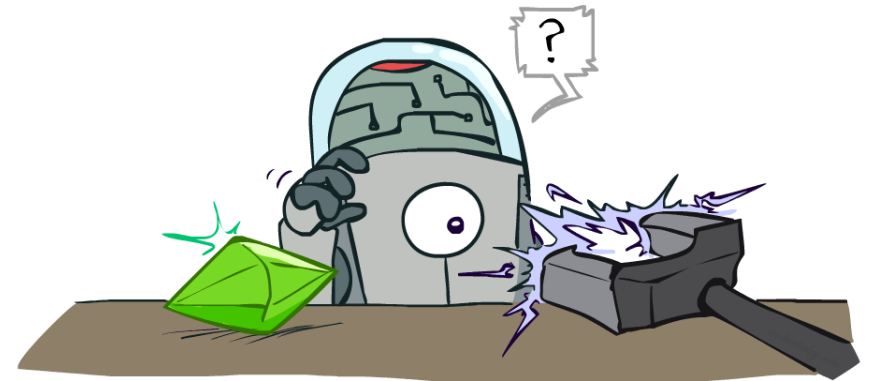


---

# CSE 473: Introduction to Artificial Intelligence

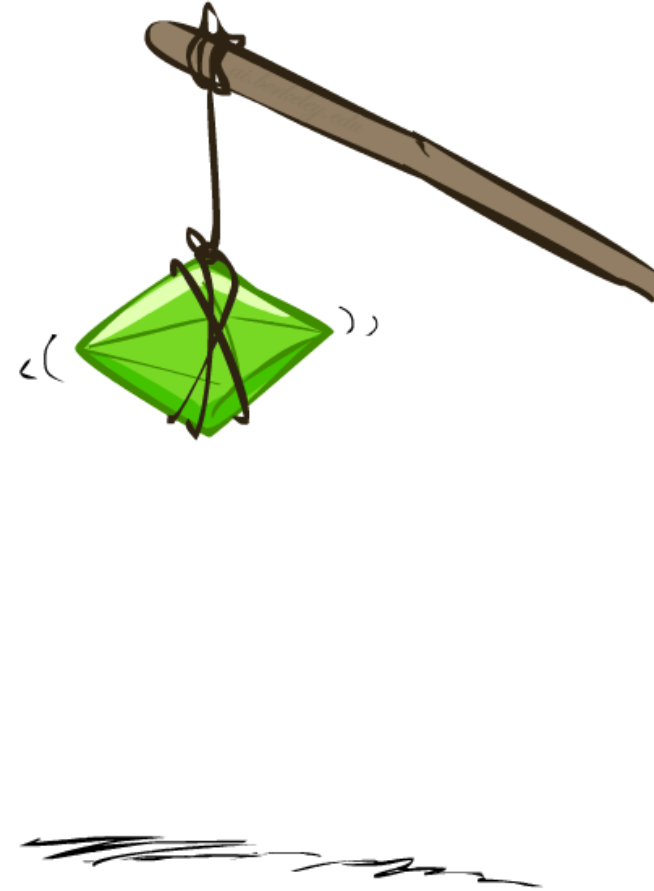
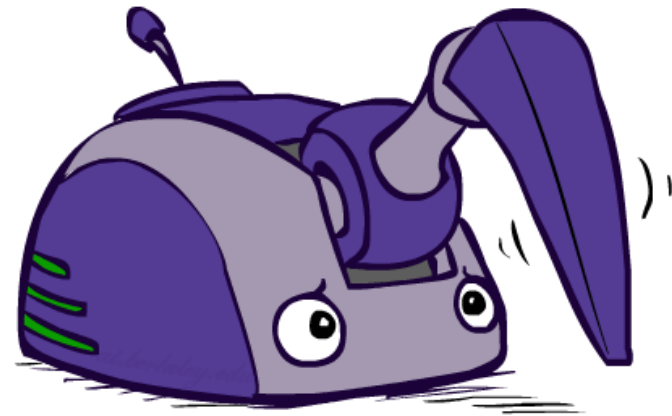
Hanna Hajishirzi  
Reinforcement Learning

slides adapted from  
Dan Klein, Pieter Abbeel [ai.berkeley.edu](http://ai.berkeley.edu)  
And Dan Weld, Luke Zettlemoyer



