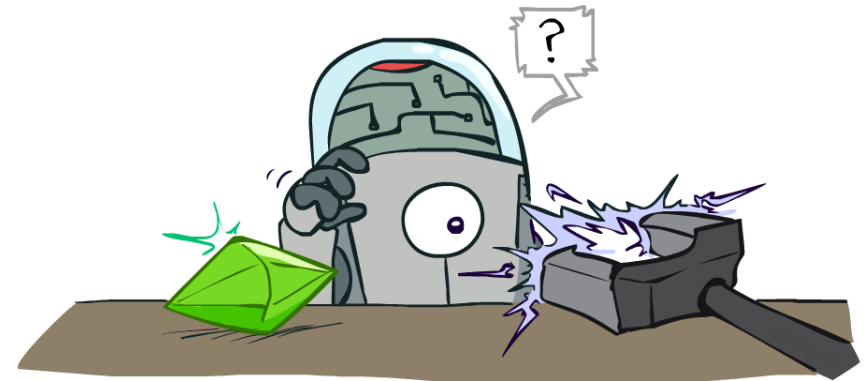

CSE 573:

Artificial Intelligence

Hanna Hajishirzi

Reinforcement Learning

slides adapted from
Dan Klein, Pieter Abbeel ai.berkeley.edu
And Dan Weld, Luke Zettlemoyer



MDPs Recap

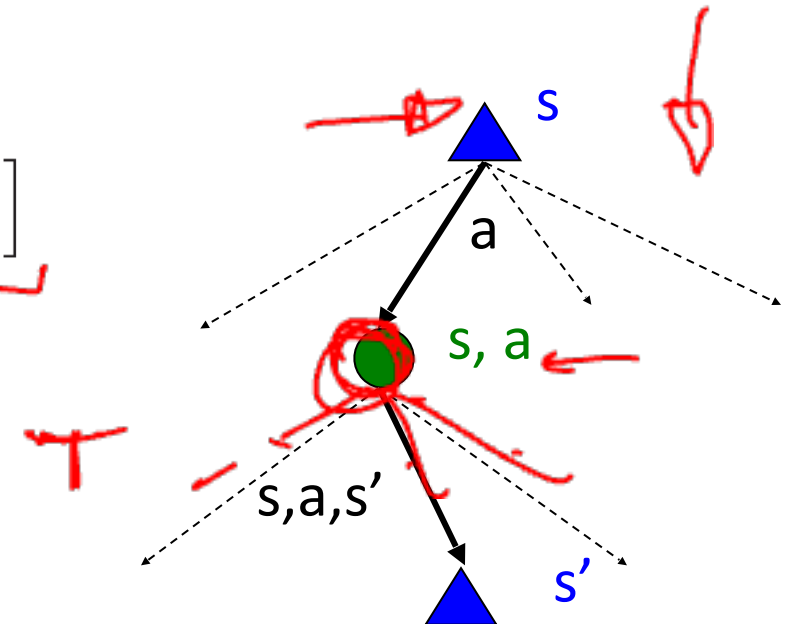
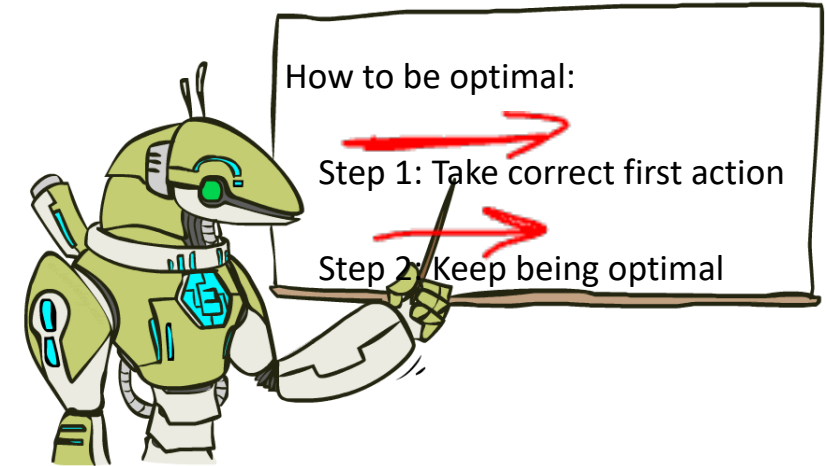
The Bellman Equations

- Definition of “optimal utility” via expectimax recurrence gives one-step lookahead relationship amongst optimal utility values

$$V^*(s) = \max_a Q^*(s, a)$$

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

- The $V^*(s) = \max_a \sum T(s, a, s') [R(s, a, s') + \gamma V^*(s')]$ value in a way, we use s' over and over.



Value Iteration

- Bellman equations **characterize** the optimal values:

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

(A red arrow points to the left side of the equation.)

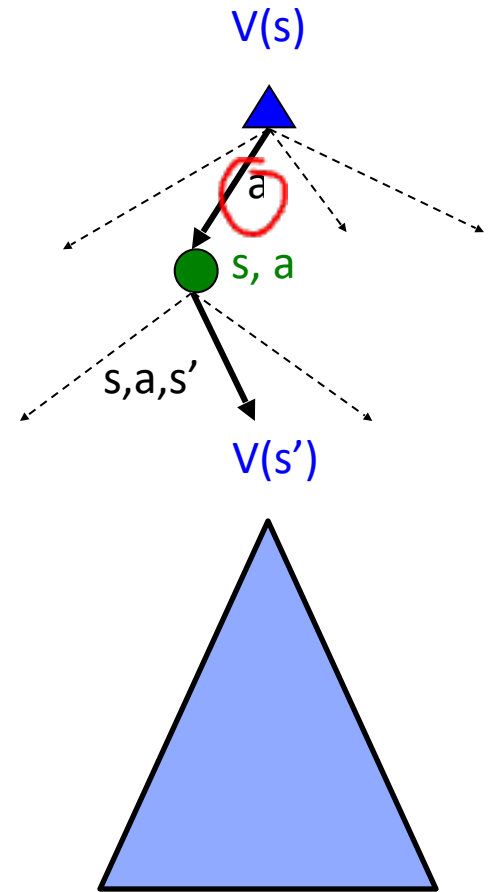
- Value iteration **computes** them:

$V_1 = 0$

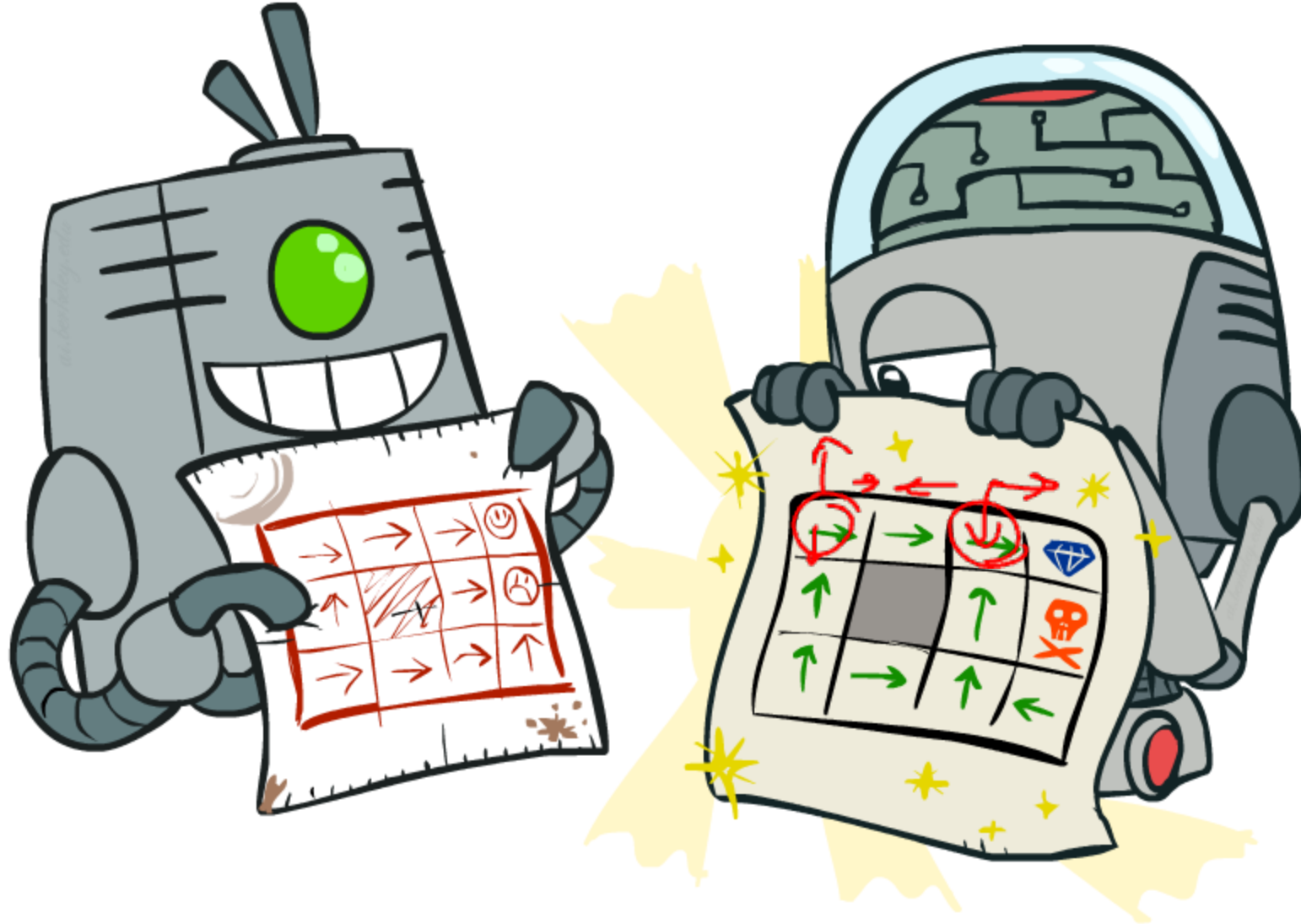
$$V_{k+1}(s) \leftarrow \max_a \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V_k(s')]$$

(Handwritten red annotations: a circle around $V_{k+1}(s)$, a box around γ , and the text $IS^2(A)$ below the equation.)

