**Slide 1**

# Reinforcement Learning

CSE 415: Introduction to Artificial Intelligence
University of Washington
Winter 2018

Presented by S. Tanimoto, University of Washington,
based on material by Dan Klein and Pieter Abbeel -- University of California.

---

**Slide 2**

# Outline

- Planning vs Learning
- Model-Based Learning
- Direct Evaluation
- Sample-Based Policy Evaluation
- Temporal Difference Learning
- Active Reinforcement Learning
- Q-Learning
- Exploration vs Exploitation
- Regret

2

---

**Slide 3**

# 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 |

---

**Slide 4**

# Example: Model-Based Learning

**Input Policy π**

Assume: γ = 1

**Observed Episodes (Training)**

Episode 1
B, east, C, -1
C, east, D, -1
D, exit, x, +10

Episode 2
B, east, C, -1
C, east, D, -1
D, exit, x, +10

Episode 3
E, north, C, -1
C, east, D, -1
D, exit, x, +10

Episode 4
E, north, C, -1
C, east, A, -1
A, exit, x, -10

**Learned Model**

$\hat{T}(s, a, s')$

T(B, east, C) = 1.00
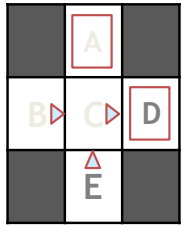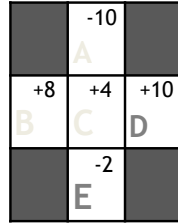T(C, east, D) = 0.75
T(C, east, A) = 0.25
...

$\hat{R}(s, a, s')$

R(B, east, C) = -1
R(C, east, D) = -1
R(D, exit, x) = +10
...

## Example: Direct Evaluation

Input Policy π    Observed Episodes (Training)    Output Values



Episode 1

B, east, C, -1
C, east, D, -1
D, exit, x, +10

Episode 2

B, east, C, -1
C, east, D, -1
D, exit, x, +10

Episode 3

E, north, C, -1
C, east,  D, -1
D, exit, x, +10

Episode 4

E, north, C, -1
C, east,  A, -1
A, exit,  x, -10

*Assume:* $\gamma = 1$

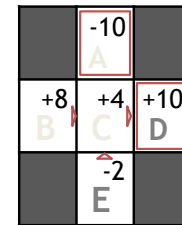Output grid:
-10 (A)
+8 (B), +4 (C), +10 (D)
-2 (E)

5

---

## Problems with Direct Evaluation

- What's good about direct evaluation?
  - It's easy to understand
  - It doesn't require any knowledge of T, R
  - It eventually computes the correct average values, using just sample transitions

- What bad about it?
  - It wastes information about state connections
  - Each state must be learned separately
  - So, it takes a long time to learn

Output Values



-10 (A)
+8 (B), +4 (C), +10 (D)
-2 (E)

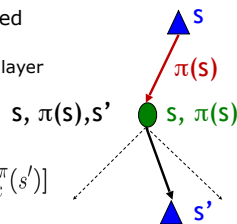*If B and E both go to C under this policy, how can their values be different?*

6

---

## Why Not Use Policy Evaluation?

- Simplified Bellman updates calculate V for a fixed policy:
  - Each round, replace V with a one-step-look-ahead layer over V

$$V_0^\pi(s) = 0$$

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

  - This approach fully exploited the connections between the states
  - Unfortunately, we need T and R to do it!

- Key question: how can we do this update to V without knowing T and R?
  - In other words, how to we take a weighted average without knowing the weights?