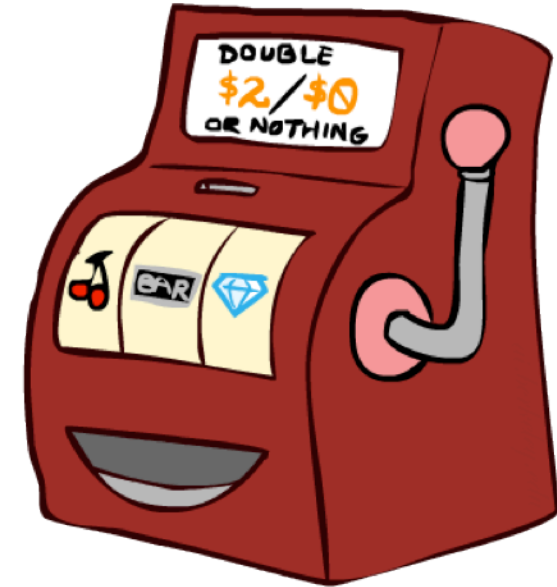
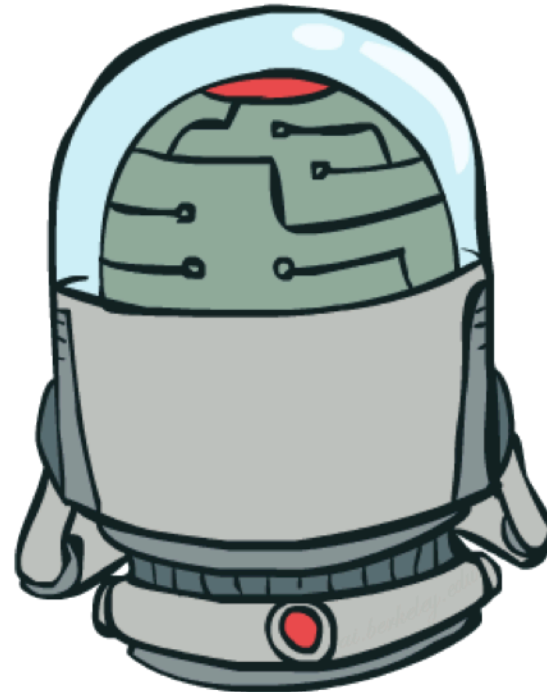
# Reinforcement Learning

