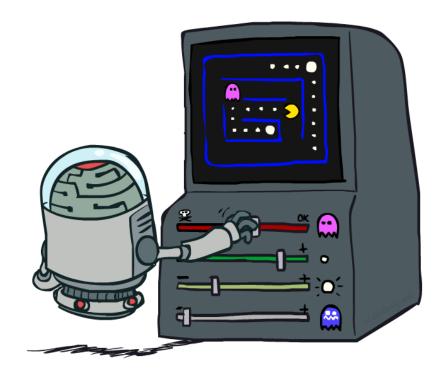
AI: Representation and Problem Solving

Reinforcement Learning II



Instructor: Pat Virtue

Slide credits: CMU AI and http://ai.berkeley.edu

Overview: MDPs and Reinforcement Learning

Known MDP: Offline Solution

Value Iteration / Policy Iteration

R(5 a 5) T(5 a 5)

Unknown MDP: Online Learning

Model-Based

Estimate MDP T(s,a,s') and R(s,a,s') from samples of environment

Model-Free

Passive Reinforcement Learning

- Direct Evaluation (simple)
- TD Learning

Active Reinforcement Learning

Q-Learning

- We are given a policy It to use.

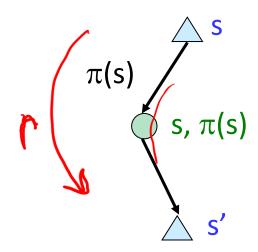
Online Learning Model-free Learning

Passive Reinforcement Learning
Temporal Difference Learning

Temporal Difference Learning

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



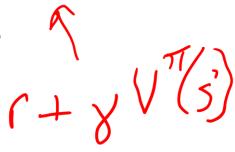
Temporal difference learning of values

- Policy still fixed, still doing evaluation!
- Move values toward value of whatever successor occurs: running average

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

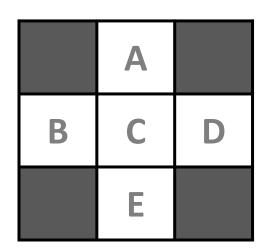
Update to V(s):
$$V^{\pi}(s) \leftarrow (1-\alpha)V^{\pi}(s) + (\alpha)sample$$