7

---

## Sample-Based Policy Evaluation?

- We want to improve our estimate of V by computing these averages:

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

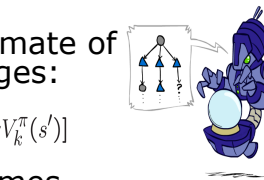- Idea: Take samples of outcomes s' (by doing the action!) and average

$$sample_1 = R(s, \pi(s), s_1') + \gamma V_k^\pi(s_1')$$

$$sample_2 = R(s, \pi(s), s_2') + \gamma V_k^\pi(s_2')$$

$$\dots$$

$$sample_n = R(s, \pi(s), s_n') + \gamma V_k^\pi(s_n')$$

$$V_{k+1}^\pi(s) \leftarrow \frac{1}{n} \sum_i sample_i$$
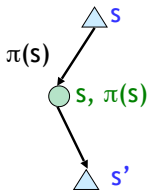


*Almost! But we can't rewind time to get sample after sample from state s.*

8

## 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

$\pi(s)$

s

s, $\pi$(s)

s'

Sample of V(s): $sample = R(s, \pi(s), s') + \gamma V^\pi(s')$

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

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

---

## Exponential Moving Average

- Exponential moving average
  - The running interpolation update:

$$\bar{x}_n = (1-\alpha) \cdot \bar{x}_{n-1} + \alpha \cdot x_n$$

  - Makes recent samples more important:

$$\bar{x}_n = \frac{x_n + (1-\alpha) \cdot x_{n-1} + (1-\alpha)^2 \cdot x_{n-2} + \dots}{1 + (1-\alpha) + (1-\alpha)^2 + \dots}$$

  - Forgets about the past (distant past values were wrong anyway)

- Decreasing learning rate (alpha) can give converging averages

---

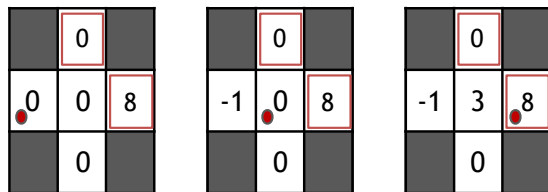## Example: Temporal Difference Learning

States | Observed Transitions

B, east, C, -2 | C, east, D, -2

| | 0 | | | 0 | | | 0 | |
|---|---|---|---|---|---|---|---|---|
| 0 | 0 | 8 | -1 | 0 | 8 | -1 | 3 | 8 |
| | 0 | | | 0 | | | 0 | |

A

B C D

E

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

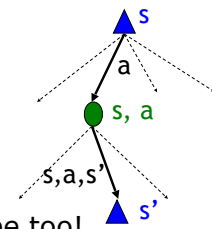$$V^\pi(s) \leftarrow (1-\alpha)V^\pi(s) + \alpha \left[ R(s, \pi(s), s') + \gamma V^\pi(s') \right]$$

---

## Problems with TD Value Learning

- TD value learning 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:

$\pi(s) = \arg\max_a Q(s, a)$

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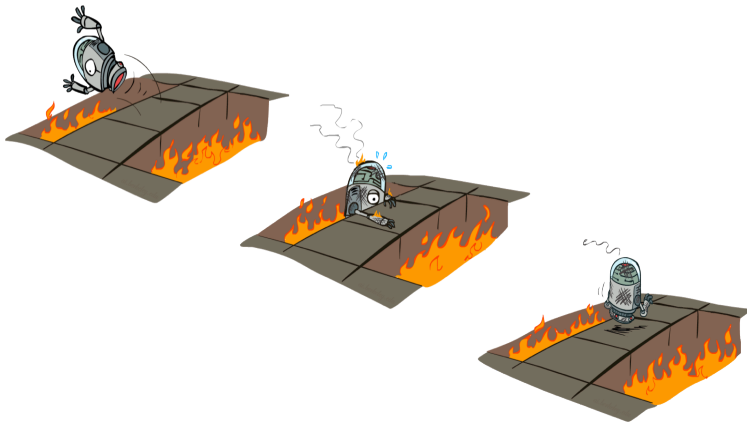
s

a

s, a

s,a,s'

s'

- Idea: learn Q-values, not values
- Makes action selection model-free too!

