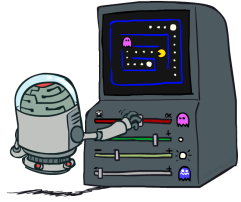


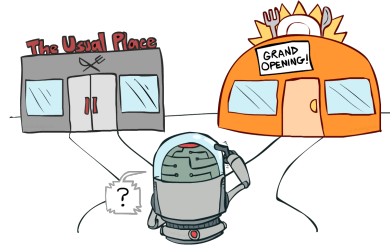
CS 473: Artificial Intelligence Reinforcement Learning II



Dieter Fox / University of Washington

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

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 $1-\epsilon$, 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
 - Another solution: exploration functions



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



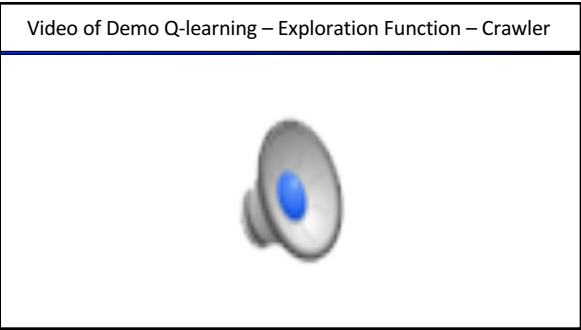
Video of Demo Q-learning – Epsilon-Greedy – Crawler



Exploration Functions

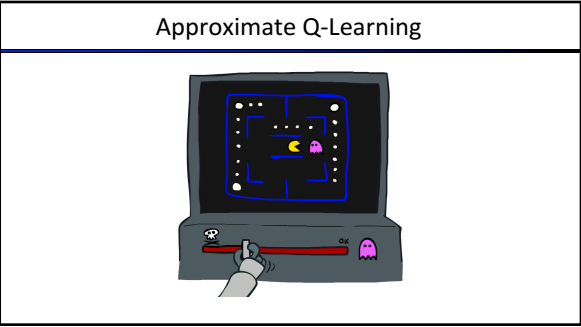
- When to explore?
 - Random actions: explore a fixed amount
 - Better idea: explore areas whose badness is not (yet) established, eventually stop exploring
- Exploration function
 - Takes a value estimate u and a visit count n , and returns an optimistic utility, e.g. $f(u, n) = u + k/n$
 - Regular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





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



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

[demo - RL pacman]

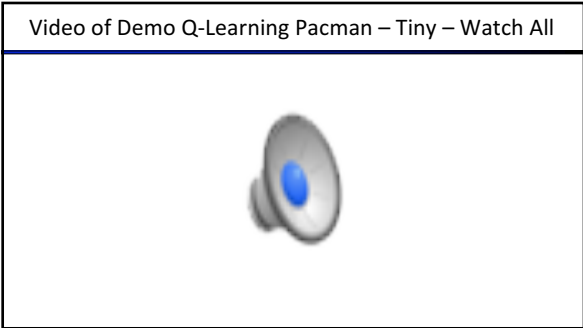
Example: Pacman

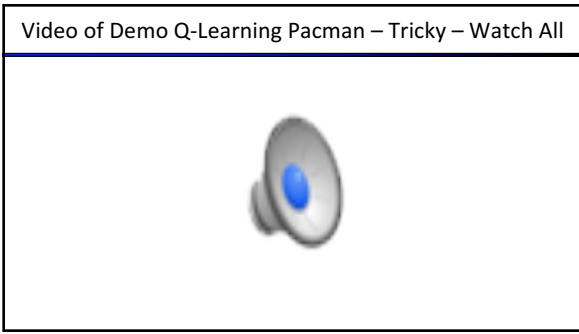
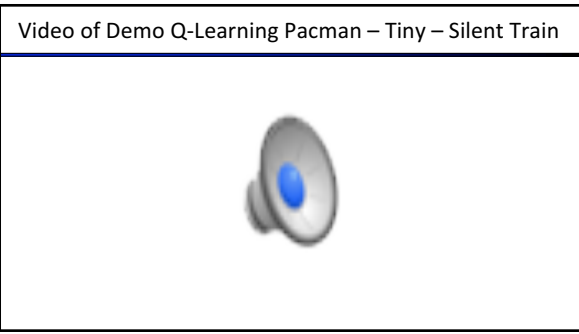
Let's say we discover through experience that this state is bad:

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** (aka "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 / (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!