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

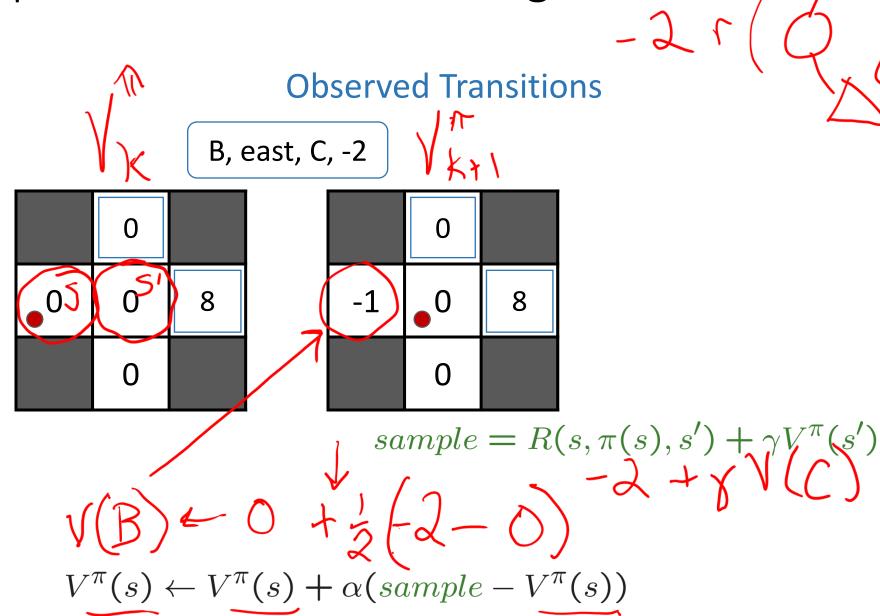


Example: Temporal Difference Learning

States

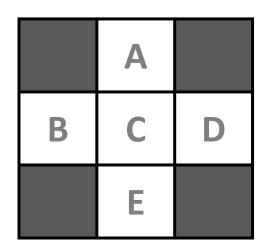


Assume: $\gamma = 1$, $\alpha = 1/2$

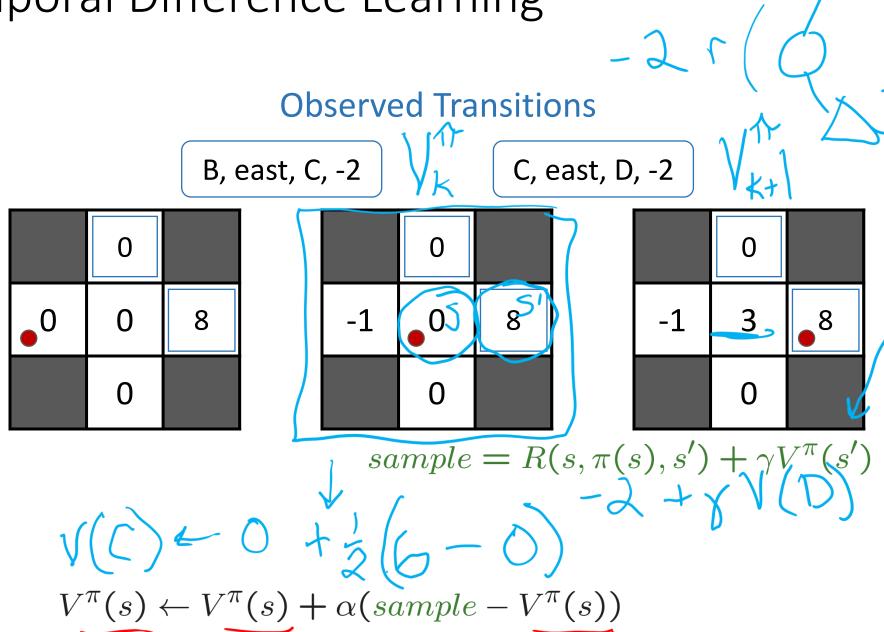


Example: Temporal Difference Learning

States



Assume: $\gamma = 1$, $\alpha = 1/2$



Temporal Difference Learning

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

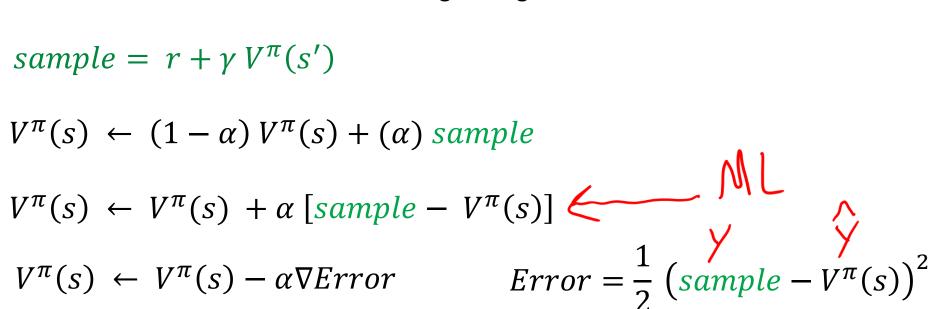
- Policy still fixed, still doing evaluation!
- Move values toward value of whatever successor occurs: running average

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

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

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

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



$$\pi(s)$$
 s
 s
 s
 s'

ML detour: Quick Calculus Quiz

$$f(x) = \frac{1}{2}(y - x)^2$$

What is
$$\frac{df}{dx}$$
? $= \left(y - x \right) \cdot \left(-1 \right)$

ML detour: Gradient Descent

Goal: find x that minimizes f(x)

- 1. Start with initial guess, x_0
- 2. Update x by taking a step in the direction that f(x) is changing fastest (in the negative direction) with respect to x:
 - $x \leftarrow x \alpha \nabla_x f$, where α is the step size or learning rate
- 3. Repeat until convergence