## Active Reinforcement Learning

13

---

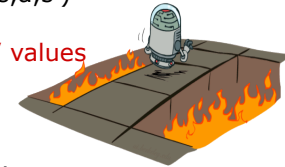## Active Reinforcement Learning

- Full reinforcement learning: optimal policies (like value iteration)
  - You don't know the transitions T(s,a,s')
  - You don't know the rewards R(s,a,s')
  - You choose the actions now
  - Goal: learn the optimal policy / values

- In this case:
  - Learner makes choices!
  - Fundamental tradeoff: exploration vs. exploitation
  - This is NOT offline planning! You actually take actions in the world and find out what happens…

14

---

## Detour: Q-Value Iteration

- Value iteration: find successive (depth-limited) values
  - Start with $V_0(s) = 0$, which we know is right
  - Given $V_k$, calculate the depth k+1 values for all states:

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

- But Q-values are more useful, so compute them instead
  - Start with $Q_0(s,a) = 0$, which we know is right
  - Given $Q_k$, calculate the depth k+1 q-values for all q-states:

$$Q_{k+1}(s,a) \leftarrow \sum_{s'} T(s,a,s') \left[ R(s,a,s') + \gamma \max_{a'} Q_k(s',a') \right]$$
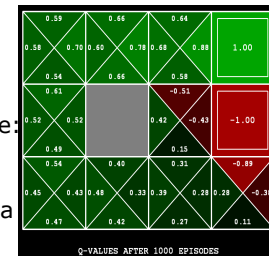
15

---

## Q-Learning

- Q-Learning: sample-based Q-value iteration
$$Q_{k+1}(s,a) \leftarrow \sum_{s'} T(s,a,s') \left[ R(s,a,s') + \gamma \max_{a'} Q_k(s',a') \right]$$



Q-VALUES AFTER 1000 EPISODES

- Learn Q(s,a) values as you go
  - Receive a sample (s,a,s',r)
  - Consider your old estimate $Q(s,a)$
  - Consider your new sample estimate:
    $$sample = R(s,a,s') + \gamma \max_{a'} Q(s',a')$$
  - Incorporate the new estimate into a running average:
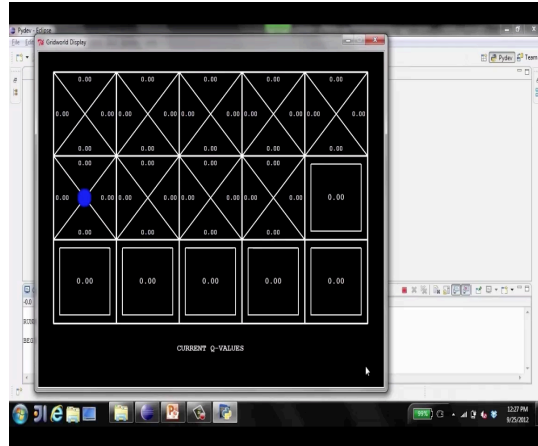    $$Q(s,a) \leftarrow (1-\alpha)Q(s,a) + (\alpha) \left[ sample \right]$$

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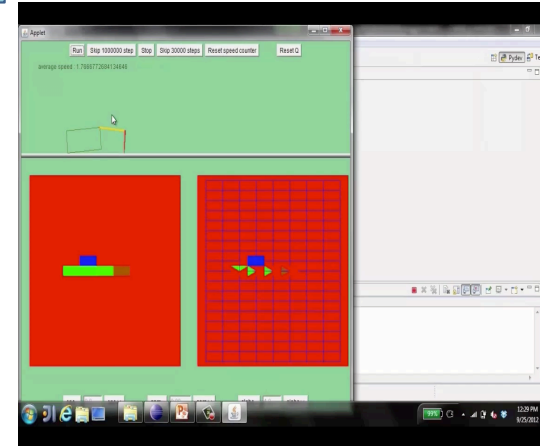[Demo: Q-learning – crawler (L10D3)]