---



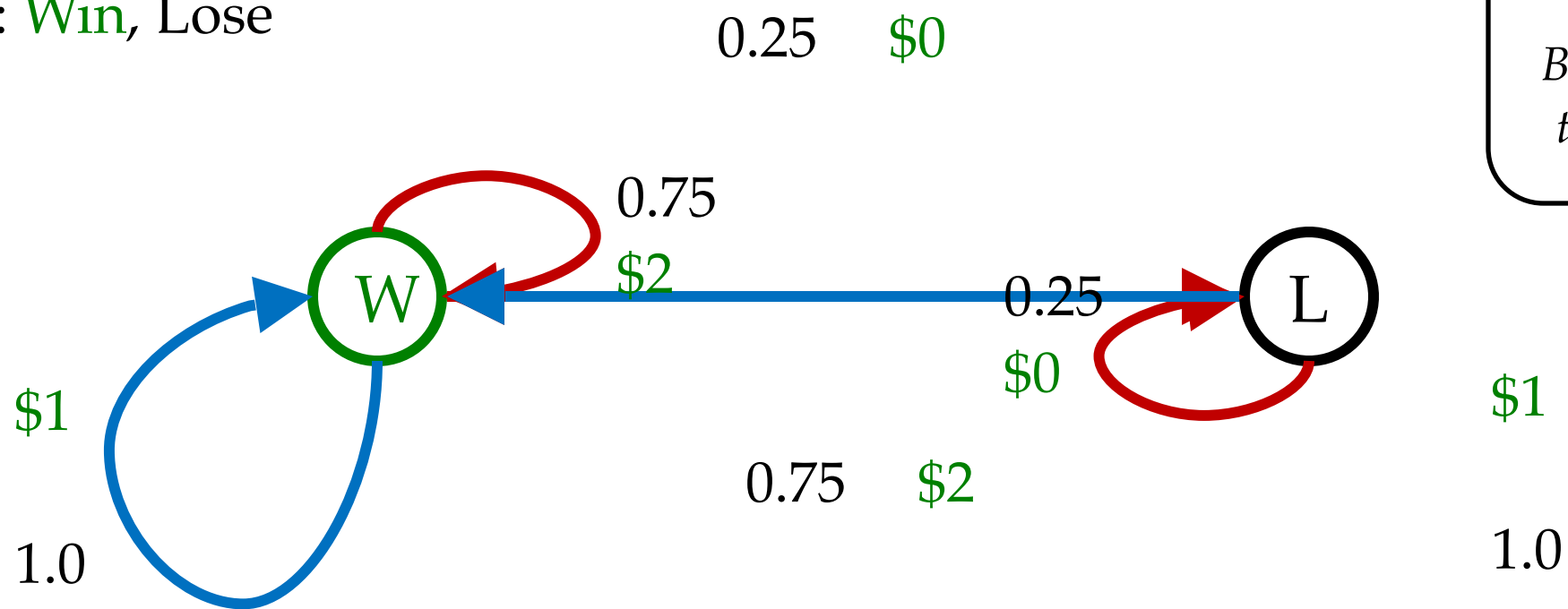
# Double Bandits

---



# Double-Bandit MDP

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



No discount  
10 time steps  
Both states have  
the same value

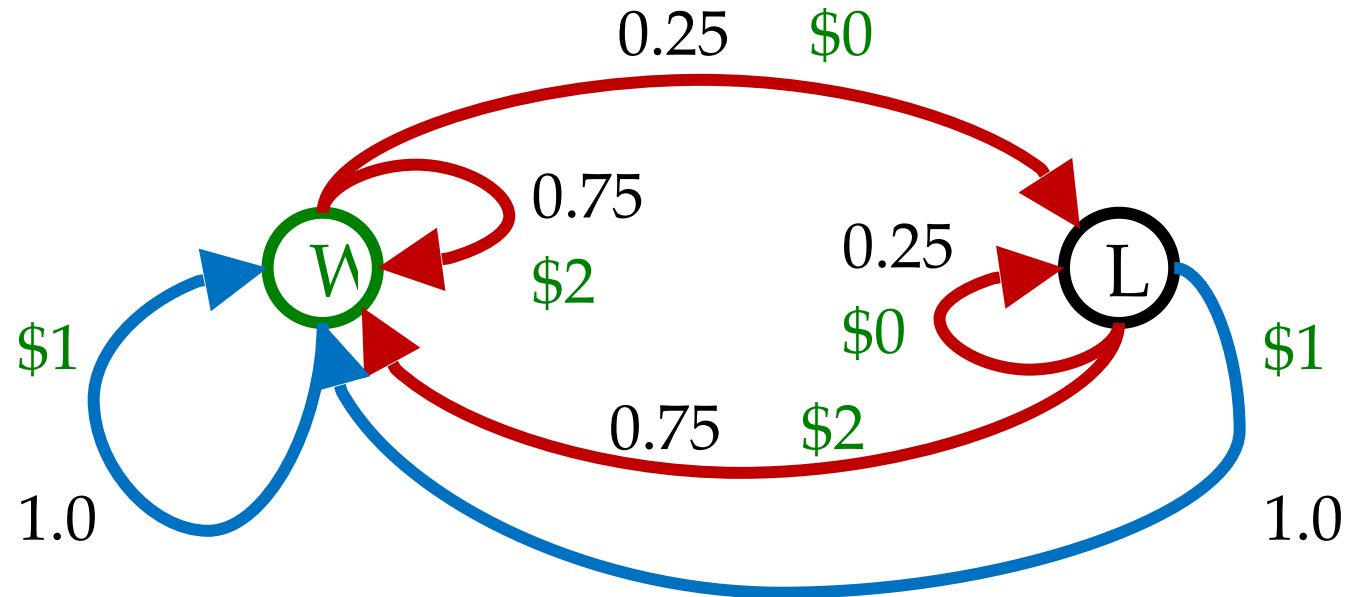


# Offline Planning

- 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!