- Value iteration is just a fixed point solution method
 - ... though the V_k vectors are also interpretable as time-limited values



Policy Methods



Policy Evaluation

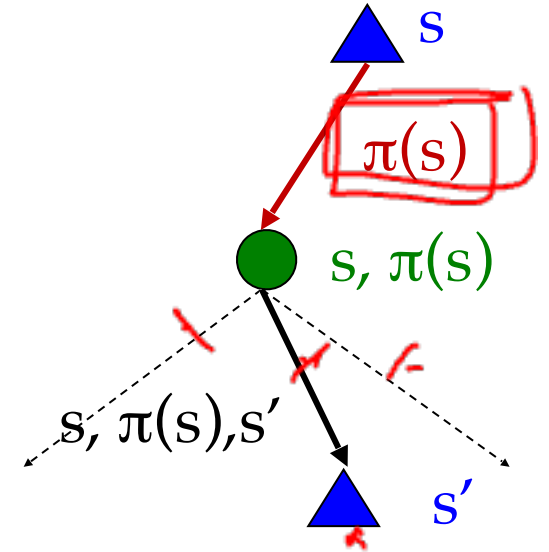
- Another basic operation: compute the utility of a state s under a fixed (generally non-optimal) policy

- Define the utility of a state s , under a fixed policy π :

$V^\pi(s)$ = expected total discounted rewards starting in s and following π

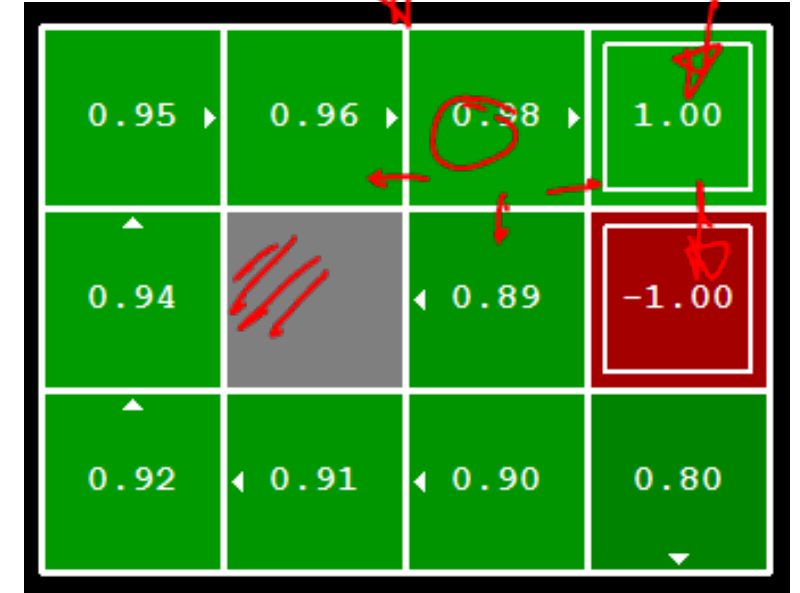
- Recursive relation (one-step look-ahead / Bellman equation):

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



Policy Extraction

- Let's imagine we have the optimal values $V^*(s)$
- How should we act?
 - It's not obvious!
- We need to do a mini-expectimax (one step)



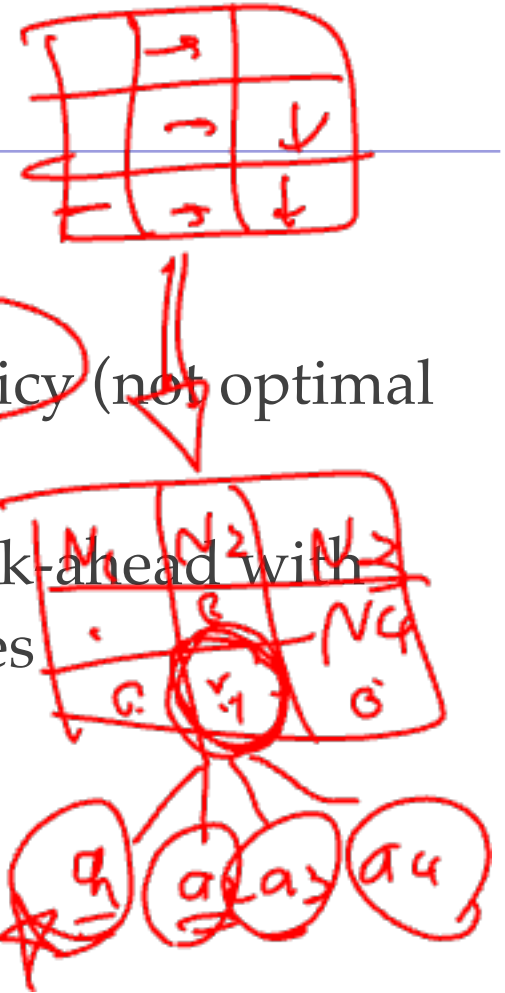
$$\pi^*(s) = \arg \max_a \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V^*(s')]$$

Red handwritten annotations: a circle around a with an arrow pointing to the $\arg \max$ operator; a long arrow from $V^*(s')$ back to $V^*(s)$ in the text below; and a wavy line above $V^*(s')$.

- This is called **policy extraction**, since it gets the policy implied by the values

Policy Iteration *

- Alternative approach for optimal values:
 - **Step 1: Policy evaluation:** calculate utilities for some fixed policy (not optimal utilities!) until convergence
 - **Step 2: Policy improvement:** update policy using one-step look-ahead with resulting converged (but not optimal!) utilities as future values
 - Repeat steps until policy converges
- This is **policy iteration**
 - It's still optimal!
 - Can converge (much) faster under some conditions



Policy Iteration

- Evaluation: For fixed current policy π , find values with policy evaluation:

- Iterate until values converge:

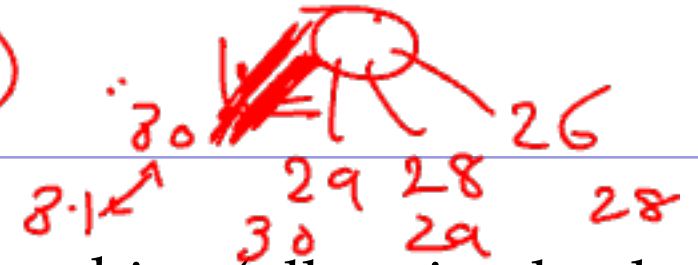
$$\rightarrow V_{k+1}^{\pi_i}(s) \leftarrow \sum_{s'} T(s, \pi_i(s), s') [R(s, \pi_i(s), s') + \gamma V_k^{\pi_i}(s')]$$

- Improvement: For fixed values, get a better policy using policy extraction

- One-step look-ahead:

$$\pi_{i+1}(s) = \arg \max_a \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V^{\pi_i}(s')]$$

Comparison



- Both value iteration and policy iteration compute the same thing (all optimal values)

- In value iteration:

- Every iteration updates both the values and (implicitly) the policy
- We don't track the policy, but taking the max over actions implicitly recomputes it

- ~~In policy iteration:~~

- We do several passes that update utilities with fixed policy (each pass is fast because we consider only one action, not all of them)
- After the policy is evaluated, a new policy is chosen (slow like a value iteration pass)
- The new policy will be better (or we're done)

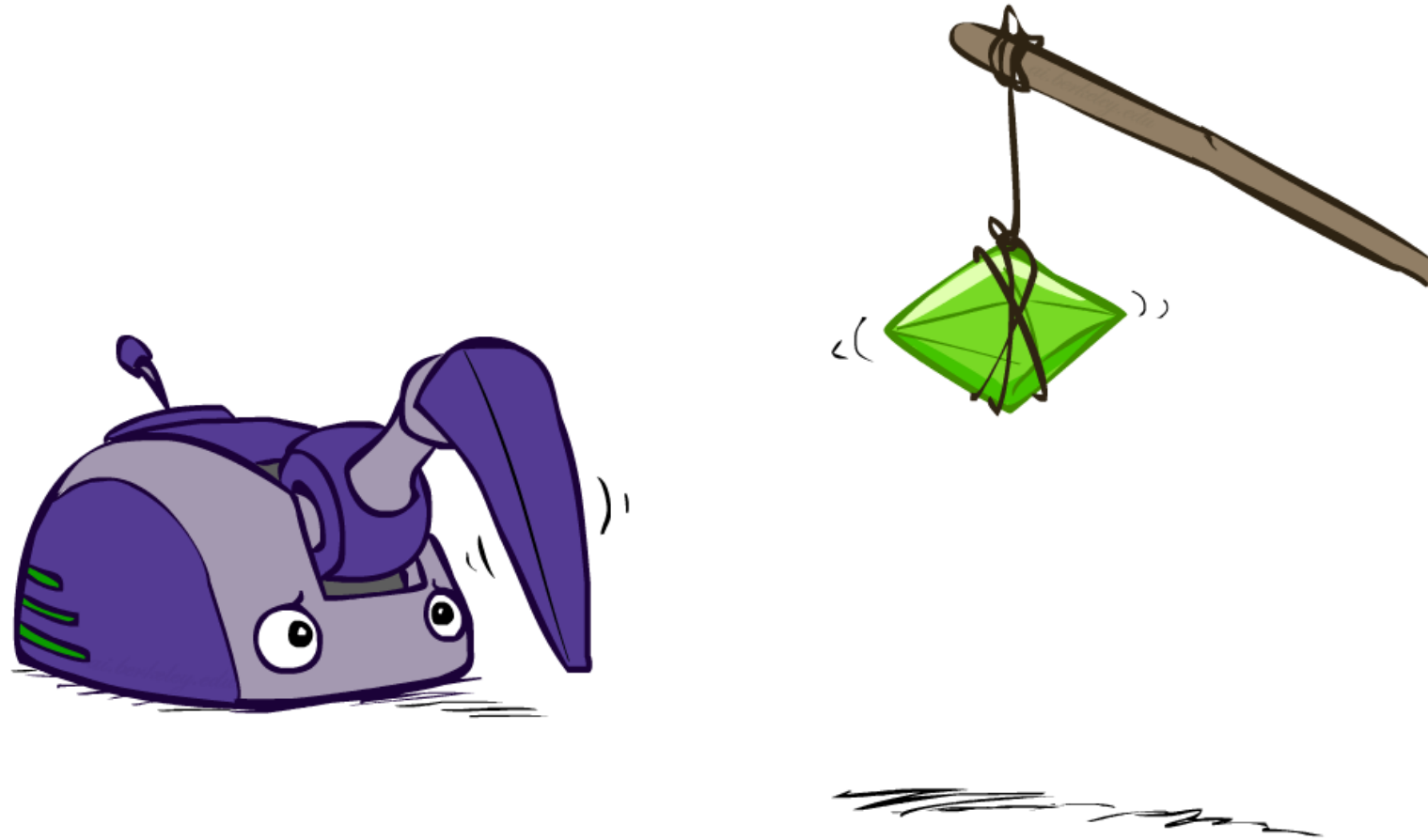
- Both are dynamic programs for solving MDPs

Handwritten equation showing the Bellman optimality equation: $V_k = f(V_{k-1}) = f(V_{k-2})$. Below it, $V_0 = 0$ is written.