TD goal: find value(s), V, that minimizes difference between sample(s)

and V

$$V \leftarrow V - \alpha \nabla_V Error$$

$$Error(V) = \frac{1}{2} (sample - V)^2$$

$$f(x) = \frac{1}{2}(y - x)^2$$

$$\frac{df}{dx} = -(y - x)$$

$$\frac{df}{dx} = -(y - x)$$

Temporal Difference Learning

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

- Policy still fixed, still doing evaluation!
- Move values toward value of whatever successor occurs: running average

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

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

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

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

$$V^{\pi}(s) \leftarrow V^{\pi}(s) - \alpha \nabla Error$$

$$Error = \frac{1}{2} \left(sample - V^{\pi}(s) \right)^{2}$$

$$\pi(s)$$
 s
 s
 s
 s'

Poll 1

TD update:

$$V^{\pi}(s) = V^{\pi}(s) + \alpha [r + \gamma V^{\pi}(s') - V^{\pi}(s)]$$

Which converts TD values into a policy?

A) Value iteration:

$$V_{k+1}(s) = \max_{a} \sum_{s'} P(s'|s,a) [R(s,a,s') + \gamma V_k(s')], \quad \forall$$

B) Q-iteration:

$$V_{k+1}(s) = \max_{a} \sum_{s'} P(s'|s,a) [R(s,a,s') + \gamma V_k(s')], \quad \forall s$$

$$Q_{k+1}(s,a) = \sum_{s'} P(s'|s,a) [R(s,a,s') + \gamma \max_{a'} Q_k(s',a')], \quad \forall s, a$$

C) Policy extraction:

$$\pi_V(s) = \underset{a}{\operatorname{argmax}} \sum_{s'} P(s'|s,a) [R(s,a,s') + \gamma V(s')], \quad \forall s$$

D) Policy evaluation:

$$V_{k+1}^{\pi}(s) = \sum_{s'} P(s'|s,\pi(s))[R(s,\pi(s),s') + \gamma V_k^{\pi}(s')], \quad \forall s$$

E) Policy improvement:
$$\pi_{new}(s) = \underset{a}{\operatorname{argmax}} \sum_{s'} P(s'|s,a) [R(s,a,s') + \gamma V^{\pi_{old}}(s')], \quad \forall s'$$

F) None of the above

Problems with TD Value Learning

TD value leaning is a model-free way to do policy evaluation, mimicking Bellman updates with running sample averages

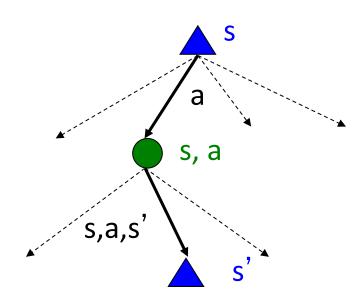
However, if we want to turn values into a (new) policy, we're sunk:

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

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

Idea: learn Q-values, not values

Makes action selection model-free too!



MDP/RL Notation

$$V(s) = \max_{a} \sum_{s'} P(s'|s,a)V(s')$$

$$V^*(s) = \max_{a} \sum_{s'} P(s'|s,a) [R(s,a,s') + \gamma V^*(s')]$$

$$V_{k+1}(s) = \max_{a} \sum_{s'} P(s'|s,a) [R(s,a,s') + \gamma V_k(s')], \quad \forall s'$$

$$Q_{k+1}(s,a) = \sum_{s'} P(s'|s,a) [R(s,a,s') + \gamma \max_{a'} Q_k(s',a')], \quad \forall \, s,a$$

$$\pi_V(s) = \underset{a}{\operatorname{argmax}} \sum_{s'} P(s'|s, a) [R(s, a, s') + \gamma V(s')], \quad \forall s$$

$$V_{k+1}^{\pi}(s) = \sum_{s'} P(s'|s,\pi(s))[R(s,\pi(s),s') + \gamma V_k^{\pi}(s')], \quad \forall s$$

$$\pi_{new}(s) = \underset{a}{\operatorname{argmax}} \sum_{s'} P(s'|s, a) [R(s, a, s') + \gamma V^{\pi_{old}}(s')], \quad \forall s'$$

$$V^{\pi}(s) = V^{\pi}(s) + \alpha [r + \gamma V^{\pi}(s') - V^{\pi}(s)]$$

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



-> Agent gets to choose next action

Online Learning Model-free Learning



Active Reinforcement Learning Q-Learning

Q-Learning

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

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

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

Instead, compute average as we go

- Receive a sample transition (s,a,r,s')
- The sample value is then:

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

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

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



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 (!)