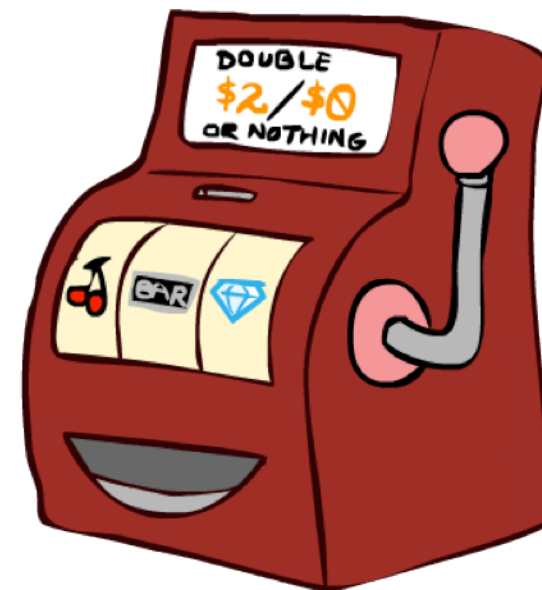
*No discount*  
*10 time steps*

	Value
Play Red	15
Play Blue	10



# Let's Play!

---

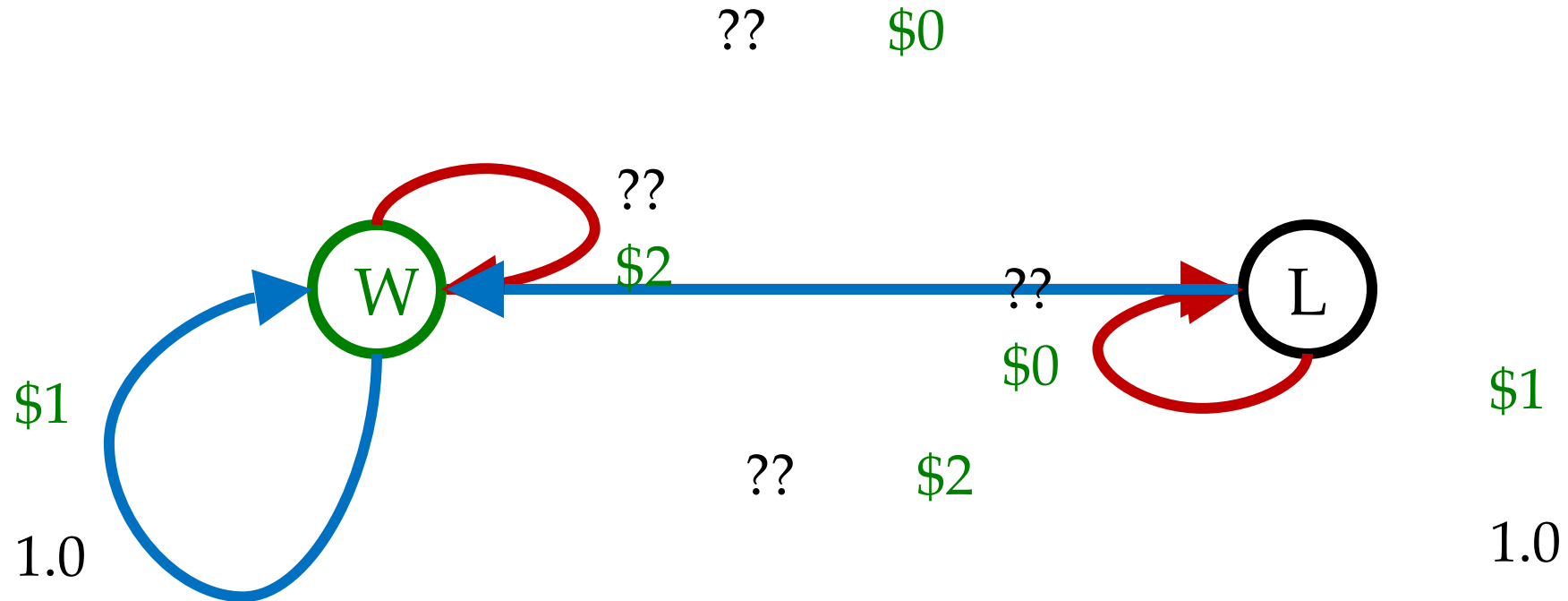


\$2 \$2 \$0 \$2 \$2

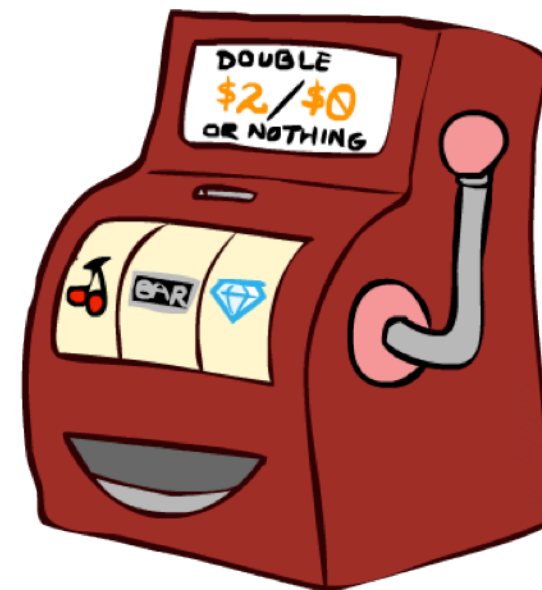