Summary: MDP Algorithms

- So you want to....
 - Compute optimal values: use value iteration or policy iteration
 - Compute values for a particular policy: use policy evaluation
 - Turn your values into a policy: use policy extraction (one-step lookahead)
- These all look the same!
 - They basically are – they are all variations of Bellman updates
 - They all use one-step lookahead expectimax fragments
 - They differ only in whether we plug in a fixed policy or max over actions

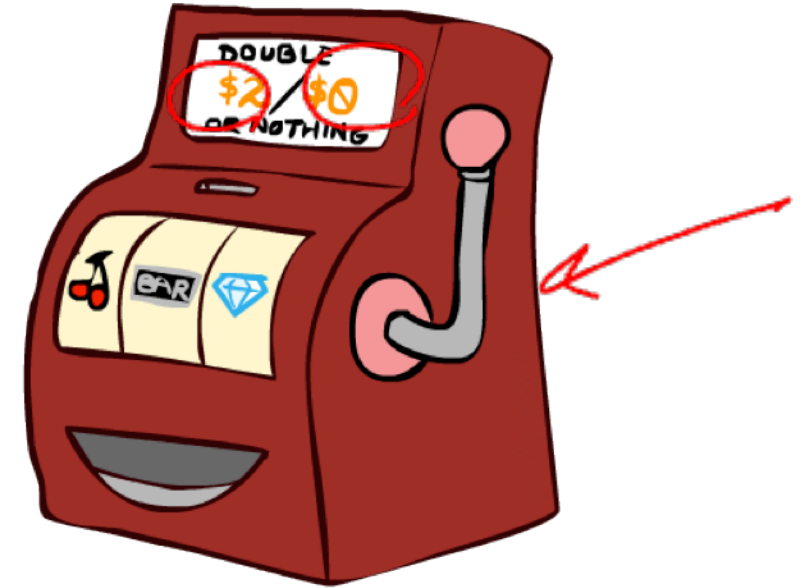
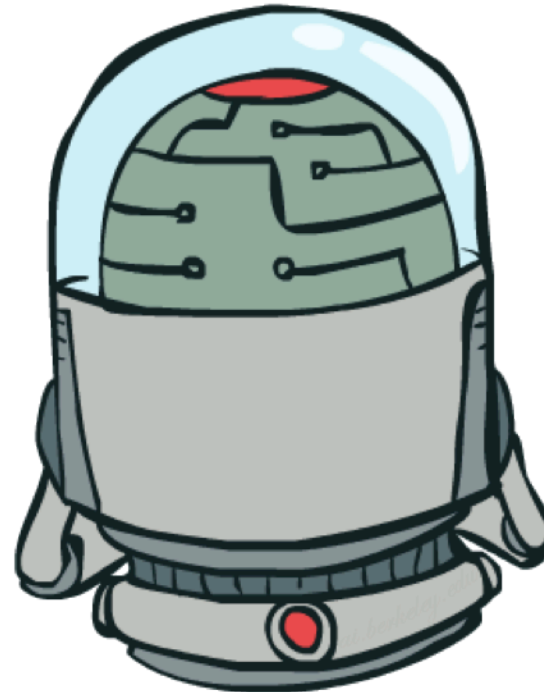
Reinforcement Learning



MDP
offline planning

Double Bandits

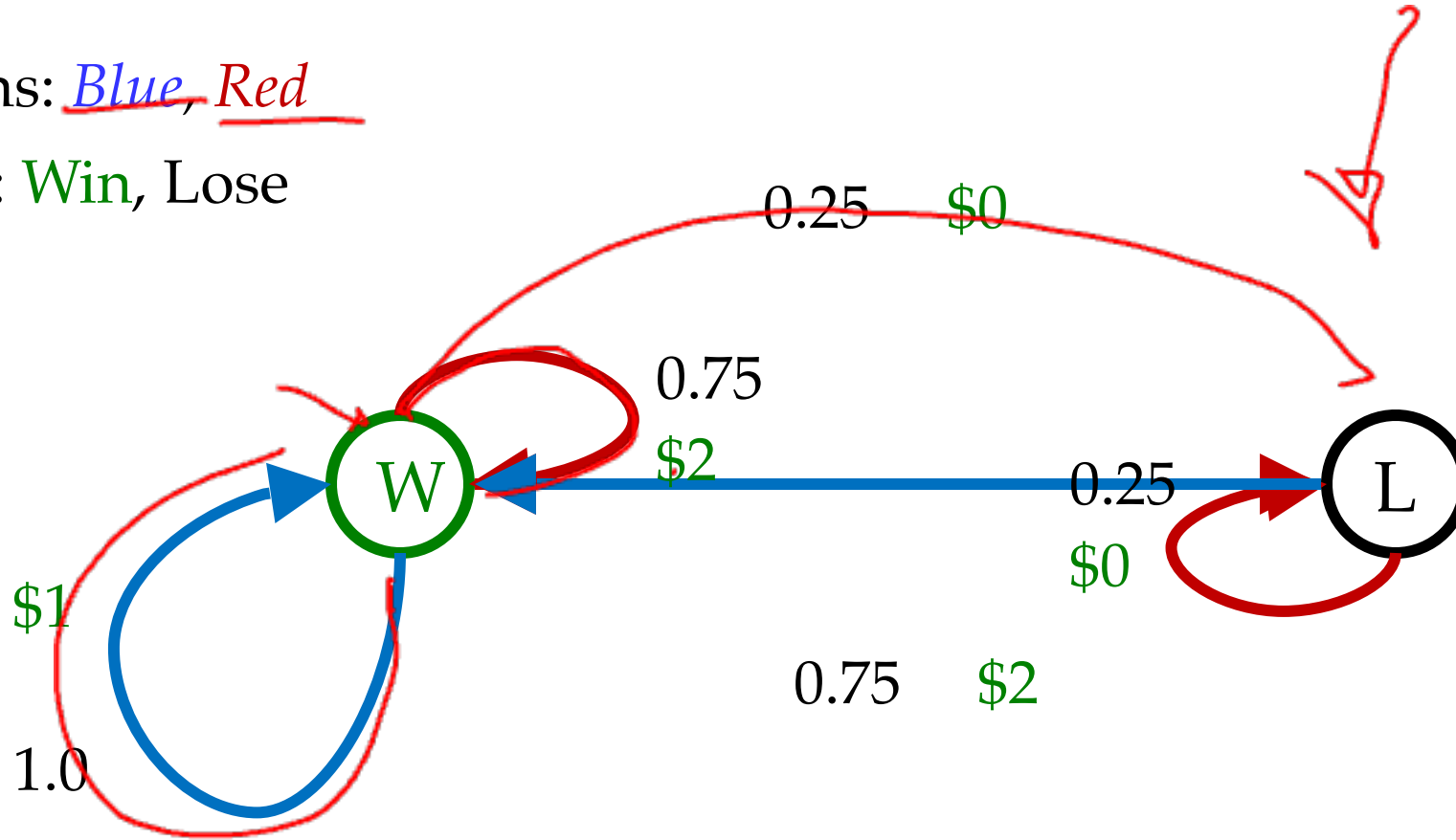
RL
online learning



0.95
 $\frac{3}{4}$ $\frac{1}{4}$

Double-Bandit MDP

- Actions: Blue, Red
- States: Win, Lose



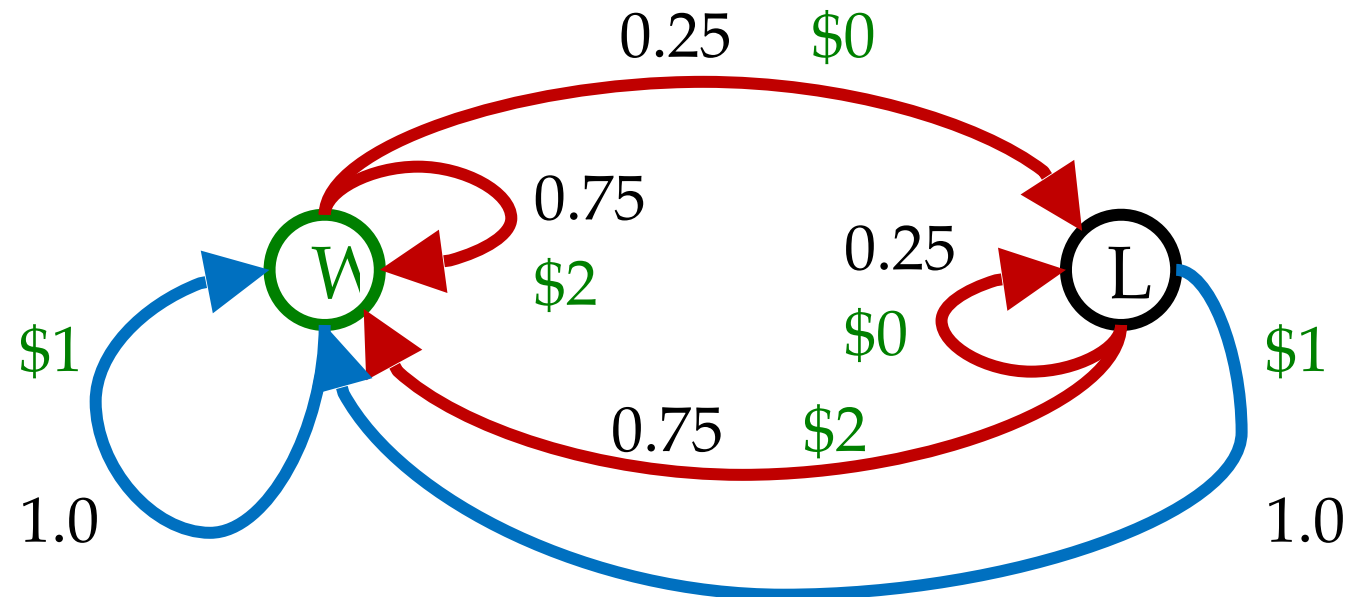
No discount
10 time steps
Both states have
the same value

$$10 \times \left(\frac{3}{4} \times 2 + \frac{1}{4} \times 0 \right) = \$15$$

Offline Planning

- Solving MDPs is offline planning
 - You determine all quantities through computation
 - You need to know the details of the MDP
 - You do not actually play the game!

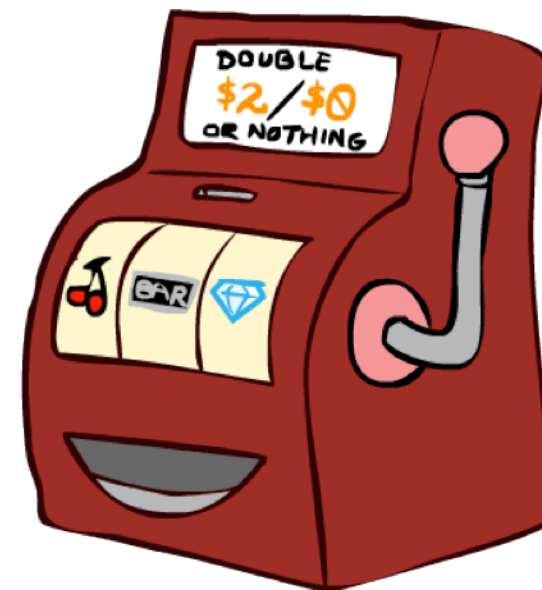
No discount
10 time steps



Let's Play!