16

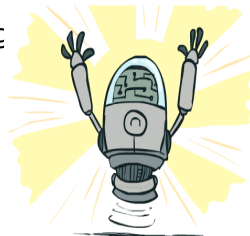## Video of Demo Q-Learning -- Gridworld

## Video of Demo Q-Learning -- Crawler

## 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
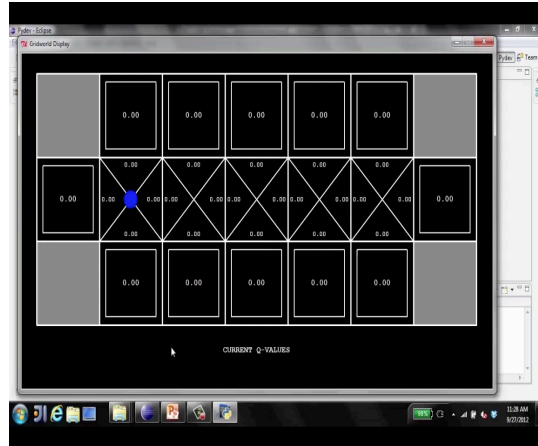  - Basically, in the limit, it doesn't matter how you select actions (!)

## How to Explore?

- Several schemes for forcing exploration
  - Simplest: random actions ($\varepsilon$-greedy)
    - Every time step, flip a coin
    - With (small) probability $\varepsilon$, act randomly
    - With (large) probability $1-\varepsilon$, 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 $\varepsilon$ over time
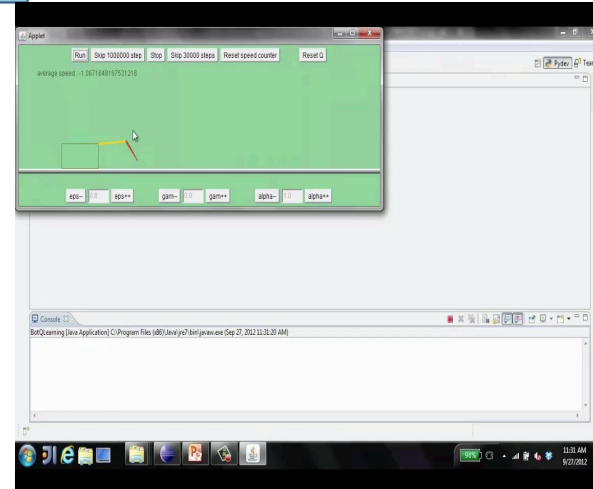    - Another solution: exploration functions

## Slide 21

Re Reinforcement Learning

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



21

## Slide 22

Re Reinforcement Learning

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



22

## Slide 23

Re Reinforcement Learning

# 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.
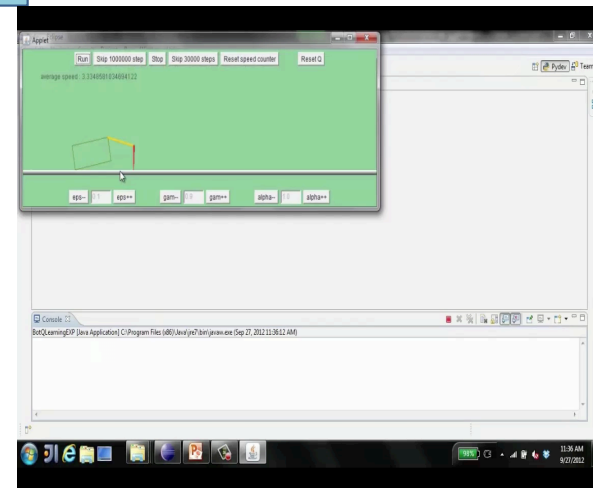
$$f(u,n) = u + k/n$$

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

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

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

23

## Slide 24

Re Reinforcement Learning

### Video of Demo Q-learning – Exploration Function – Crawler



24

# 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