# CIS 471/571(Fall 2020): Introduction to Artificial Intelligence

## Lecture 10: Learning (Part

Thanh H. Nguyen

Source: http://ai.berkeley.edu/home.html

# Reinforcement Learning

Agent

State: s
Reward: r

Actions: a

Environ

- Basic idea:
  - Receive feedback in the form of rewards
  - Agent's utility is defined by the reward function
  - Must (learn to) act so as to maximize expected rewards
  - All learning is based on observed samples of outcomes!

# Example: Learning to Walk



Initial

A Learning Trial

After Learning
[1K Trials]

[Kohl and Stone, ICRA 2004]

# Example: Learning to Walk

[Kohl and Stone, ICRA 2004]                    Initial

# Example: Learning to Walk

[Kohl and Stone, ICRA 2004]

Training

# Example: Learning to Walk

[Kohl and Stone, ICRA 2004]          Finished

# Reinforcement Learning

- Still assume a Markov decision process (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,

- Still looking for a policy $\pi(s)$

- New twist: don't know T or R
  - I.e. we don't know which states are good or what the actions do
  - Must actually try out actions and states to learn

# Offline (MDPs) vs. Online (RL)

Offline
Solution

Online
Learning

# Model–Based Learning

Assignment Project Exam Help

https://eduassistpro.github.io/

Add WeChat edu_assist_pro

# Model-Based Learning

- Model-Based Idea:
  - Learn an approximate model based on experiences
  - Solve for values as if the learned model were correct

- Step 1: Learn empirical
  - Count outcomes s' for each
  - Normalize to give an estimate of
  - Discover each $\widehat{R}(s, a, s')$ when we exp , s')
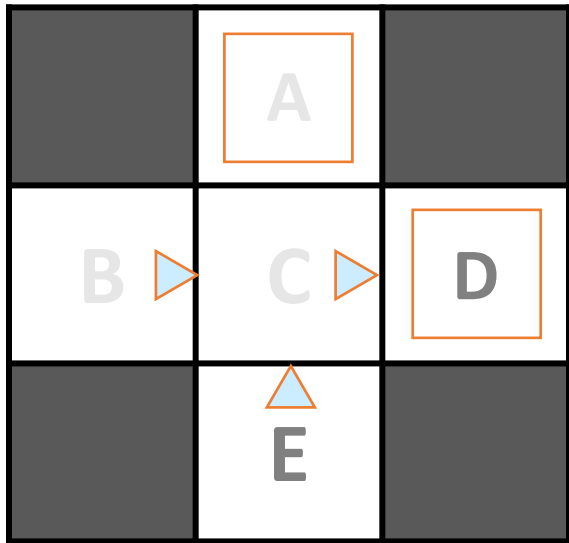
- Step 2: Solve the learned MDP
  - For example, use value iteration, as before

# Example: Model-Based Learning

**Input Policy**

**π**



*Assume:* $\gamma = 1$

**Observed Episodes (Training)**

Episode 1 <span style="color:red">Assignment Project Exam Help</span> Episode 2

B, e <span style="color:red">https://eduassistpro.github.io/</span>
C, e
D, exit, x, +10 C, -1
D, -1
<span style="color:red">Add WeChat edu_assist_pro</span> A, +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: Expected Age

Goal: Compute expected age of UO students

**Known P(A)**

Assignment Project Exam Help
$= 0.35 \times 20 + ...$

https://eduassistpro.github.io/

Without P(A), instead collec $[a_1, a_2, ... a_N]$

Add WeChat edu_assist_pro

**Unknown P(A): "Model Based"**

$$\hat{P}(a) = \frac{\text{num}(a)}{N}$$

$$E[A] \approx \sum_a \hat{P}(a) \cdot a$$

Why does this work? Because eventually you learn the right model.

**Unknown P(A): "Model Free"**

$$E[A] \approx \frac{1}{N} \sum_i a_i$$

Why does this work? Because samples appear with the right frequencies.