Overview: MDPs and Reinforcement Learning

Known MDP: Offline Solution

Value Iteration / Policy Iteration

Unknown MDP: Online Learning

Model-Based

Estimate MDP T(s,a,s') and R(s,a,s') from samples of environment

Model-Free

Passive Reinforcement Learning

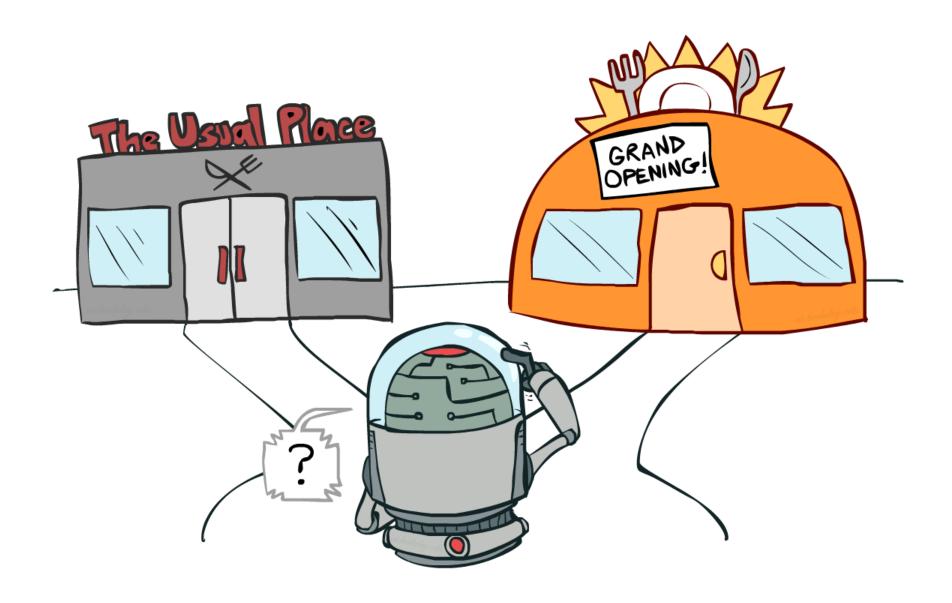
- Direct Evaluation (simple)
- TD Learning

Active Reinforcement Learning

Q-Learning

Demo Q-Learning Auto Cliff Grid

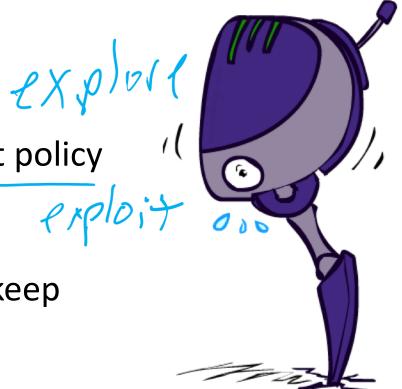
Exploration vs. Exploitation



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
- 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



Demo Q-learning – Manual Exploration – Bridge Grid

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

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

[Demo: exploration – Q-learning – crawler – exploration function (L11D4)]

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.

$$f(u,n) = u + k/(n+1)$$

Regular Q-Update: $Q(s,a) = Q(s,a) + \alpha \left[r + \gamma \max_{a'} Q(s',a') - Q(s,a)\right]$

Modified Q-Update: $Q(s, a) = Q(s, a) + \alpha \left[r + \gamma \max_{a'} f(Q(s', a'), N(s', a')) - Q(s, a)\right]$