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

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

S

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

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

$a = \text{NORTH}$

$r = -500$

s'

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

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

difference = -501

$w_{DOT} \leftarrow 4.0 + \alpha [-501] 0.5$

$w_{GST} \leftarrow -1.0 + \alpha [-501] 1.0$

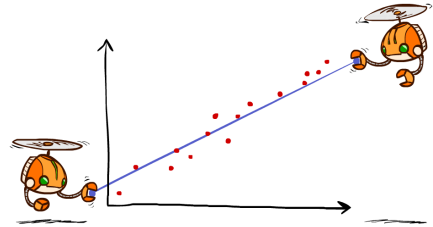
$$Q(s, a) = 3.0 f_{DOT}(s, a) - 3.0 f_{GST}(s, a)$$

[Demo: approximate Q-learning pacman (131D10)]

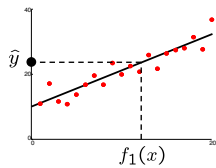
Video of Demo Approximate Q-Learning -- Pacman



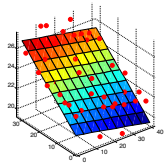
Q-Learning and Least Squares



Linear Approximation: Regression*



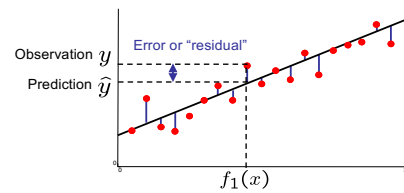
Prediction:
 $\hat{y} = w_0 + w_1 f_1(x)$



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

Optimization: Least Squares*

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



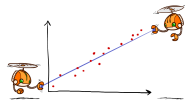
Minimizing Error*

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

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

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

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

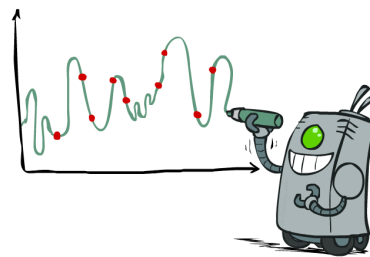


Approximate q update explained:

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

"target" "prediction"

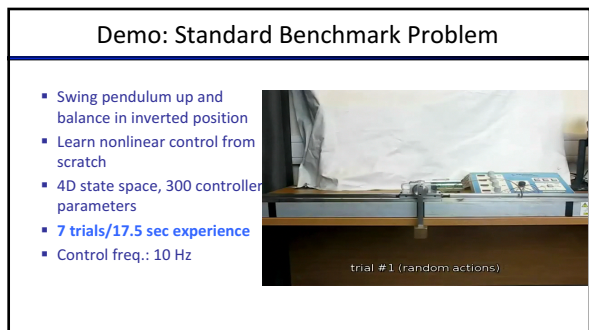
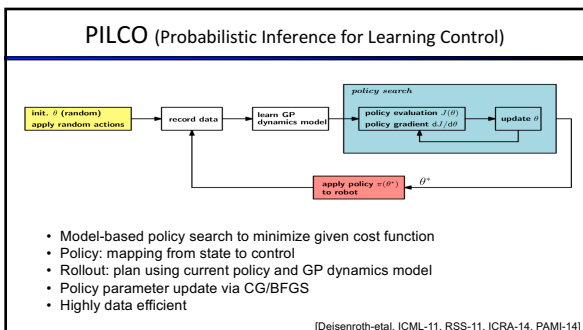
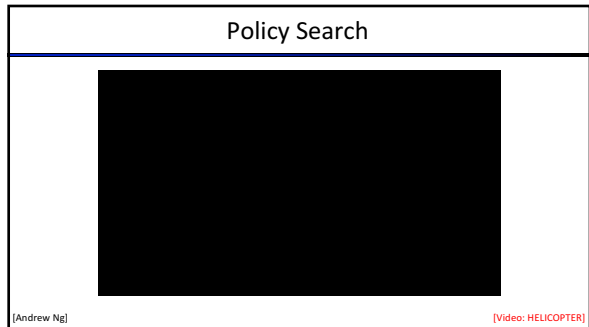
Overfitting: Why Limiting Capacity Can Help*





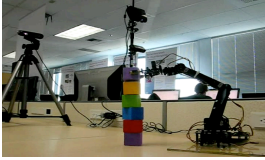
- ### Policy Search
- Problem: often the feature-based policies that work well (win games, maximize utilities) aren't the ones that approximate V / Q best
 - E.g. your value functions from project 2 were probably horrible estimates of future rewards, but they still produced good decisions
 - Q-learning's 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:
 - 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!
 - If there are a lot of features, this can be impractical
 - Better methods exploit lookahead structure, sample wisely, change multiple parameters...



