea: Compute all averages over T using sample outcomes

(final, compute estimate of T)
i.e., $\sum_s T(s, a, s')$

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

Goal	Technique
Compute V^*, Q^*, π^*	VI/PI on approx. MDP
Evaluate a fixed policy π	PE on approx. MDP

Unknown MDP: Model-Free

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

analogs of VI/PI and PE for MDPs

Model-Free Learning

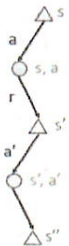
Model-free (temporal difference) learning

- Experience world through episodes

$(s, a, r, s', a', r', s'', a'', r'', s''', \dots)$

- Update estimates each transition (s, a, r, s')

- Over time, updates will mimic Bellman updates



Q-Learning

- We'd like to do Q-value updates to each Q-state:

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

- But can't compute this update without knowing T, R

- Instead, compute average as we go

- Receive a sample transition (s, a, r, s')

- This sample suggests

$$Q(s, a) \approx r + \gamma \max_{a'} Q(s', a')$$

- But we want to average over results from (s, a) (Why?)

- So keep a running average

$$Q(s, a) \leftarrow (1 - \alpha)Q(s, a) + (\alpha) [r + \gamma \max_{a'} Q(s', a')]$$

layer of optimality

any 1 sample may be misrepresentative

Q-Learning Properties

- Amazing result: Q-learning converges to optimal policy -- even if you're acting suboptimally!

- This is called off-policy learning

- Caveats:

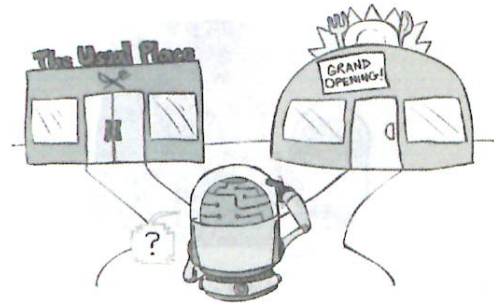
- You have to explore enough
- You have to eventually make the learning rate small enough
- ... but not decrease it too quickly
- Basically, in the limit, it doesn't matter how you select actions (!)



[Demo: Q-learning - auto - cliff grid (L11011)]

Video of Demo Q-Learning Auto Cliff Grid

Exploration vs. Exploitation



How to Explore?

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

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 $1-\epsilon$, act on current policy
- Problems with random actions?
 - You do eventually explore the space, but keep thrashing around once learning is done
 - One solution: lower ϵ over time — *still good*
 - Another solution: exploration functions



[Demo: Q-learning - manual exploration - bridge grid (L1102)]
[Demo: Q-learning - epsilon-greedy - crawler (L1103)]

• Our first foray into choosing actions

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

- Even though Crawler learned how to move forward w/ $\epsilon = 0.8$, this ϵ is so high that ~~it took~~ the probability of a chain of actions leading to forward movement was very low, so Crawler didn't move forward. Once ϵ was set to 0, Crawler followed the optimal actions for the policy it learned so far and moved forward although slowly compared to if it had trained for longer.

Exploration Functions

When to explore?

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

k is a bonus!

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: $Q(s, a) \leftarrow \alpha R(s, a, s') + (1 - \alpha) \max_{a'} Q(s', a')$

Modified Q-Update: $Q(s, a) \leftarrow \alpha R(s, a, s') + (1 - \alpha) \max_{a'} (Q(s', a') + \frac{k}{N(s', a')})$

- Note: this propagates the "bonus" back to states that lead to unknown states as well!

→ trying things that are unknown that lead to unknown states
• α means arg in RHS quantity w/ weight α

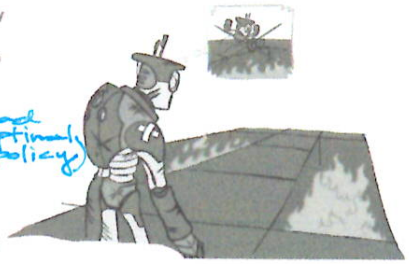


Video of Demo Q-learning – Exploration Function – Crawler

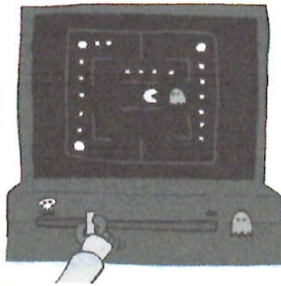
- w/ $\epsilon = 0.1$ + an exploration func., it quickly ~~learns~~ learns some good behavior more quickly adopts good behavior than before even though this good behavior is far from optimal.
- w/ exploration func., but quickly learns some states are bad + stops consciously exploring them (still may visit bc though bc $\epsilon > 0$).