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



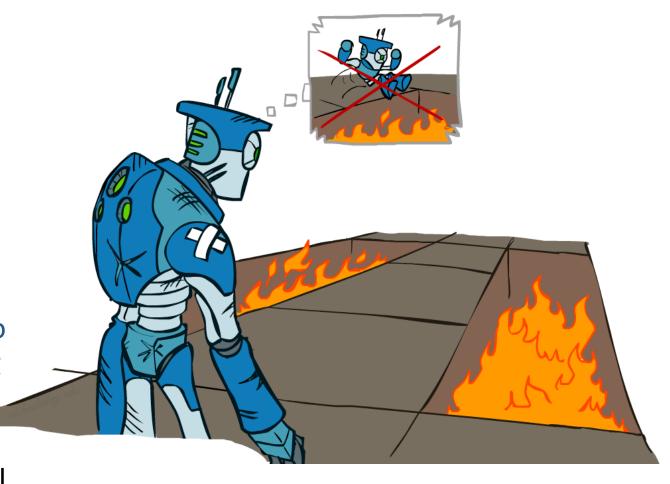
Regret

Even if you learn the optimal policy, you still make mistakes along the way!

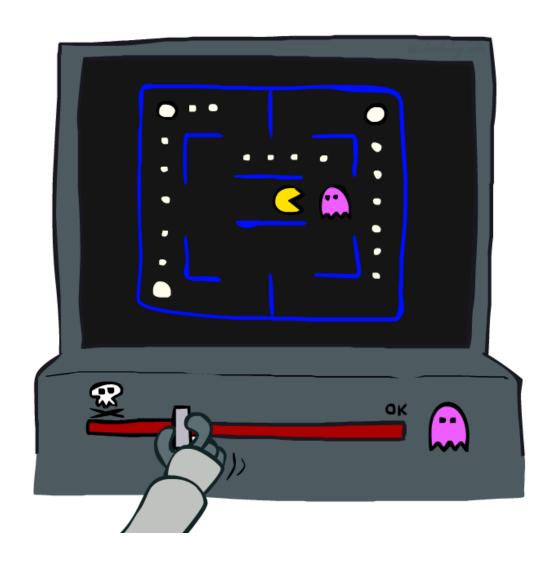
Regret is a measure of your total mistake cost: the difference between your (expected) rewards, including youthful suboptimality, and optimal (expected) rewards

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



Approximate Q-Learning



Approximate Q-Learning

What happens when we change Candy Grab to start with 1000 pieces?

Pieces Available	Take 1	Take 2
2	0%	100%
3	2%	0%
4	75 %	2%
5	4%	68%
6	5%	6%
7	60%	5%

Generalizing Across States

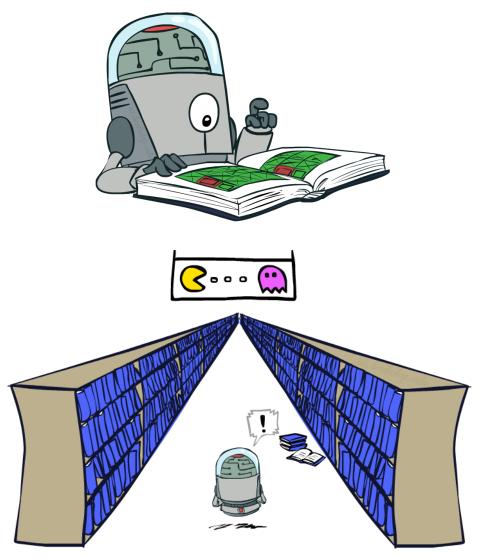
Basic Q-Learning keeps a table of all q-values

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:

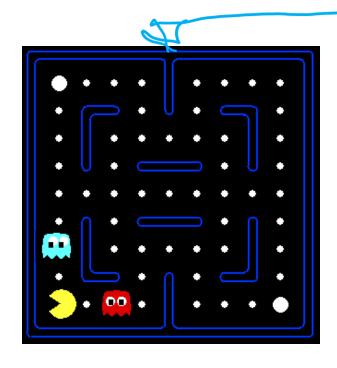
- Learn about some small number of training states from experience
- Generalize that experience to new, similar situations
- This is a fundamental idea in machine learning, and we'll see it over and over again