# Model-Free Learning

Assignment Project Exam Help

https://eduassistpro.github.io/

Add WeChat edu_assist_pro

# Passive Reinforcement Learning

Assignment Project Exam Help

https://eduassistpro.github.io/

Add WeChat edu_assist_pro

# Passive Reinforcement Learning

- Simplified task: policy evaluation
  - Input: a fixed policy $\pi(s)$
  - You don't know the transitions T(s,a,s')
  - You don't know the rewar
  - Goal: learn the state valu https://eduassistpro.github.io/

- In this case:
  - Learner is "along for the ride"
  - No choice about what actions to take
  - Just execute the policy and learn from experience
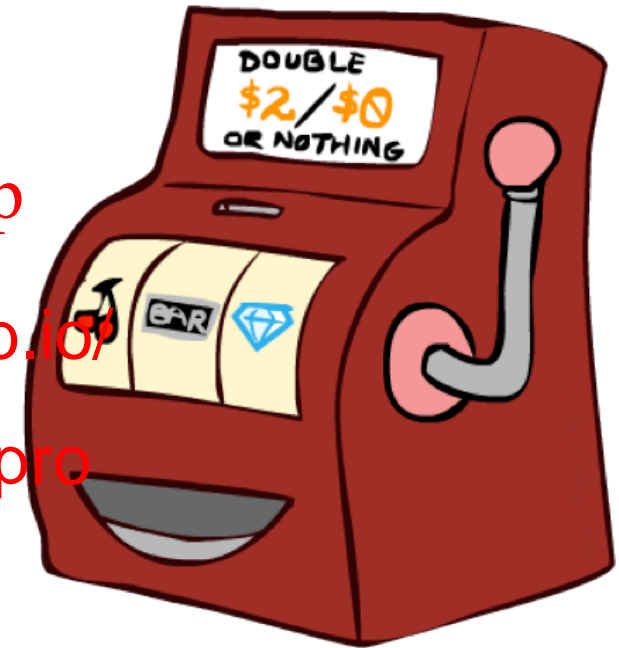  - This is NOT offline planning!  You actually take actions in the world.

# Direct Evaluation

- Goal: Compute values for each state under $\pi$

- Idea: Average together observed sample values
  - Act according to $\pi$
  - Every time you visit a state, write down sum of discounted rewards turned out
  - Average those samples

- This is called direct evaluation

# Example: Direct Evaluation

## Input Policy π



*Assume: γ = 1*

## Observed Episodes (Training)

### Episode 1

B, e
C, e
D, exit, x, +10

### Episode 2

C, -1
D, -1
+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

## Output Values

# 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 t
    values, using just sampl
- 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



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

# Why Not Use Policy Evaluation?

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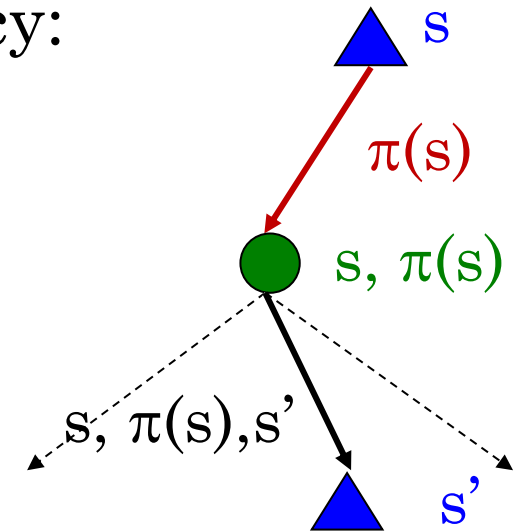  - Each round, replace V with a one-step-look-ahead layer over V

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

  - 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?

# Sample-Based Policy Evaluation?

- We want to improve our estimate of V by computing these averages:
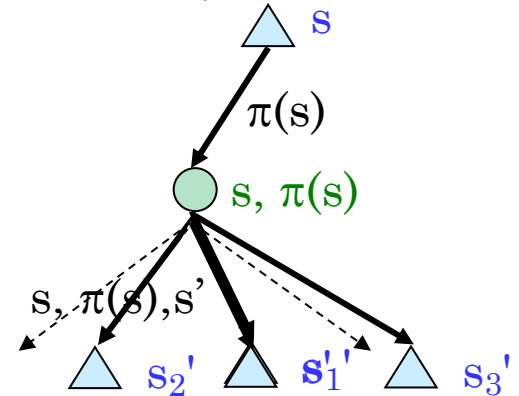
- Idea: Take samples of g the action!) and average

$$sample_2 = R(s, \pi(s), s'_2) + \gamma V_k^\pi(s'_2)$$

$$\dots$$

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

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

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

# Temporal Difference Learning

Assignment Project Exam Help

https://eduassistpro.github.io/

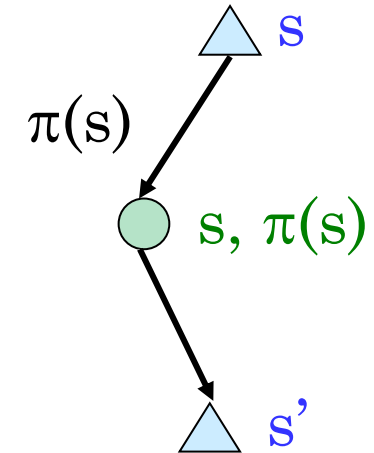Add WeChat edu_assist_pro

# 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 learn
  - Policy still fixed, still doing e
  - Move values toward value of whatever succes          nning average