\$10



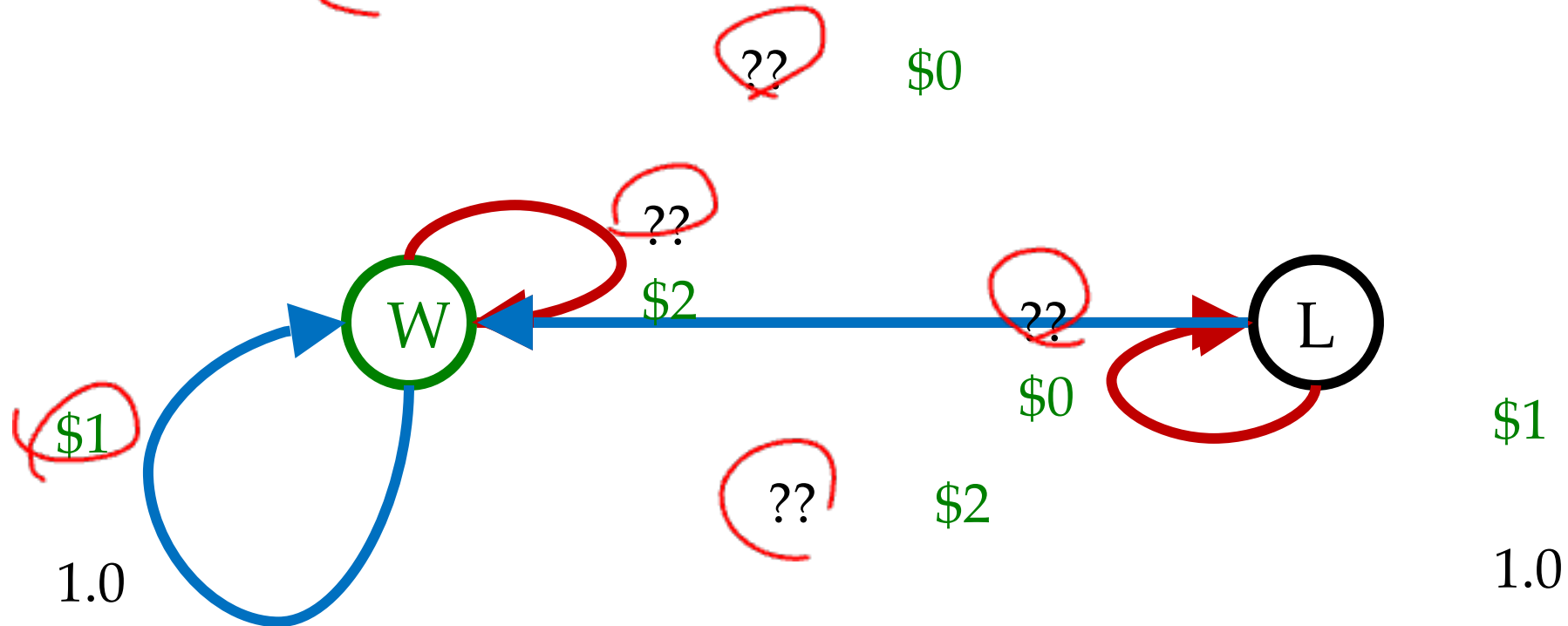
10

\$2	\$2	\$0	\$2	\$2
<u>\$2</u>	<u>\$2</u>	\$0	<u>\$0</u>	<u>\$0</u>
<u>—</u>	<u>—</u>			

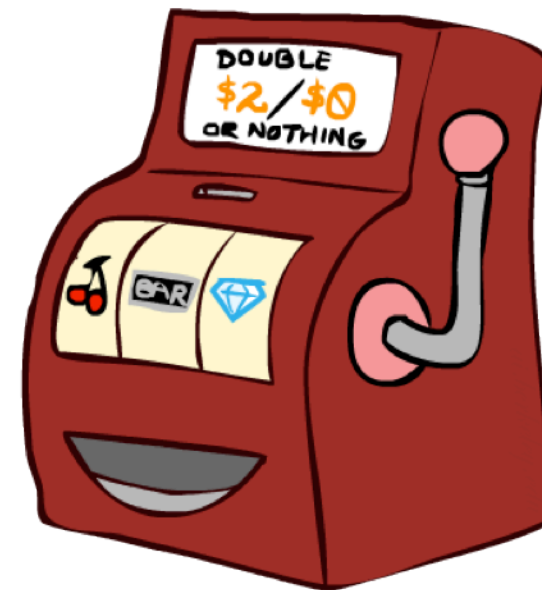
\$12

Online Planning

- Rules changed! Red's win chance is different.



Let's Play!



\$0 \$0 \$2 \$0
\$0 \$2 \$2 \$0 \$0
\$0

What Just Happened?

- That wasn't planning, it was learning!
 - Specifically, reinforcement learning
 - There was an MDP, but you couldn't solve it with just computation
 - You needed to actually act to figure it out
- Important ideas in reinforcement learning that came up
 - Exploration: you have to try unknown actions to get information
 - Exploitation: eventually, you have to use what you know
 - Regret: even if you learn intelligently, you make mistakes
 - Sampling: because of chance, you have to try things repeatedly
 - Difficulty: learning can be much harder than solving a known MDP

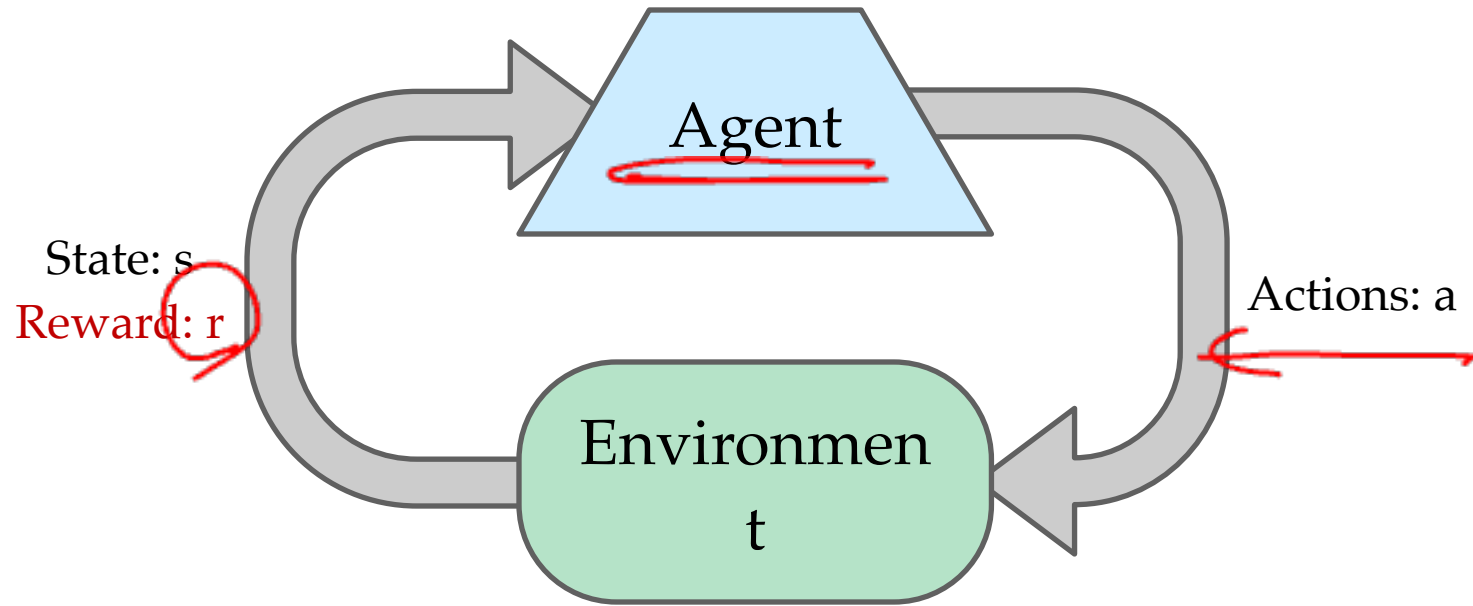


Reinforcement Learning

- 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



Reinforcement Learning



- 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]