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 (if you had optimal policy)
- 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



Approximate Q-Learning

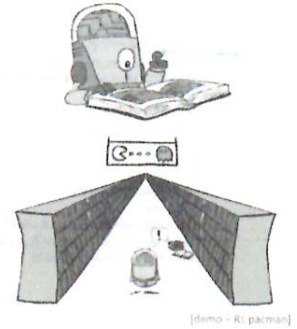


- when you learn a ghost is bad, you should transfer that to other relevant states ^{info.}

Generalizing Across States

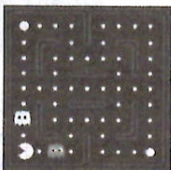
- Basic Q-Learning keeps a table of all q-values
- In realistic situations, we cannot possibly learn about every single state!
 - Too many states to visit them all in training
 - Too many states to hold the q-tables in memory
- Instead, we want to generalize:
 - Learn about some small number of training 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

• Faster learning

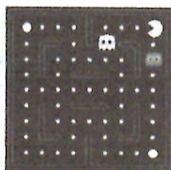


Example: Pacman

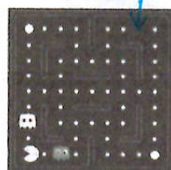
Let's say we discover through experience that this state is bad:



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



Or even this one!



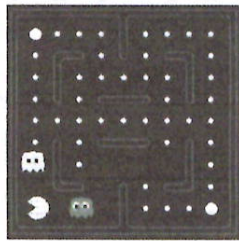
missing a dot here

[Demo Q-learning - pacman - tiny - watch all ([1105])]
 [Demo Q-learning - pacman - tiny - silent train ([1106])]
 [Demo Q-learning - pacman - tricky - watch all ([1107])]

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

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
 - $1 / (\text{dist to dot})^2$
 - Is Pacman in a tunnel? {0/1}
 - 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:

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

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

- Advantage: our experience is summed up in a few powerful numbers
- Disadvantage: states may share features but actually be very different in value!
 - 2 ghosts on same side is ~~diff~~ much less dangerous than better state than 2 ghosts ~~not~~ on diff sides, like in image to the left.
 - Must make features ~~that~~ ^{such that} discriminate diff. utility ~~states~~ ^{well} ~~action pairs~~ ^{state / action pairs}

Approximate Q-Learning

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

- Q-learning with linear Q-functions:

$$\text{transition} = (s, a, r, s')$$

$$\text{difference} = [r + \gamma \max_{a'} Q(s', a')] - Q(s, a)$$

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

$$w_i \leftarrow w_i + \alpha [\text{difference}] f_i(s, a)$$

Approximate Q's



can't do this know bc we have features

do this instead

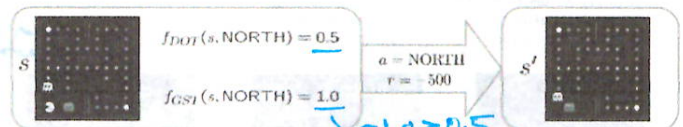
→ **Generalization**
• this is just a gradient descent update

"activation" of the feature

• How is convergence to the optimal policy affected by approx using approximate Q-Learning?
Do the features need to explain the underlying process entirely?