Controlling a Low-Cost Robotic Manipulator

- **Low-cost** system (\$500 for robot arm and Kinect)
- **Very noisy**
- No sensor information about robot's joint configuration used
- **Goal:** Learn to stack tower of 5 blocks from scratch
- Kinect camera for tracking block in end-effector
- State: coordinates (3D) of block center (from Kinect camera)
- 4 controlled DoF
- 20 learning trials for stacking 5 blocks (5 seconds long each)
- Account for **system noise**, e.g.,
 - Robot arm
 - Image processing



Playing Atari with Deep Reinforcement Learning

Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Alex Graves, Iason Achilleopoulos, John Graves, Martin Riedmiller

DeepMind Technologies

[arXiv:1312.5680v1 [cs.LG], 2013] [arXiv:1312.5680v2 [cs.LG], 2014]

Abstract

We present the first deep learning model to successfully learn control policies for nearly the full domain of popular Atari games using reinforcement learning. The model is a convolutional neural network, trained with a variant of Q-learning, whose input is raw pixels and whose output is a value function for each game state. We show that the resulting policy generalizes across games, learning to play 49 out of 49 games, with an achievement of the maximum score in 26 games. We include qualitative and quantitative analysis of the generalization.

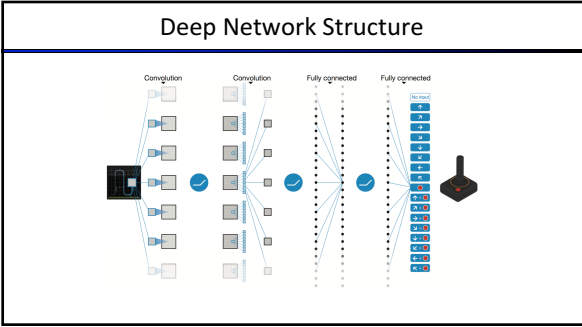
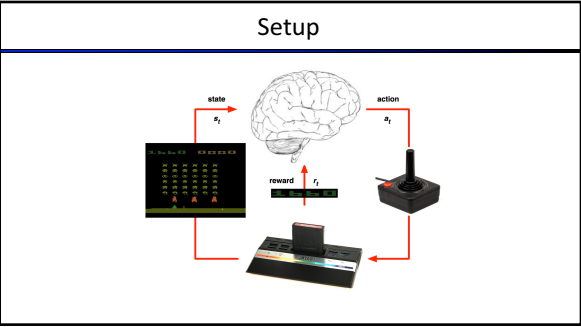
1 Introduction

Learning to control agents directly from high-dimensional sensory inputs like vision and speech is one of the long-standing challenges of computational learning theory. Many successful algorithms have been proposed in the domain of control, but most of them require hand-crafted features combined with some form of search or greedy optimization. Thanks to the development of such algorithms based on ideas from the theory of reinforcement learning, it is now possible to train high-level policies from raw sensory data in a wide range of control tasks, including convolutional networks, multiple representation and dual representation methods, and model-based methods, which can learn to control a wide range of tasks, such as playing Atari games, without any human-curated knowledge.

Recent advances in deep learning have made it possible to extract high-level features from raw sensory data using deep learning architectures. In this paper, we combine these advances with reinforcement learning to learn control policies for a wide range of Atari games. We use a deep convolutional neural network to extract features from raw pixels, and a deep Q-network to learn the value function for each game state. The resulting policy generalizes across games, learning to play 49 out of 49 games, with an achievement of the maximum score in 26 games. We include qualitative and quantitative analysis of the generalization.

Deep reinforcement learning presents several challenges. First, a deep learning perspective directly over-regularized learning prevents generalization to the non-regularized target domain. Second, the high-dimensional sensory input makes it difficult to learn a value function. Third, the high-dimensional sensory input makes it difficult to learn a policy. We address these challenges by introducing a deep convolutional neural network and a deep Q-network. Another issue is that deep learning architectures are not designed for reinforcement learning. We address this by introducing a deep convolutional neural network and a deep Q-network. Finally, we address the issue of generalization by introducing a deep convolutional neural network and a deep Q-network.

The paper demonstrates that a convolutional neural network can overcome these challenges to learn control policies for a wide range of Atari games, with a value function that generalizes across games. The network is



Deep learning: Representation

Sequence of functions parameterized via w_i

$$x \xrightarrow{w_1} h_1 \xrightarrow{w_2} h_2 \cdots h_{n-1} \xrightarrow{w_n} h_n \xrightarrow{w_{n+1}} y$$

$$\frac{\partial h_1}{\partial x} \leftarrow \frac{\partial h_2}{\partial h_1} \leftarrow \frac{\partial h_3}{\partial h_2} \cdots \frac{\partial h_n}{\partial h_{n-1}} \leftarrow \frac{\partial y}{\partial h_n} \leftarrow \frac{\partial L}{\partial y}$$

$$\downarrow \quad \downarrow \quad \downarrow \quad \downarrow \quad \downarrow$$

$$\frac{\partial h_1}{\partial w_1} \quad \frac{\partial h_2}{\partial w_2} \quad \frac{\partial h_{n-1}}{\partial w_{n-1}} \quad \frac{\partial h_n}{\partial w_n} \quad \frac{\partial y}{\partial w_{n+1}}$$