Example: Learning to Walk



Initial

Example: Learning to Walk



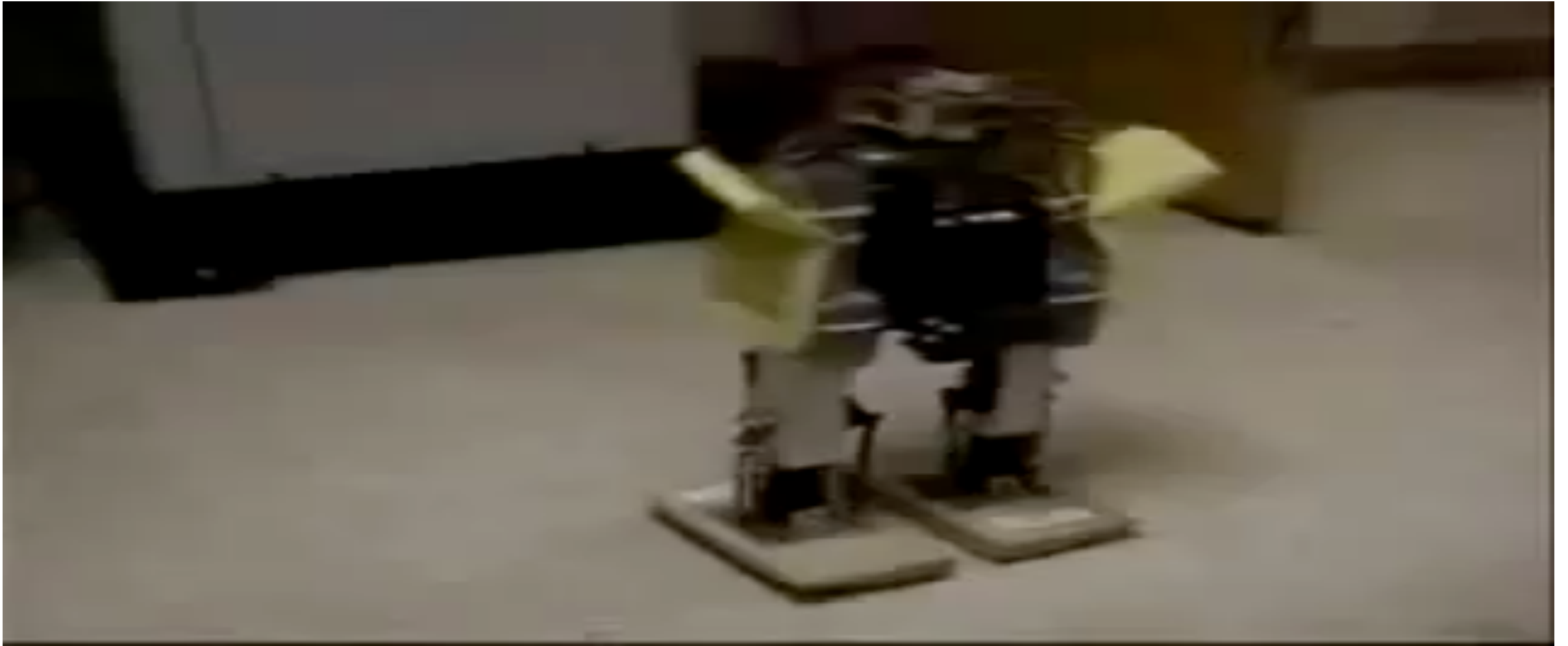
Training

Example: Learning to Walk



Finished

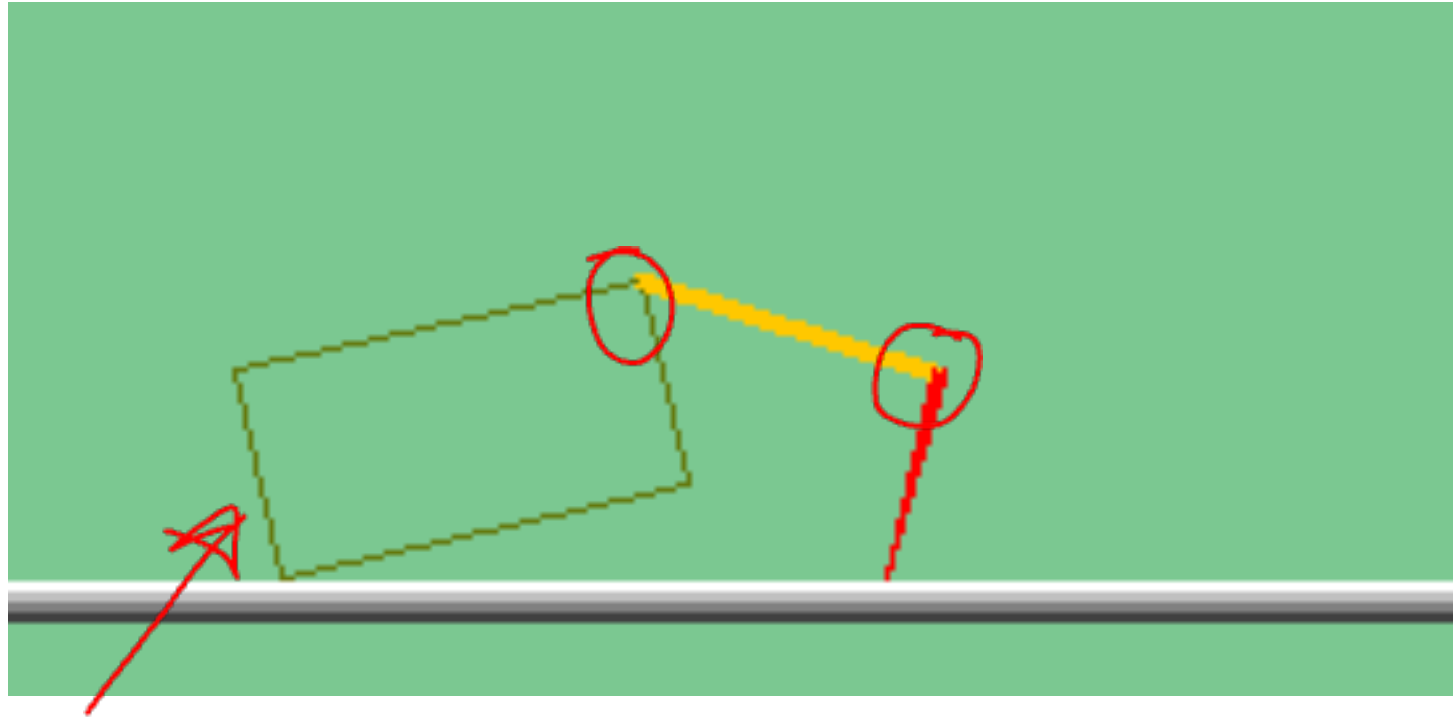
Example: Toddler Robot



Robotics Rubik Cube_e

- <https://www.youtube.com/watch?v=x4O8pojMF0w>

The Crawler!



Video of Demo Crawler Bot

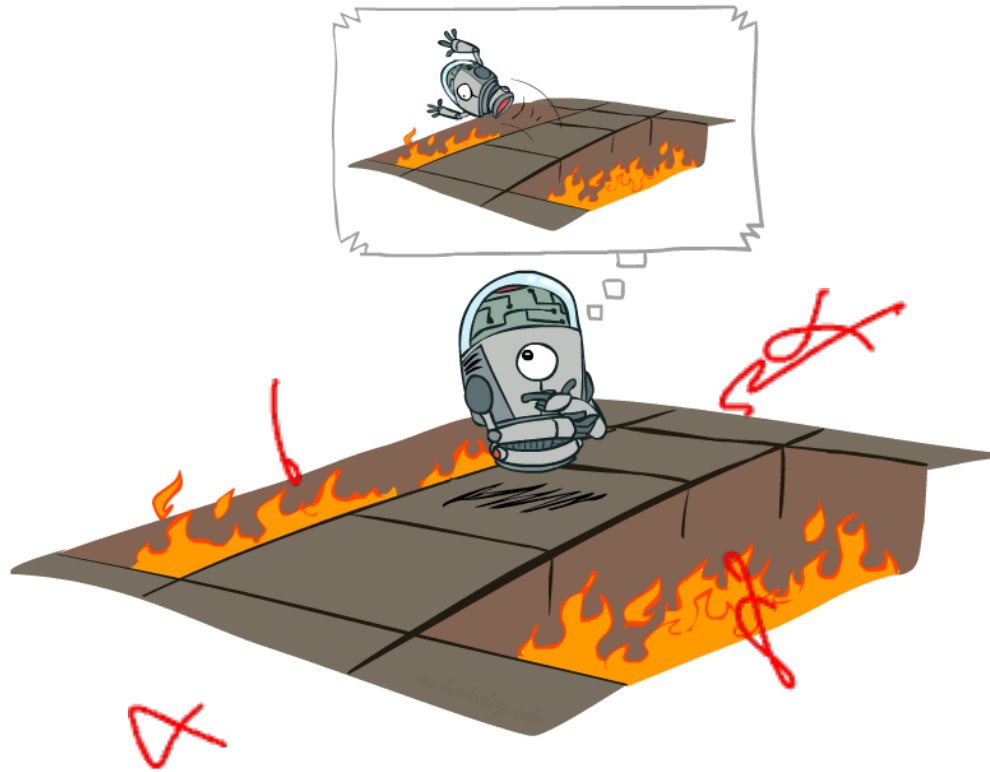


Reinforcement Learning

- 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)

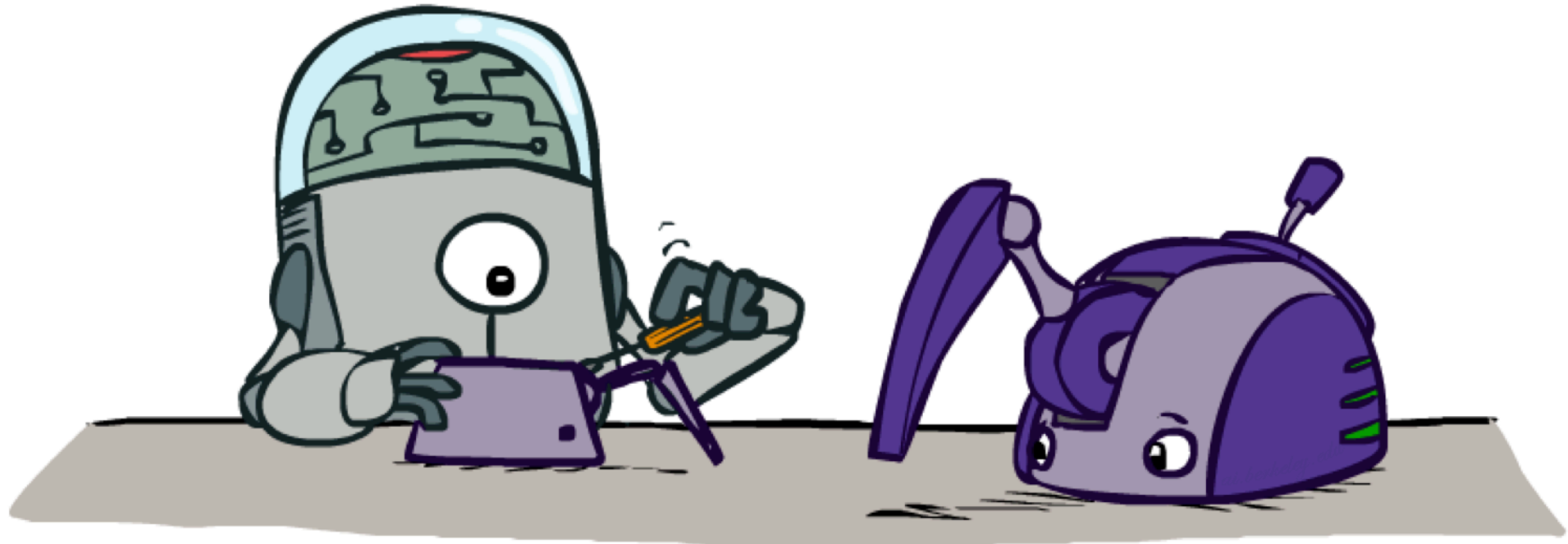


Offline Solution



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 , $\hat{T}(s, a, s')$
- Normalize to $\hat{R}(s, a, s')$ estimate of
- 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



T

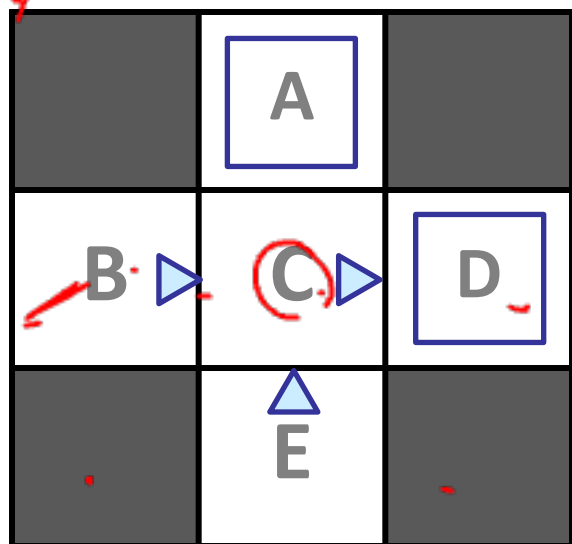


R

$=$

Example: Model-Based Learning

Input Policy π



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
 ...

Model based → Model free

Analogy: Expected Age

40% 22
20% 23

Goal: Compute expected age of cse573 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"

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

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

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

$\hat{P}(4) = \frac{3}{5}$
 $\hat{P}(5) = \frac{2}{5}$
 $\frac{3}{5} \times 4 + \frac{2}{5} \times 5$

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.