Example: Q-Pacman

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



$$Q(s, \text{NORTH}) = +1$$

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

$$\text{difference} = -501 \Rightarrow \begin{aligned} w_{DOT} &\leftarrow 4.0 + \alpha [-501] 0.5 \\ w_{GST} &\leftarrow -1.0 + \alpha [-501] 1.0 \end{aligned}$$

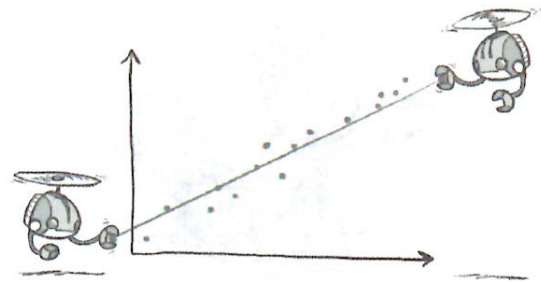
$$Q(s, a) = 3.0 f_{DOT}(s, a) - 3.0 f_{GST}(s, a)$$

[Demo: approximate Q-learning pacman (11/10/11)]

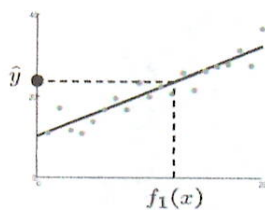
ghost feature "more active" than dot feature
credit assignment

Video of Demo Approximate Q-Learning -- Pacman

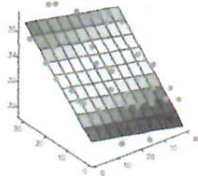
Q-Learning and Least Squares



Linear Approximation: Regression*



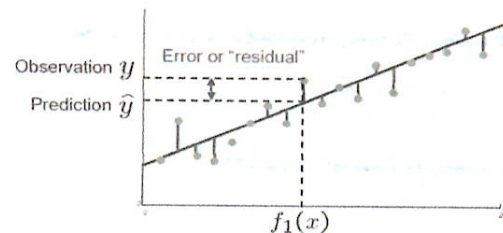
Prediction:
 $\hat{y} = w_0 + w_1 f_1(x)$



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

Optimization: Least Squares*

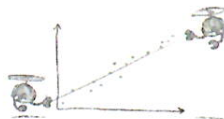
$$\text{total error} = \sum_i (y_i - \hat{y}_i)^2 = \sum_i \left(y_i - \sum_k w_k f_k(x_i) \right)^2$$



Minimizing Error*

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

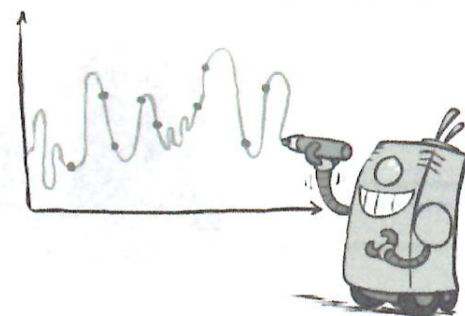
$$\begin{aligned} \text{error}(w) &= \frac{1}{2} \left(y - \sum_k w_k f_k(x) \right)^2 \\ \frac{\partial \text{error}(w)}{\partial w_m} &= \left(y - \sum_k w_k f_k(x) \right) f_m(x) \\ w_m &\leftarrow w_m + \alpha \left(y - \sum_k w_k f_k(x) \right) f_m(x) \end{aligned}$$



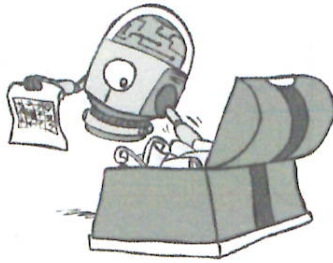
Approximate q update explained:

$$w_m \leftarrow w_m + \alpha \left[\underbrace{y}_{\text{"target"}} + \gamma \max_a Q(s', a) - \underbrace{Q(s, a)}_{\text{"prediction"}} \right] f_m(s, a)$$

Overfitting: Why Limiting Capacity Can Help*



Policy Search



- In general, Q-learning only takes you only so far
→ policy search to improve

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) → better modeling → closer accurate Q-values
 - Action selection priority: get ordering of Q-values right (prediction) → better priority ordering
 - We'll see this distinction between modeling and prediction again later in the course
 - tradeoff b/w modeling + 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:
 - 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
 - Don't need Q-val update step here
- Problems:
 - How do we tell the policy got better?
 - Need to run many sample episodes!
 - If there are a lot of features, this can be impractical
 - usually insufficient
- Better methods exploit lookahead structure, sample wisely, change multiple parameters...

Policy Search



[Andrew Ng]

[Video: HELICOPTER]

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!

