

# Overview: MDPs and Reinforcement Learning

## Known MDP: Offline Solution

Goal

Compute  $V^*$ ,  $Q^*$ ,  $\pi^*$

Evaluate a fixed policy  $\pi$

Technique

Value / policy iteration

Policy evaluation

## Unknown MDP: Model-Based

Goal

Compute  $V^*$ ,  $Q^*$ ,  $\pi^*$

Eval fixed policy  $\pi$

Technique

VI/PI on approx. MDP

PE on approx. MDP

## Unknown MDP: Model-Free

Goal

Compute  $V^*$ ,  $Q^*$ ,  $\pi^*$

Eval fixed policy  $\pi$

Technique

Q-learning

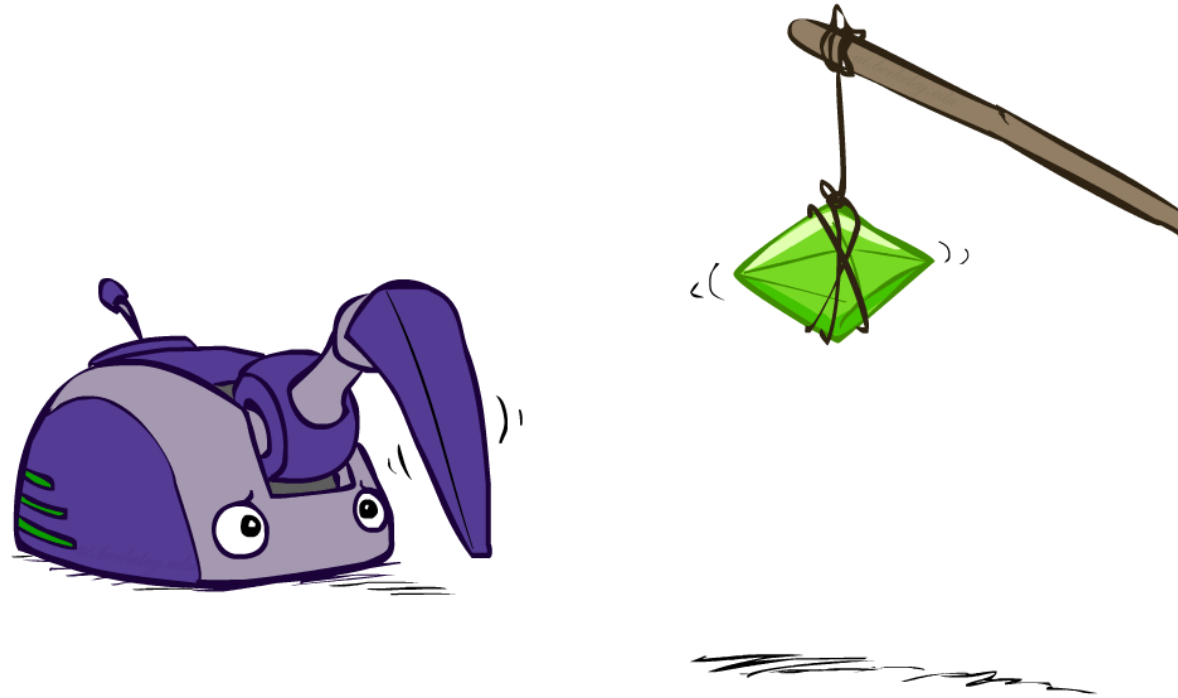
TD/Value Learning

# Wrap up MDPs

Switch to MDP II slides

# AI: Representation and Problem Solving

## Reinforcement Learning



Instructor: Pat Virtue

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

# MDP Notation

Standard expectimax:  $V(s) = \max_a \sum_{s'} P(s'|s, a)V(s')$

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

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

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

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

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

# Poll (unused)

Which of the following are used in policy iteration?

A. Value iteration: 
$$V_{k+1}(s) = \max_a \sum_{s'} P(s'|s, a) [R(s, a, s') + \gamma V_k(s')], \quad \forall s$$

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

C. Policy extraction: 
$$\pi_V(s) = \operatorname{argmax}_a \sum_{s'} P(s'|s, a) [R(s, a, s') + \gamma V(s')], \quad \forall s$$

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

# Poll (unused)

Which of the following are used in policy iteration?