π(s)

s

s, π(s)

s'

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

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

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

# 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:

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

- Decreasing learning rate (alpha) can give converging averages

# Example: Temporal Difference Learning

Observed Transitions

B, east, C, -2

C, east, D, -2

Assignment Project Exam Help

https://eduassistpro.github.io/

Add WeChat edu_assist_pro

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

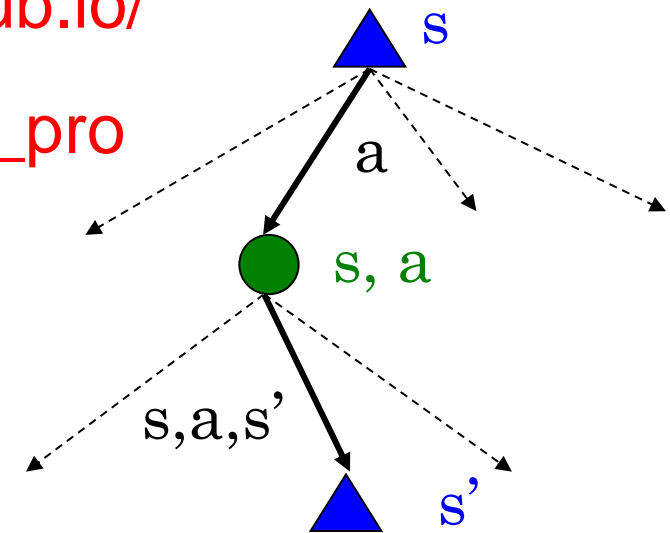$$V^{\pi}(s) \leftarrow (1-\alpha)V^{\pi}(s) + \alpha\left[R(s, \pi(s), s') + \gamma V^{\pi}(s')\right]$$

# Problems with TD Value Learning

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

- However, if we want policy, we're sunk:

s

a

s, a

s,a,s'

s'

- Idea: learn Q-values, not values

- Makes action selection model-free too!

# Active Reinforcement Learning

# 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 no
  - Goal: learn the optimal p

- 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…

# 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:

- But Q-values are more useful, so compu              stead
  - 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:

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

# Q-Learning

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

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

- Learn Q(s,a) values as
  - Receive a sample (s,a,s',
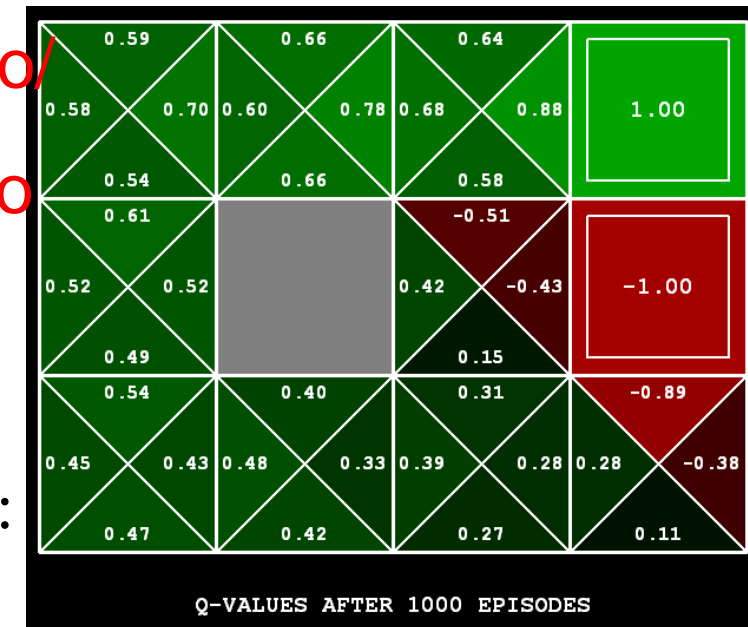  - Consider your old estimate: $Q(s,a)$
  - Consider your new sample estimate:

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

  - Incorporate the new estimate into a running average:

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



Q-VALUES AFTER 1000 EPISODES

# Q-Learning Properties

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

- This is called off-policy l

- 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 (!)