$\frac{1}{5} (4 + 4 + 4 + 5 + 5)$

CSE 573:

Artificial Intelligence

Hanna Hajishirzi

Reinforcement Learning

slides adapted from
Dan Klein, Pieter Abbeel ai.berkeley.edu
And Dan Weld, Luke Zettlemoyer



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

- Big Idea: Compute all averages over T using sample outcomes

The Story So Far: MDPs and RL

Known MDP: Offline Solution ✓

Goal

Compute V^*, Q^*, π^*

Evaluate a fixed policy π

Technique

Value / policy iteration

Policy evaluation

Unknown MDP: Model-Based

Goal

Compute V^*, Q^*, π^*

Evaluate a fixed policy π

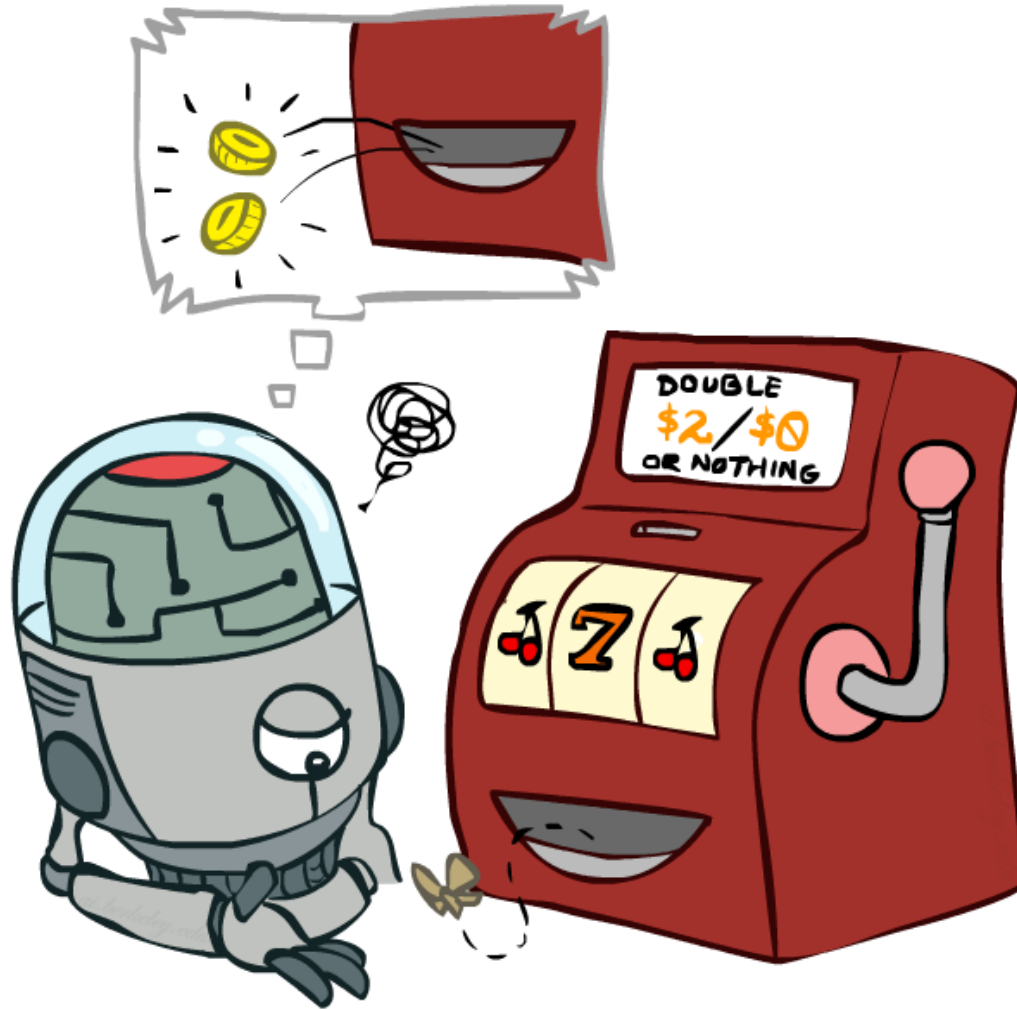
Technique

VI/PI on approx. MDP

PE on approx. MDP

Unknown MDP: Model-Free ✓

Model-Free Learning

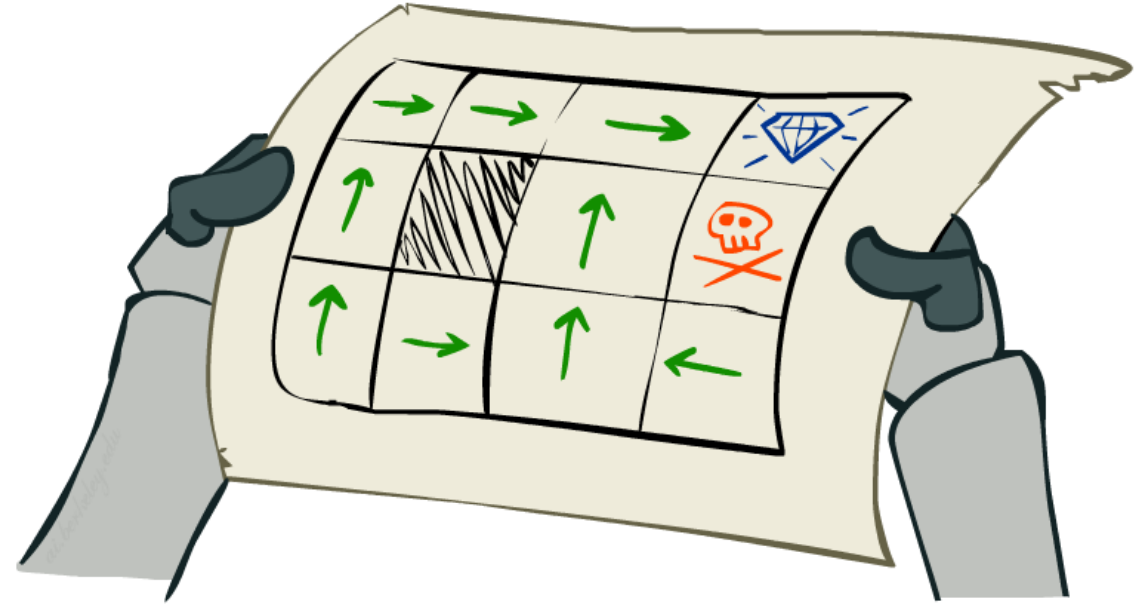


Passive Reinforcement Learning

- 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.



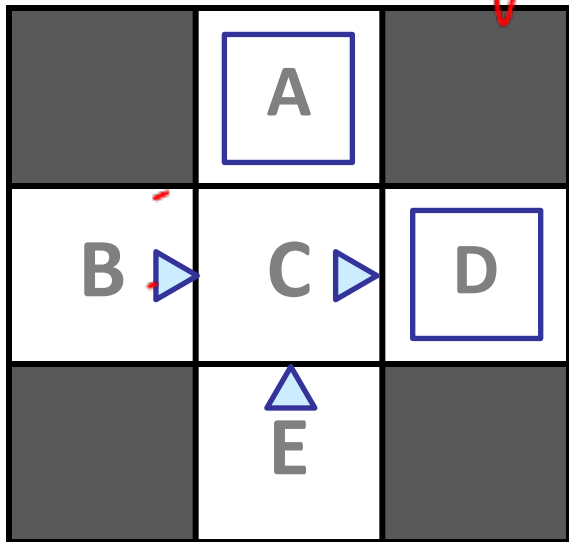
Direct Evaluation

- Goal: Compute values for each state under π
- Idea: Average together observed sample values
 - ~~Act~~ Act according to π
 - 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



Example: Direct Evaluation

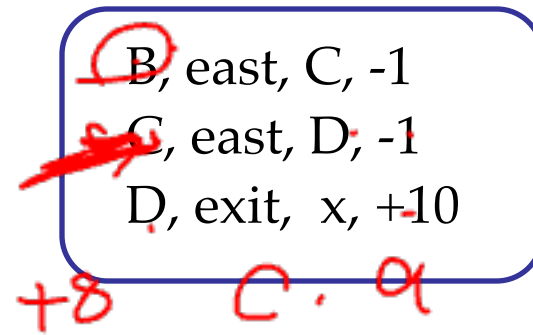
Input Policy π



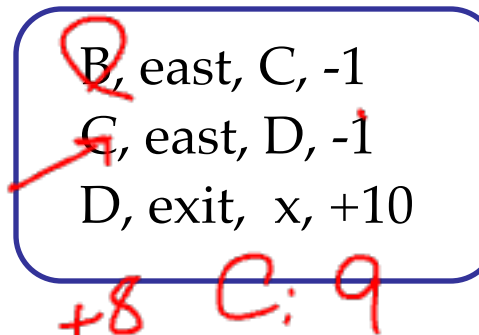
Assume: $\gamma = 1$

Observed Episodes (Training)

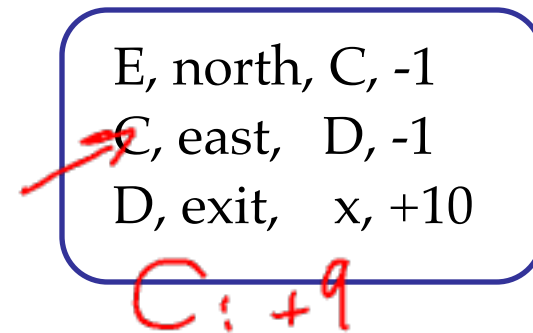
Episode 1



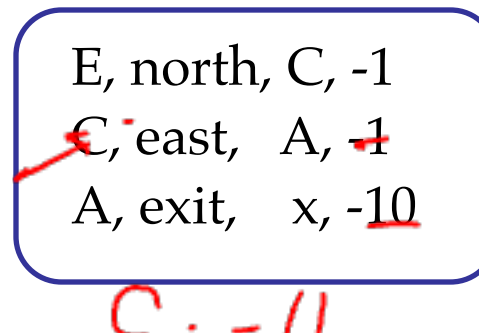
Episode 2



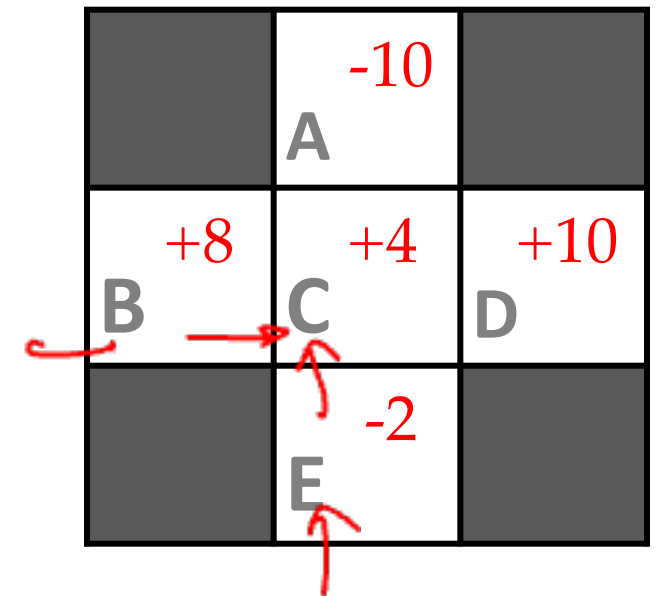
Episode 3



Episode 4



Output Values

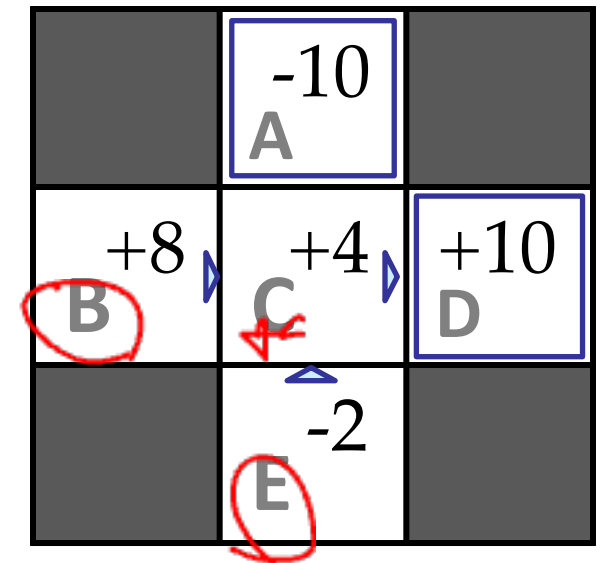


If B and E both go to C under this policy, how can their values be different?