A. Value iteration: 
$$V_{k+1}(s) = \max_a \sum_{s'} P(s'|s, a) [R(s, a, s') + \gamma V_k(s')], \quad \forall s$$

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

C. Policy extraction: 
$$\pi_V(s) = \operatorname{argmax}_a \sum_{s'} P(s'|s, a) [R(s, a, s') + \gamma V(s')], \quad \forall s$$

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

# Poll (unused)

Rewards may depend on any combination of *state*, *action*, *next state*.

Which of the following are valid formulations of the Bellman equations?

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

B.  $V^*(s) = R(s) + \gamma \max_a \sum_{s'} P(s'|s, a) V^*(s')$

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

D.  $Q^*(s, a) = R(s, a) + \gamma \sum_{s'} P(s'|s, a) \max_{a'} Q^*(s', a')$

# Poll (unused)

Rewards may depend on any combination of *state, action, next state*.

Which of the following are valid formulations of the Bellman equations?

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

✓ B.  $V^*(s) = R(s) + \gamma \max_a \sum_{s'} P(s'|s, a) V^*(s')$

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

✓ D.  $Q^*(s, a) = R(s, a) + \gamma \sum_{s'} P(s'|s, a) \max_{a'} Q^*(s', a')$



# Reinforcement learning

What if we didn't know  $P(s'|s, a)$  and  $R(s, a, s')$ ?

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

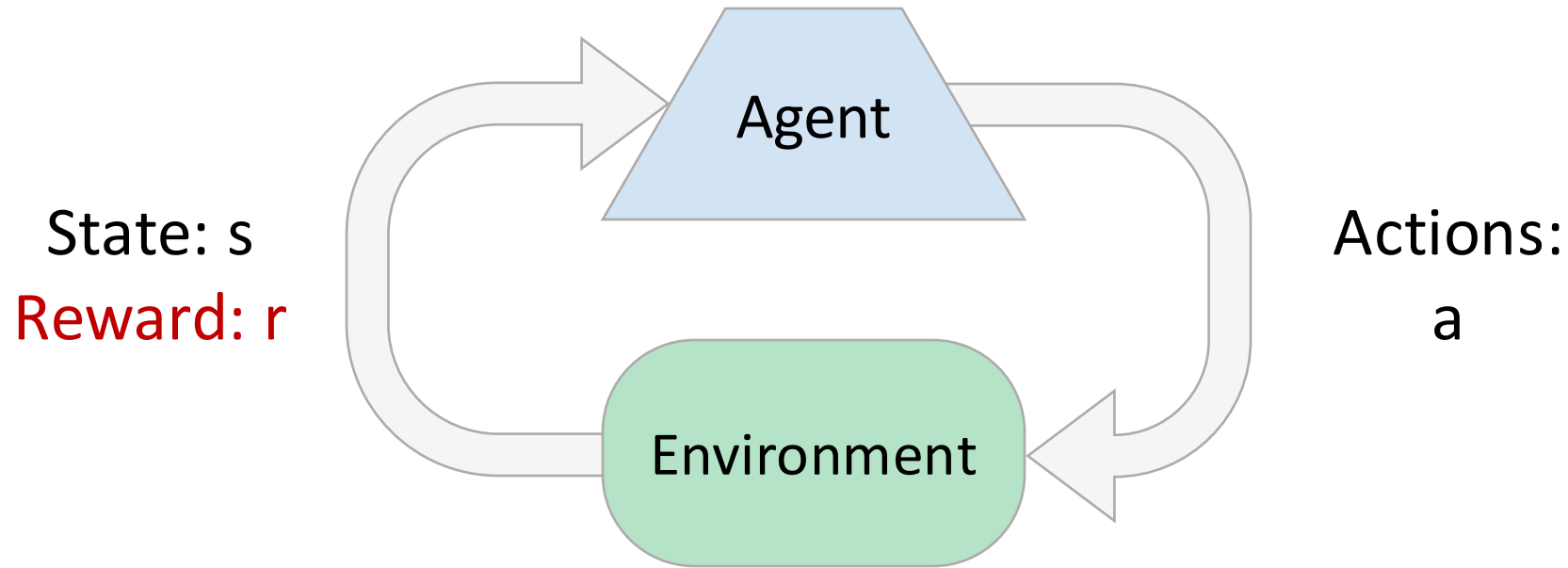
Q-iteration: 
$$Q_{k+1}(s, a) = \sum_{s'} P(s'|s, a) [R(s, a, s') + \gamma \max_{a'} Q_k(s', a')], \quad \forall s, a$$