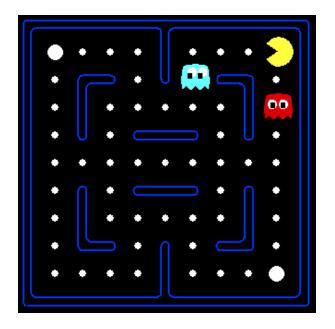


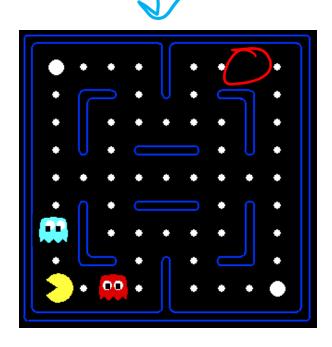
Example: Pacman

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)]

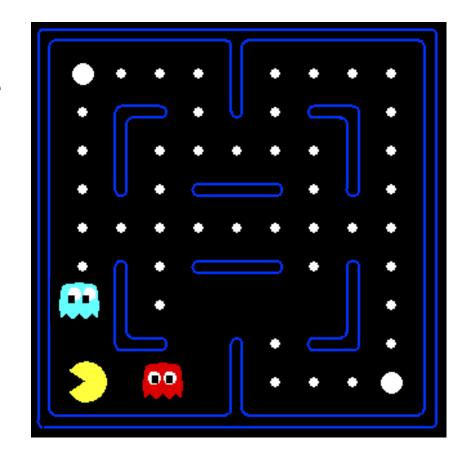
[Demo: Q-learning – pacman – tiny – silent train (L11D6)]

[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
 - 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:

$$V_{w}(s) = W_{1}f_{1}(s) + W_{2}f_{2}(s) + ... + W_{n}f_{n}(s)$$

$$Q_{w}(s,a) = W_{1}f_{1}(s,a) + W_{2}f_{2}(s,a) + ... + W_{n}f_{n}(s,a)$$

Advantage: our experience is summed up in a few powerful numbers

Disadvantage: states may share features but actually be very different in value!

Updating a linear value function

Original Q learning rule tries to reduce prediction error at s, a:

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

Instead, we update the weights to try to reduce the error at s, a:

Quick Calculus Quiz
$$Error(w) = \frac{1}{2} (y - wf(x))^{2}$$

What is
$$\frac{dError}{dw}$$
? -1 $(y-wf(x))$

Last time
$$\frac{1}{2} = \frac{1}{2} (y - x)^{2}$$

$$\frac{dError}{dx} = -(y - x)$$

Updating a linear value function

Original Q learning rule tries to reduce prediction error at s, a:

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

Instead, we update the weights to try to reduce the error at s, a:

•
$$w_i \leftarrow w_i + \alpha \cdot [r + \gamma \max_{a'} Q(s',a') - Q(s,a)] \underline{\partial Q_w(s,a)/\partial w_i}$$

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

$$Q_{w}(s,a) = w_{1}f_{1}(s,a) + w_{2}f_{2}(s,a)$$

$$\frac{\partial Q}{\partial w_{2}} = \int_{\mathcal{X}} \left(\mathcal{S}_{y}(s,a) \right) ds$$

$$Error(w) = \frac{1}{2} (y - wf(x))^2$$

$$\frac{dError}{dw} = -(y - wf(x))f(x)$$

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:

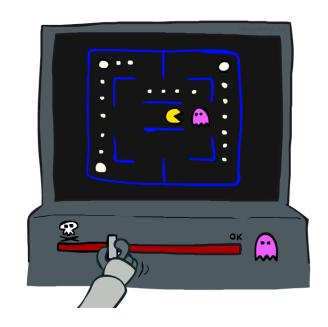
transition
$$= (s, a, r, s')$$

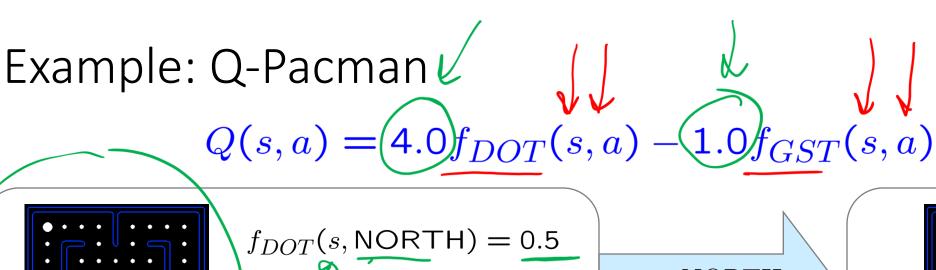
difference $= \left[r + \gamma \max_{a'} Q(s', a') \right] - Q(s, a)$
 $Q(s, a) \leftarrow Q(s, a) + \alpha$ [difference] Exact Q's
 $w_i \leftarrow w_i + \alpha$ [difference] $f_i(s, a)$ Approximate Q's

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

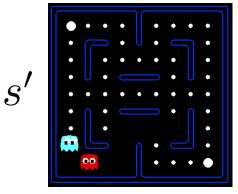




S

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

a = NORTHr = -500



$$Q(s',\cdot)=0$$

$$Q(s, NORTH) = +1$$

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

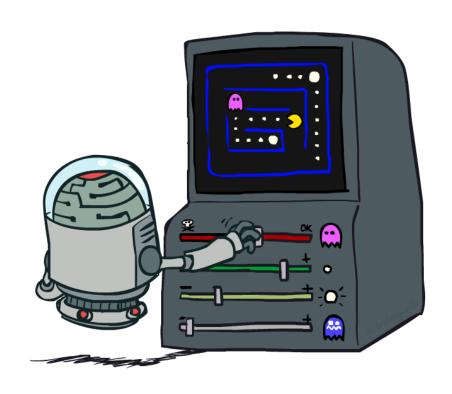
$$difference = -501$$

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

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

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

Reinforcement Learning Milestones



TDGammon

1992 by Gerald Tesauro, IBM

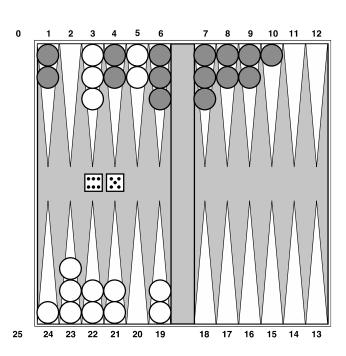
4-ply lookahead using V(s) trained from 1,500,000 games of self-play

3 hidden layers, ~100 units each

Input: contents of each location plus several handcrafted features

Experimental results:

- Plays approximately at parity with world champion
- Led to radical changes in the way humans play backgammon



Deep Q-Networks

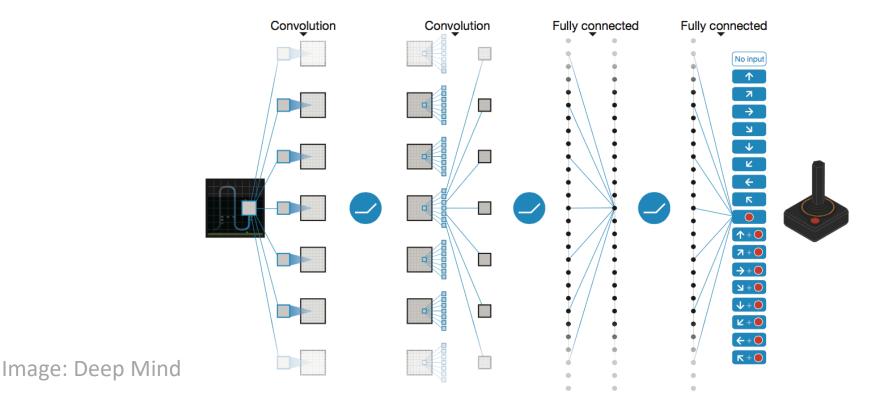
sample = $r + \gamma \max_{a'} Q_w(s',a')$ $Q_w(s,a)$: Neural network