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



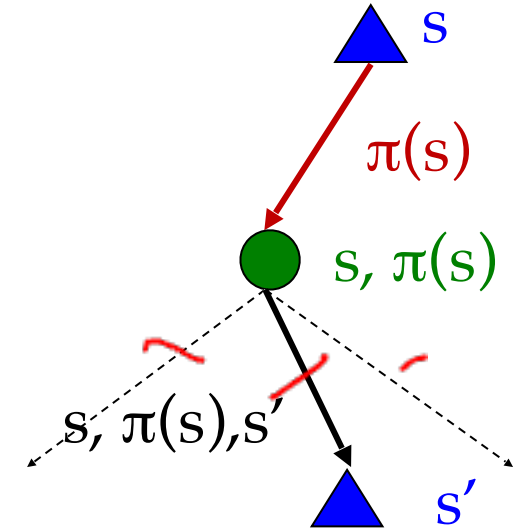
*If B and E both go to C
under this policy, how can
their values be different?*

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

$$\underline{V_0^\pi(s) = 0}$$

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



- 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 π the action) and average

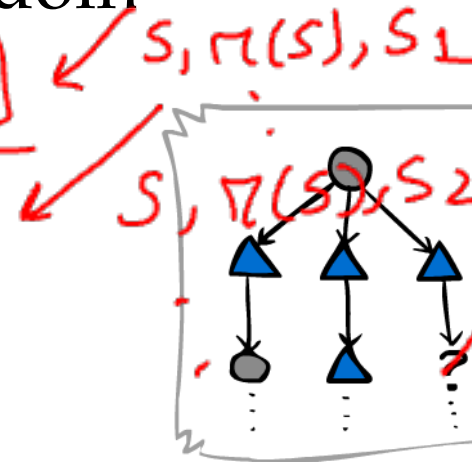
$$sample_1 = R(s, \pi(s), s'_1) + \gamma V_k^{\pi}(s'_1)$$

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

...

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

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

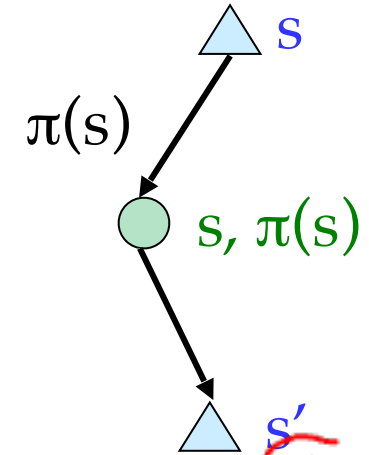


Temporal Difference Learning



- 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

$V^\pi(s)$



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



V 0.1

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

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

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

Z

- ~~Decreasing learning rate~~ (alpha) can give converging averages

$R(s,a,s')$

Example: Temporal Difference Learning

States

	A	
B	C	D
	E	

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

(1- α) $V^\pi + \alpha$ sample

$0 + 0.5 \times -2$
 $0 + 0.5 \times 6$

Observed Transitions

B, east, C, -2

C, east, D, -2

Sample: $-2 + 0$
 $-2 + 1 \times 8 = 6$

	0	
-2	0	-8
	0	

	0	
-1	0	8
	0	

	0	
-1	3	8
	0	

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

sample

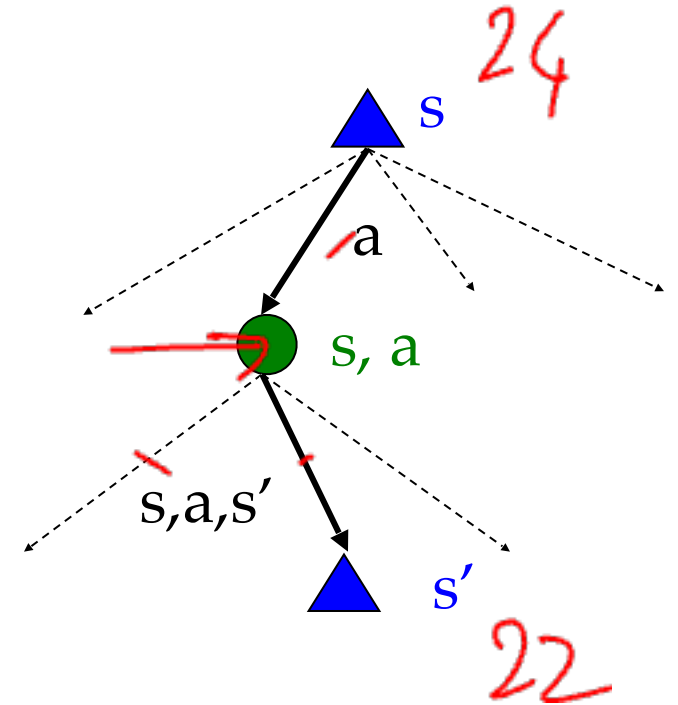
Problems with TD Value Learning

- 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:

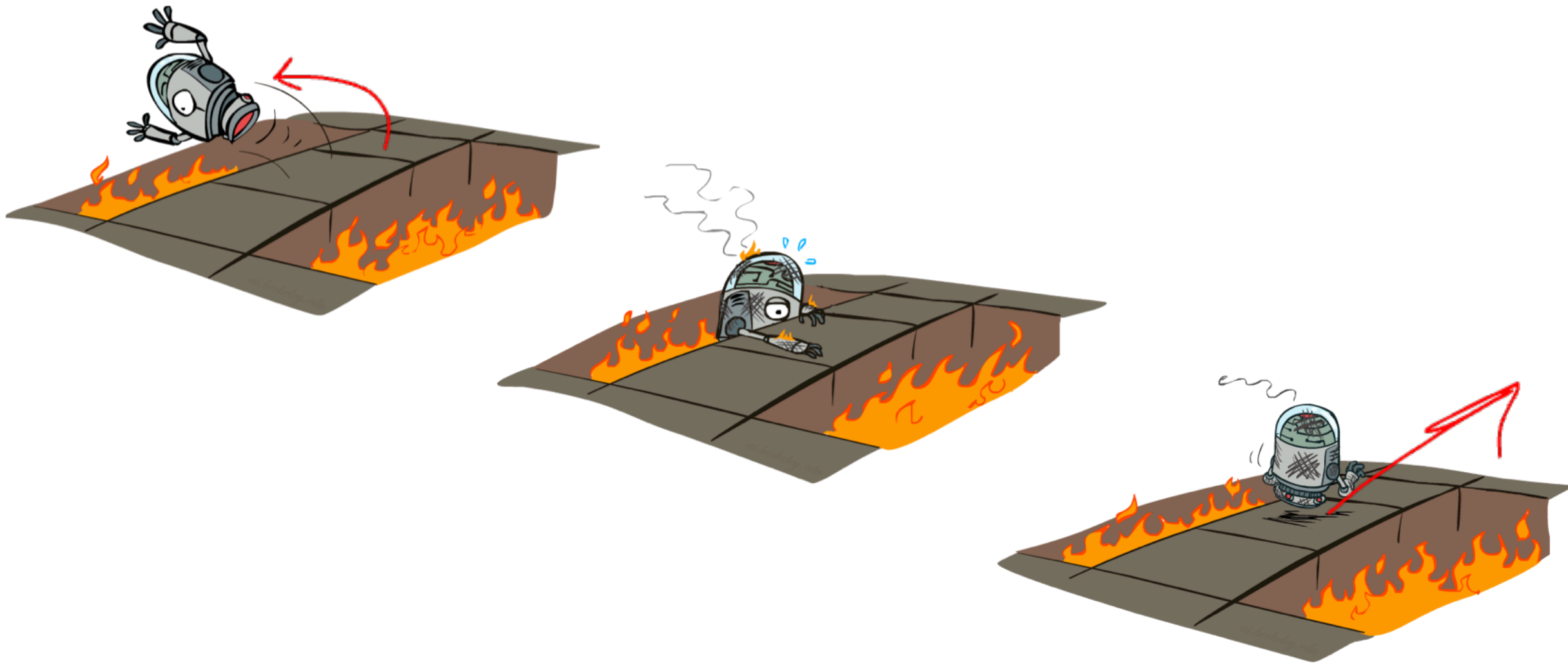
$$\pi(s) = \arg \max_a Q(s, a)$$

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

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

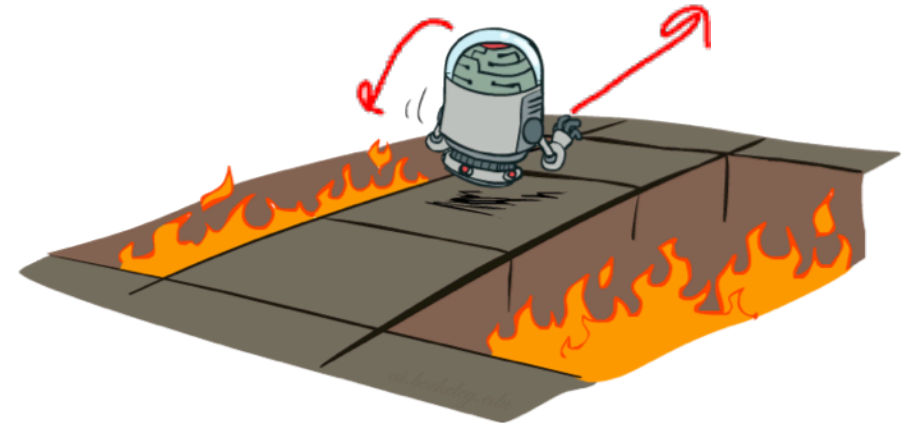


Active Reinforcement Learning



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 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...



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:

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

Handwritten red notes: $Q_{k+1}(s, a) = \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma \max_{a'} Q_k(s', a')]$



- 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 $k+1$ values for all s, a states:

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

Q-Learning

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

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

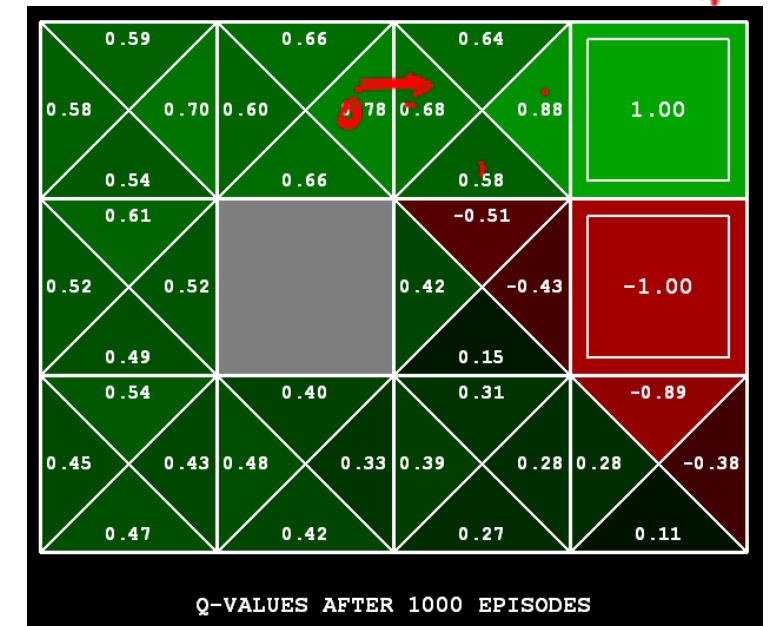
- Learn $Q(s,a)$ values as you go

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

$$sample = \underline{R(s, a, s')} + \gamma \max_{a'} Q(s', a') \text{ no longer policy evaluation!}$$