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

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

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

# Reinforcement Learning



## Basic idea:

- All learning is based on observed **samples** of rewards and next states!
- Receive feedback in the form of **rewards**
- Must (learn to) act so as to **maximize expected rewards**

# Example: Learning to Walk



Initial

# Example: Learning to Walk

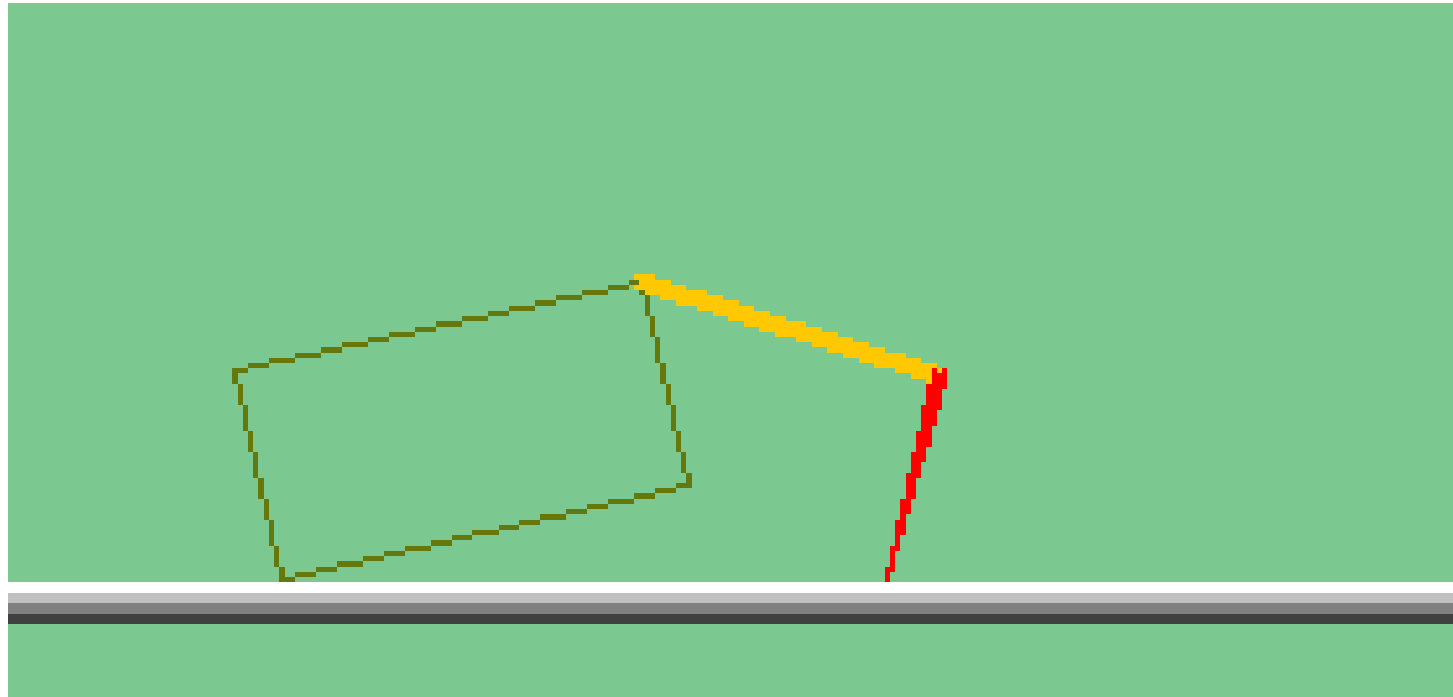