Gradients determined via chain rule / backpropagation

11/1/17 35

Deep learning: Supervised training via SGD

Given (x, y^*)

$$x \xrightarrow{w_1} h_1 \xrightarrow{w_2} h_2 \cdots h_{n-1} \xrightarrow{w_n} h_n \xrightarrow{w_{n+1}} y \longrightarrow L(y, y^*)$$

$$\frac{\partial h_1}{\partial x} \leftarrow \frac{\partial h_2}{\partial h_1} \leftarrow \frac{\partial h_3}{\partial h_2} \cdots \frac{\partial h_n}{\partial h_{n-1}} \leftarrow \frac{\partial y}{\partial h_n} \leftarrow \frac{\partial L}{\partial y}$$


$$\downarrow \quad \downarrow \quad \downarrow \quad \downarrow \quad \downarrow$$

$$\frac{\partial h_1}{\partial w_1} \quad \frac{\partial h_2}{\partial w_2} \quad \frac{\partial h_{n-1}}{\partial w_{n-1}} \quad \frac{\partial h_n}{\partial w_n} \quad \frac{\partial y}{\partial w_{n+1}}$$

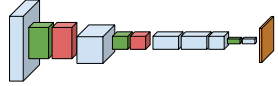
Update $w_i \leftarrow w_i - \alpha \frac{\partial L}{\partial w_i}$


11/1/17 36

Learning to Detect Hands and Parts



Input depth





Hand bounding box
(2D + hw)

$$Conv(D, g, a, c) = \sum_{k,l} D(x+i, y+j, k+l, c) \cdot K_{i,j,k,l}$$

$$ReLU(D, x, y) = \max(D(x, y), 0)$$


$$MaxPool(D, x, y, c) = \max_{i,j} (D(x+i, y+j, c))$$

$$Fully\ Connected(D, x, y, c) = \sum_{i,j} (D(x+i, y+j, c))$$

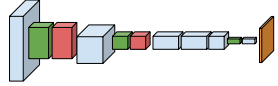
Conv + ReLU
Max Pool
Local Response Norm
Discard
Fully Connected + ReLU
Reshape


11/1/17 37

Learning to Detect Hands and Parts




Input depth






Hand bounding box
(2D + hw)




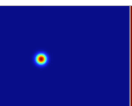
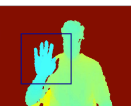
Bounding box

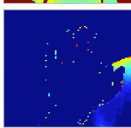

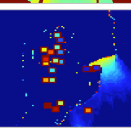


Detected parts
2D heatmap + distance

11/1/17 38

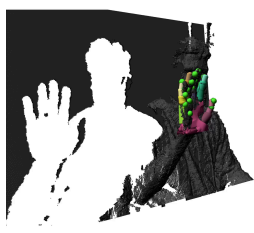
Learning to Detect Hands and Parts

11/1/17 39


Model-based refinement



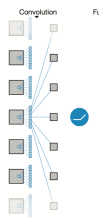
- Green dots: detected by deep net
- Colored skeleton: matched via DART

11/1/17 40

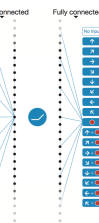
Deep Network Structure




Convolution



Convolution



Fully connected



Fully connected

11/1/17 41

Deepmind AI Playing Atari



11/1/17 42

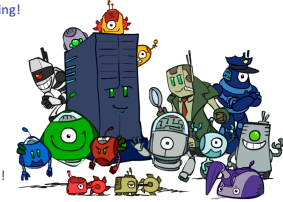
That's all for Reinforcement Learning!



- Very tough problem: How to perform any task well in an unknown, noisy environment!
- Traditionally used mostly for robotics, but becoming more widely used
- Lots of open research areas:
 - How to best balance exploration and exploitation?
 - How to deal with cases where we don't know a good state/feature representation?

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!



Midterm Topics

- Agency: types of agents, types of environments
- Search
 - Formulating a problem in terms of search
 - Algorithms: DFS, BFS, IDS, best-first, uniform-cost, A*, local
 - Heuristics: admissibility, consistency, creation
 - Constraints: formulation, search, forward checking, arc-consistency, structure
 - Adversarial: min/max, alpha-beta, expectimax
- MDPs
 - Formulation, Bellman eqns, V^* , Q^* , backups, value iteration, policy iteration

45