- Incorporate the new estimate into a running average:

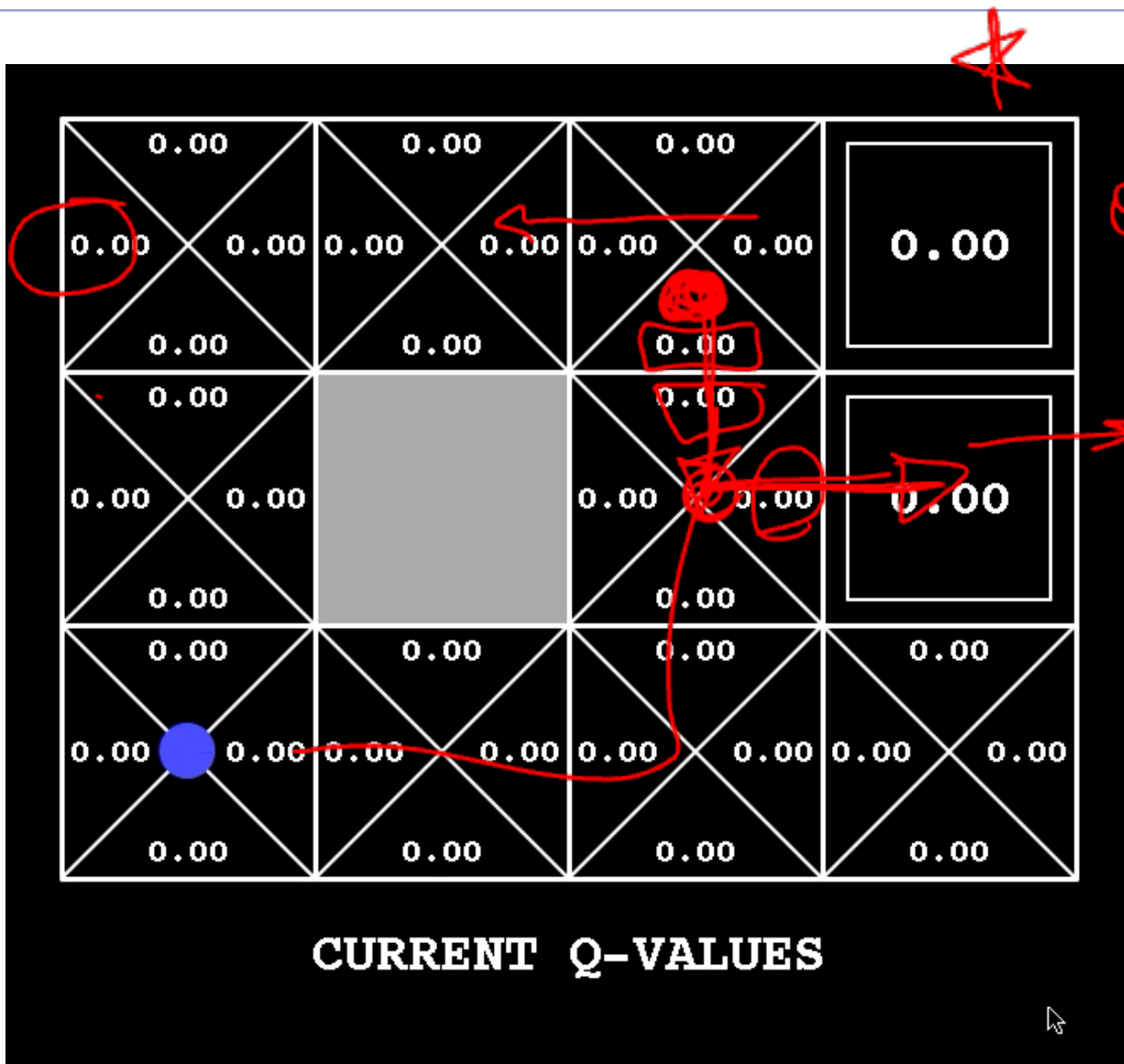
$$\underline{Q(s, a)} \leftarrow \underline{(1 - \alpha)Q(s, a) + (\alpha) [sample]}$$



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

[Demo: Q-learning – crawler (L10D3)]

Q-Learning Demo

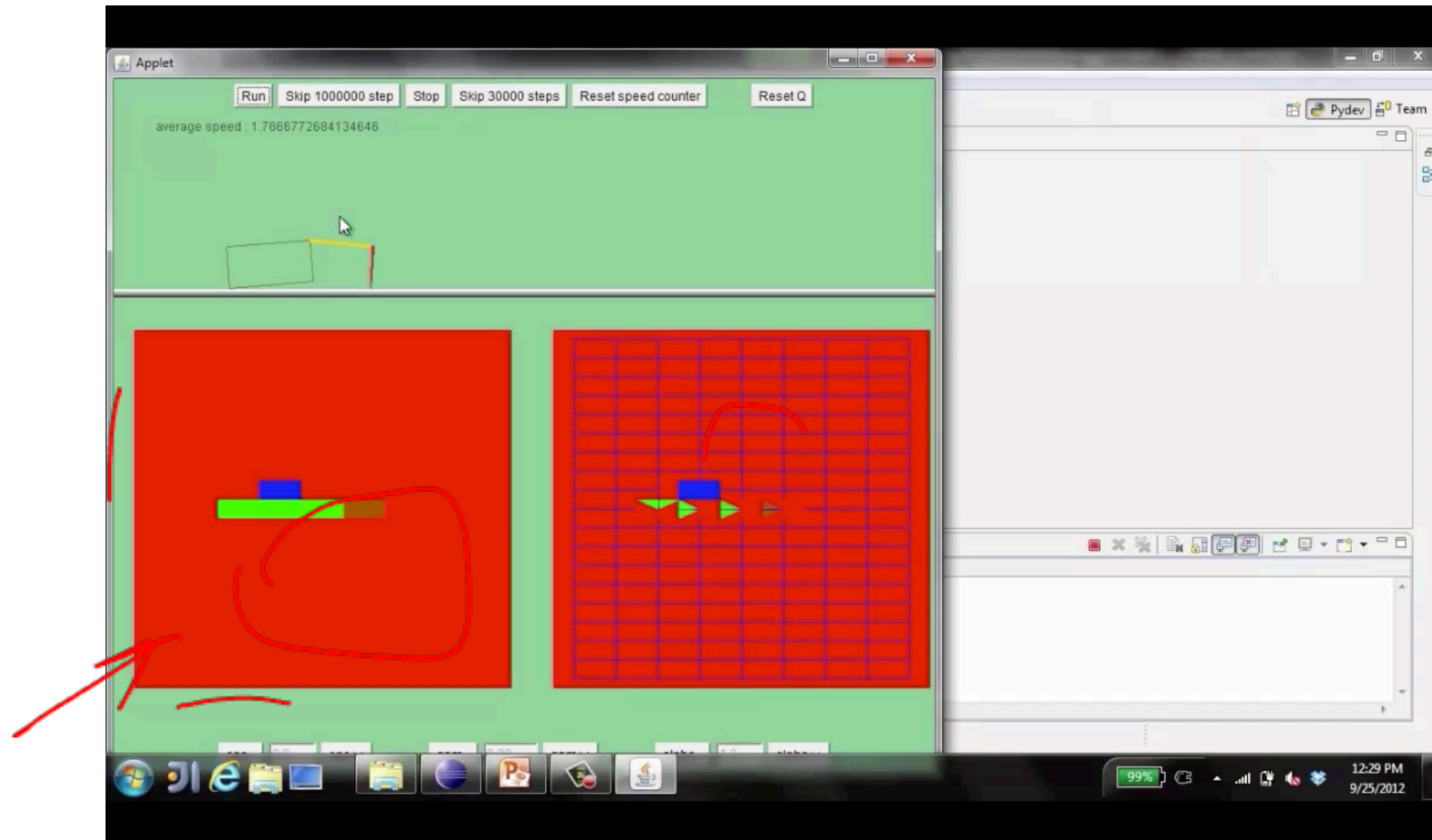
 $+2$

0.5

$$\alpha = 0.5$$

$$\delta = 1$$

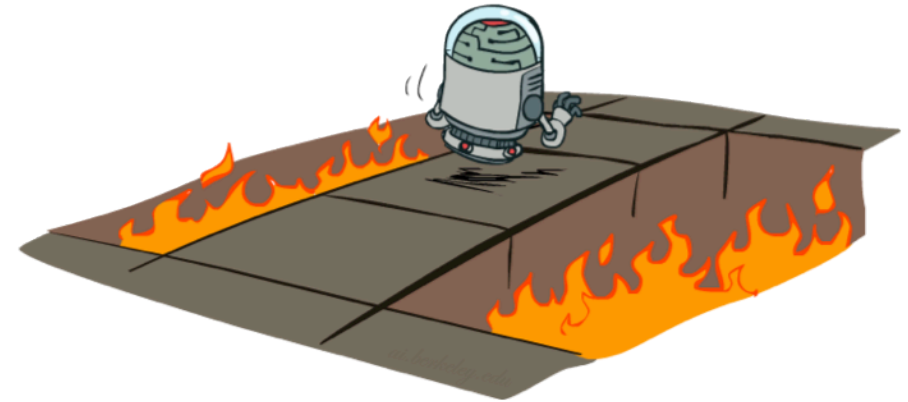
Video of Demo Q-Learning -- Crawler



Q-Learning:

act according to current optimal (and also explore...)

- 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...



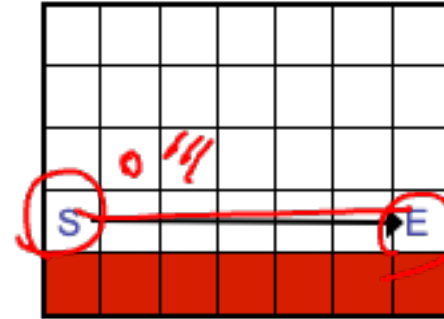
Q-Learning Properties

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

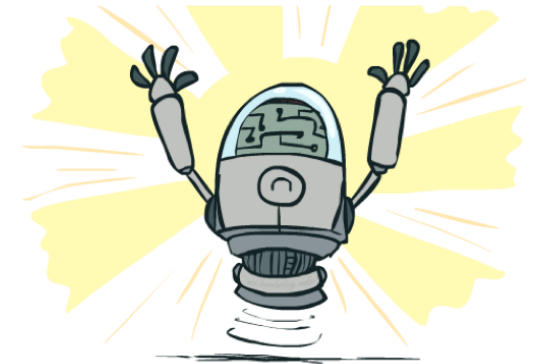
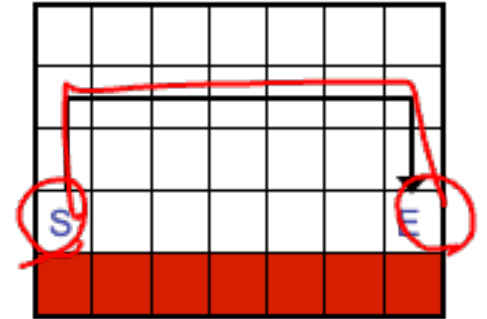
- This is called **off-policy learning**

- 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 (!)



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



Discussion: Model-Based vs Model-Free RL