Finished

# Example: Toddler Robot



# The Crawler!



# Demo Crawler Bot

# 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  $R(s,a,s')$

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

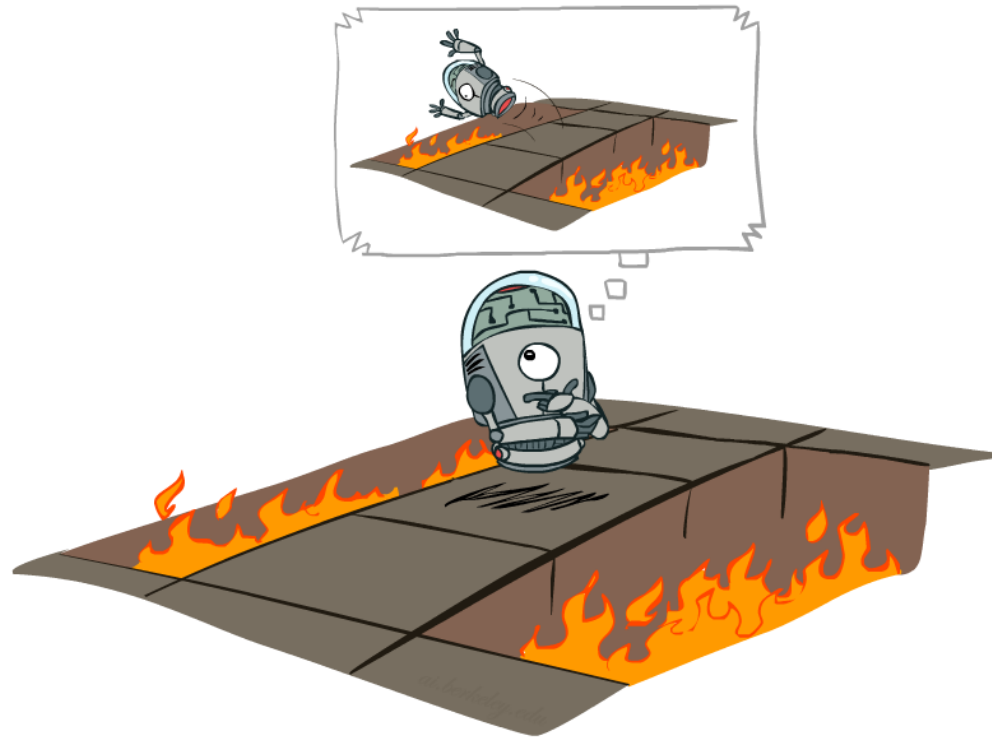


New twist: **don't know  $T$  or  $R$**

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



# Offline (MDPs) vs. Online (RL)



Offline Solution  
(Known MDP)



Online Learning  
(Unknown MDP)