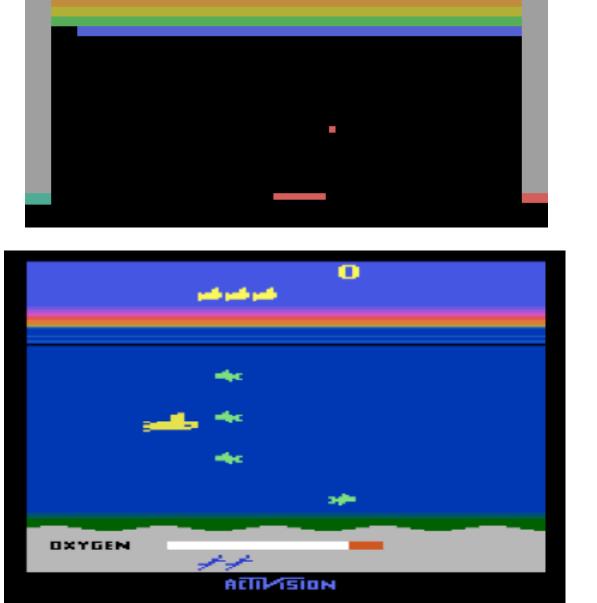
Deep Mind, 2015

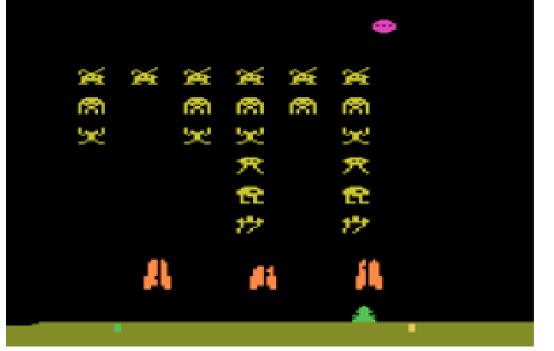
Used a deep learning network to represent Q:

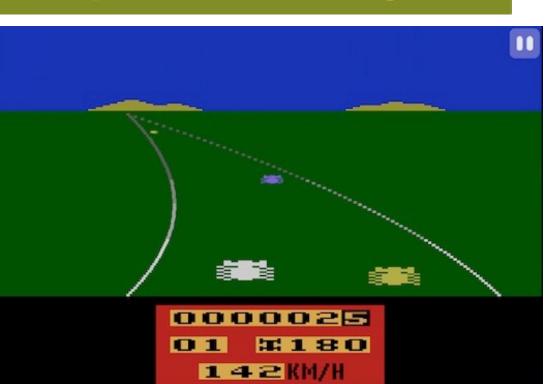
■ Input is last 4 images (84x84 pixel values) plus score

49 Atari games, incl. Breakout, Space Invaders, Seaquest, Enduro









Images: Open AI, Atari

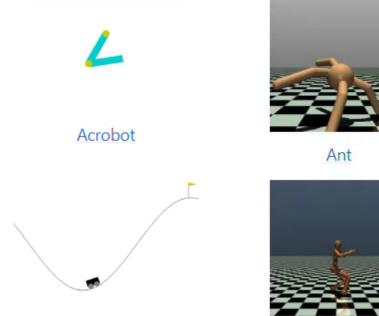
OpenAl Gym (now Gymnasium)

2016+

Benchmark problems for learning agents

Mountain Car

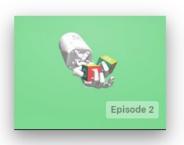
https://gymnasium.farama.org/



Humanoid



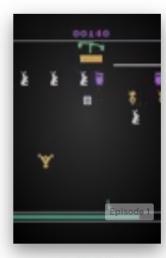




HandManipulateBlock-v0 Orient a block using a robot hand



Breakout-ram-v0 Maximize score in the game Breakout, with RAM as input



Carnival-v0 Maximize score in the game Carnival, with screen images as input

Images: Open AI/Gymnasium

AlphaGo, AlphaZero

Deep Mind, 2016+



Autonomous Vehicles?