# 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

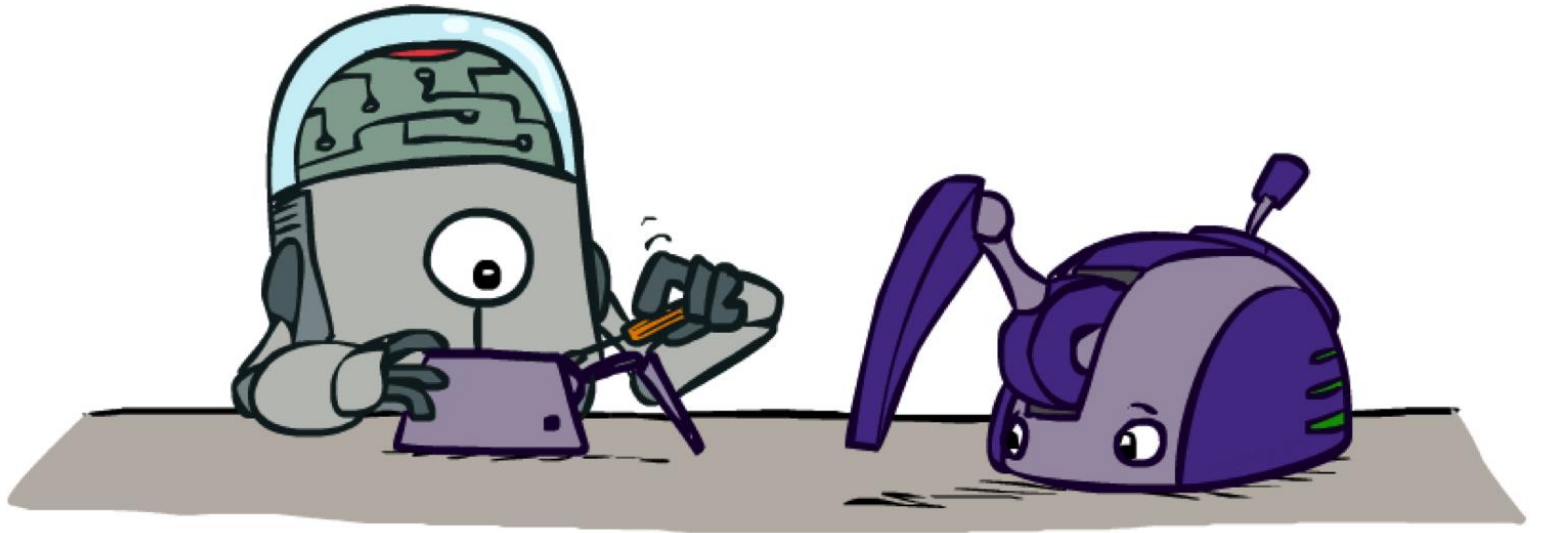
Model-Free

Passive Reinforcement Learning

- Direct Evaluation (simple)
- TD Learning

Active Reinforcement Learning

- Q-Learning



Online Learning

Model-based Learning

# 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 MDP model

- Count outcomes  $s'$  for each  $s, a$
- Normalize to give an estimate of  $\hat{T}(s, a, s')$
- Discover each  $\hat{R}(s, a, s')$  when we experience  $(s, a, s')$

## Step 2: Solve the learned MDP

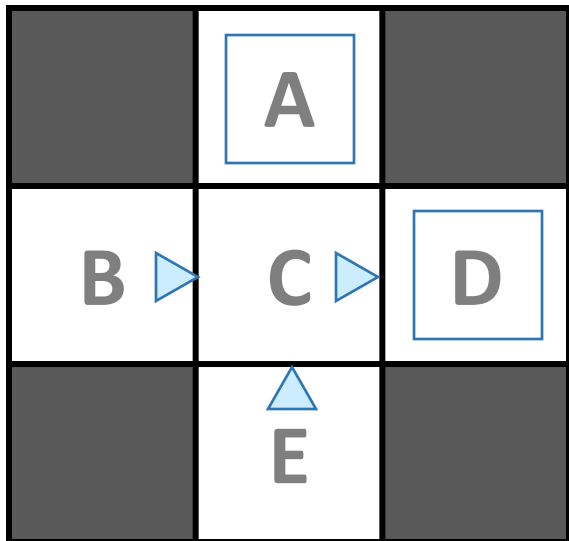
- For example, use value iteration, as before



# Example: Model-Based Learning

Episode: a sequence of states actions and rewards sampled from the environment

## Input Policy $\pi$



Assume:  $\gamma = 1$

## 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

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

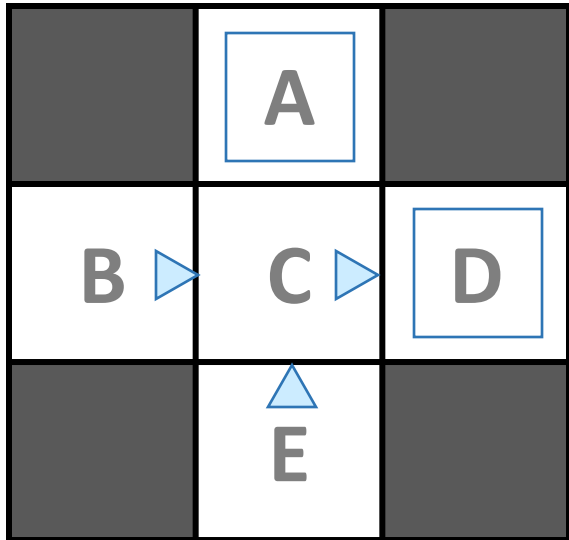
T(B, east, C) =  
T(C, east, D) =  
T(C, east, A) =  
...

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

R(B, east, C) =  
R(C, east, D) =  
R(D, exit, x) =  
...

# Example: Model-Based Learning

## Input Policy $\pi$



Assume:  $\gamma = 1$

## 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: Expected Age

Goal: Compute expected age of 15-281 students

Known  $P(A)$

$$E[A] = \sum_a P(a) \cdot a = 0.35 \times 20 + \dots$$

Without  $P(A)$ , instead collect samples  $[a_1, a_2, \dots, a_N]$

Unknown  $P(A)$ : “Model Based”

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

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

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.

# 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

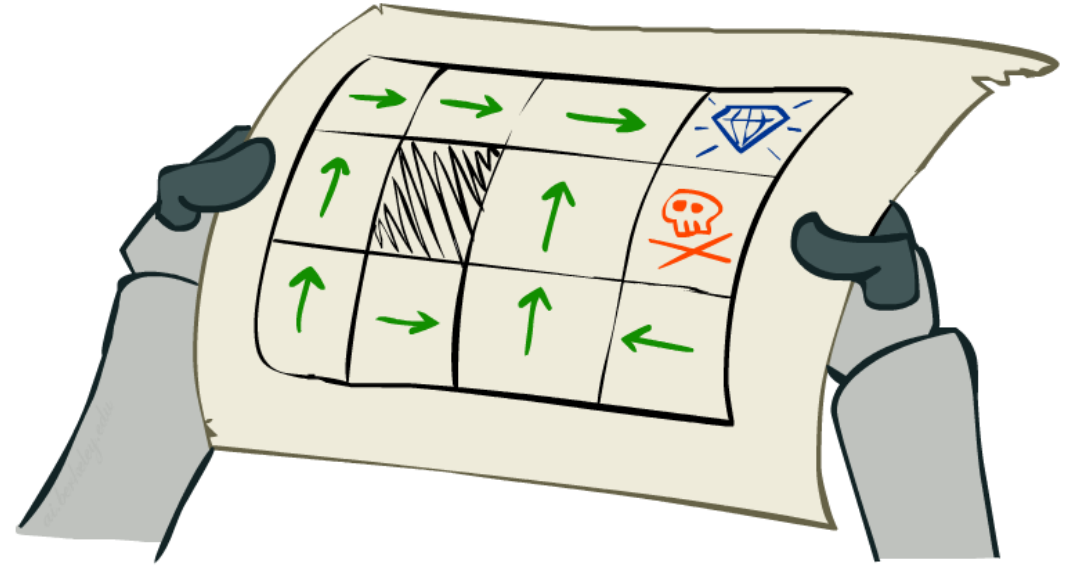
Passive Reinforcement Learning

- Direct Evaluation (simple)
- TD Learning

Active Reinforcement Learning

- Q-Learning





Online Learning

Model-free Learning

Passive Reinforcement Learning

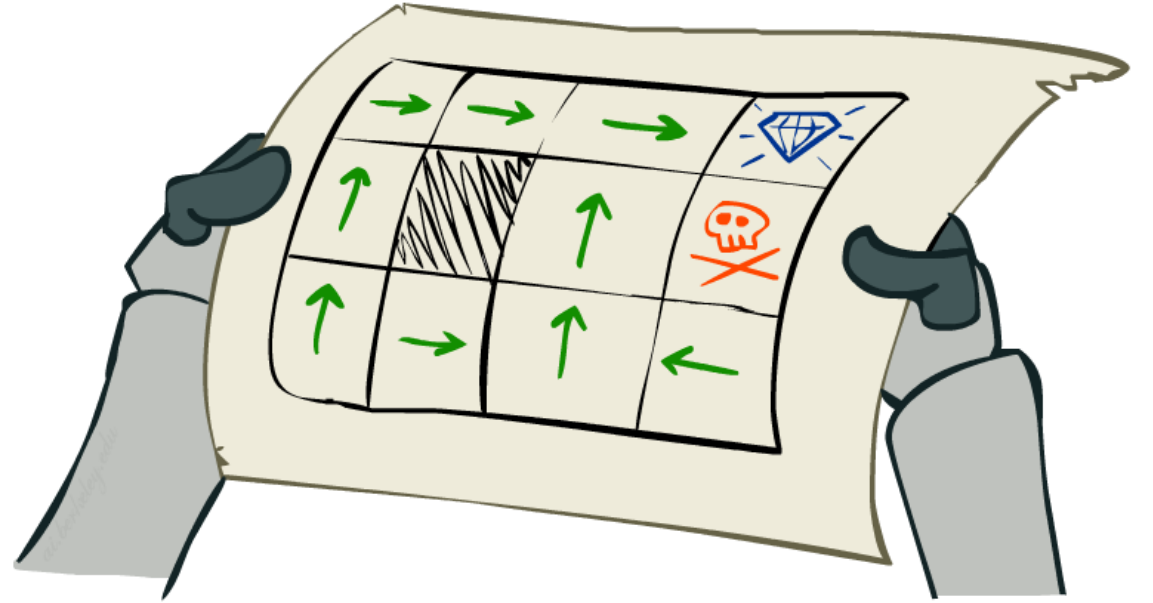
# Passive Reinforcement Learning

## 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 rewards  $R(s,a,s')$
- **Goal: learn the state values**

## 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.



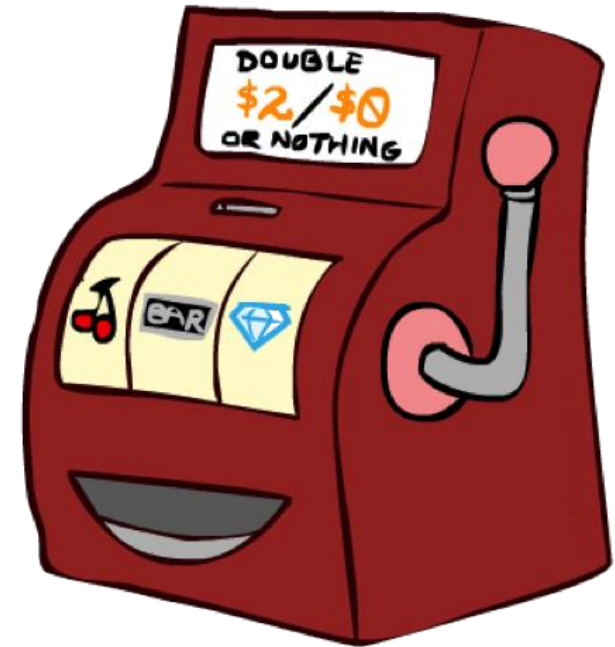
# Simple Passive Learning: Direct Evaluation

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 what the sum of discounted rewards turned out to be
- Average those samples

This is called direct evaluation



# Simple Passive Learning: Direct Evaluation

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

Idea: Average together observed sample values

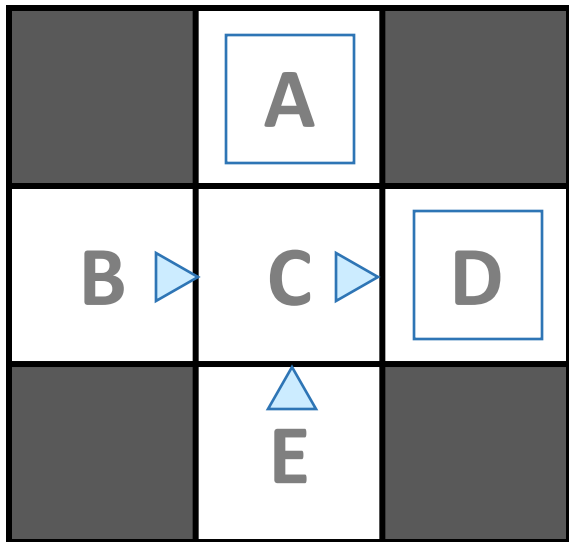
- Act according to  $\pi$
- Every time you visit a state, write down what the sum of discounted rewards turned out to be
- Average those samples

This is called direct evaluation

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

# Example: Direct Evaluation

Input Policy  $\pi$



Assume:  $\gamma = 1$

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

Output Values

	-10 A	
+8 B	+4 C	+10 D
	-2 E	

# 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	

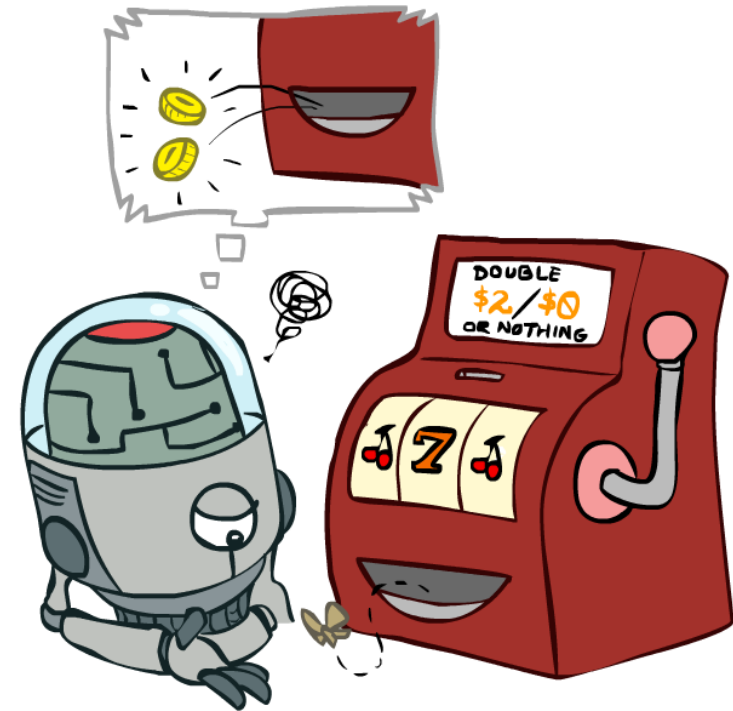
*If B and E both go to C under this policy, how can their values be different?*

Online Learning

Model-free Learning

Passive Reinforcement Learning

Temporal Difference Learning



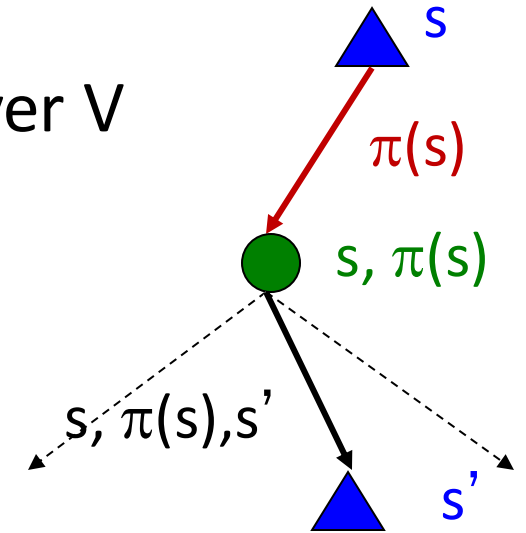
# Why Not Use Policy Evaluation?

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?



# 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

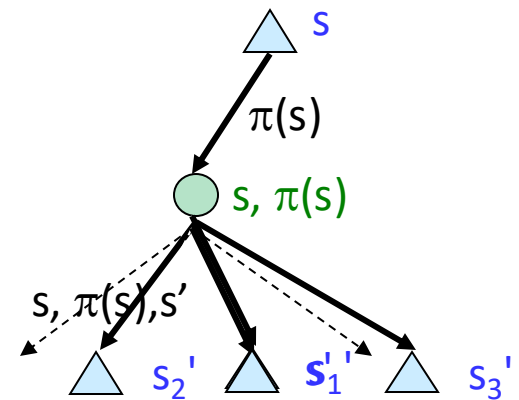
$$sample_1 = R(s, \pi(s), s'_1) + \gamma V_k^{\pi}(s'_1)$$

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

...

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

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



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

# Temporal Difference Learning

Big idea: learn from every experience!

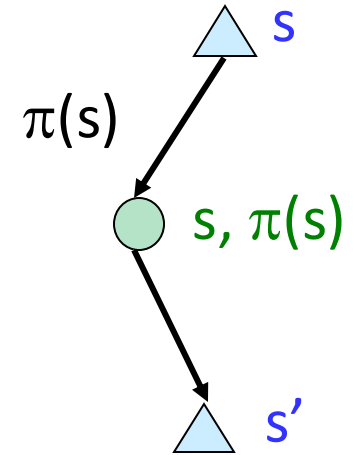
- Update  $V(s)$  each time we experience a transition  $(s, a, s', r)$
- Likely outcomes  $s'$  will contribute updates more often

Temporal difference learning of values

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

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

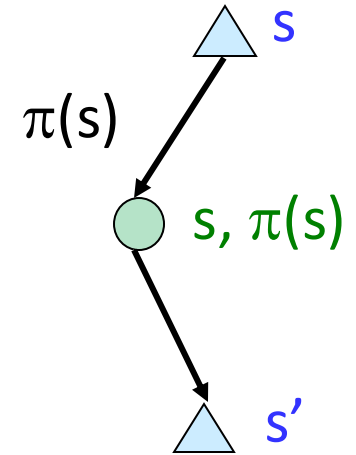
Update to  $V(s)$ :



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

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

	A	
B	C	D
	E	

Observed Transitions

B, east, C, -2

	0	
0	0	8
	0	

C, east, D, -2

	0	
-1	0	8
	0	

	0	
-1	3	8
	0	

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

# Example: Temporal Difference Learning

States

	A	
B	C	D
	E	

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

Observed Transitions

B, east, C, -2

	0	
0	0	8
	0	

C, east, D, -2

	0	
-1	0	8
	0	

	0	
-1	3	8
	0	

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

# 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:

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

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

Idea: learn Q-values, not values

Makes action selection model-free too!