\$2 \$2 \$0 \$0 \$0

# 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

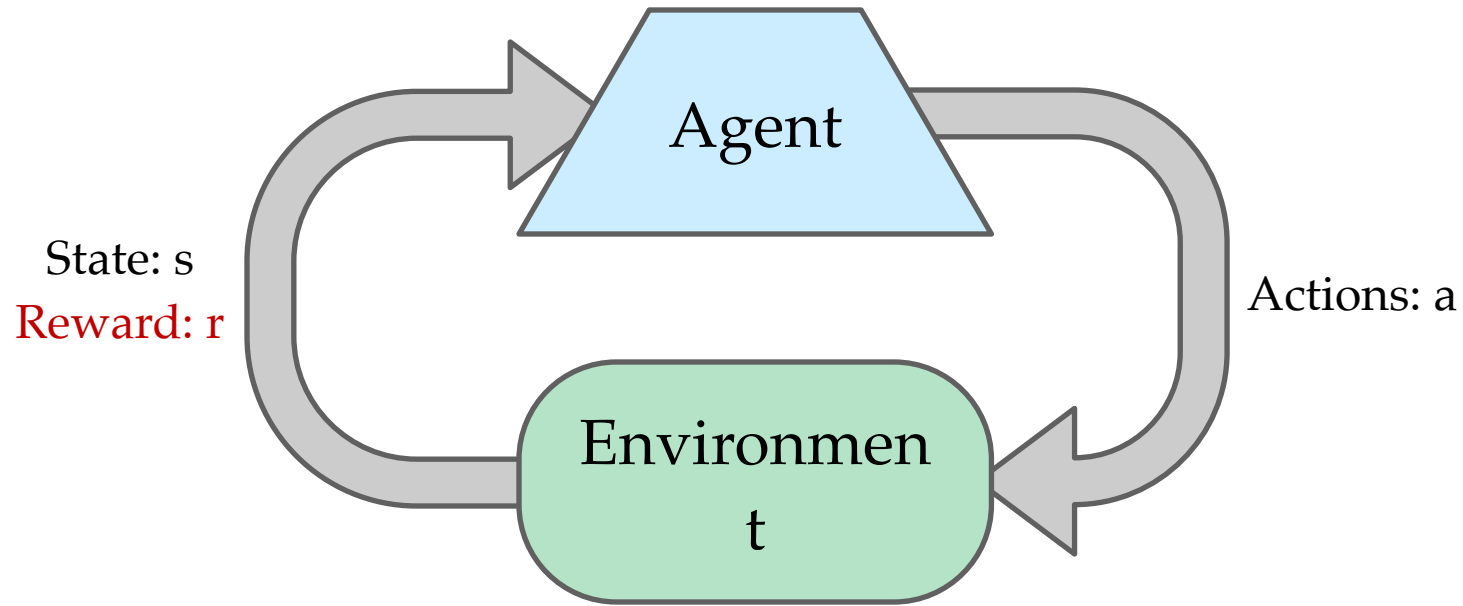
---

- 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: Toddler Robot

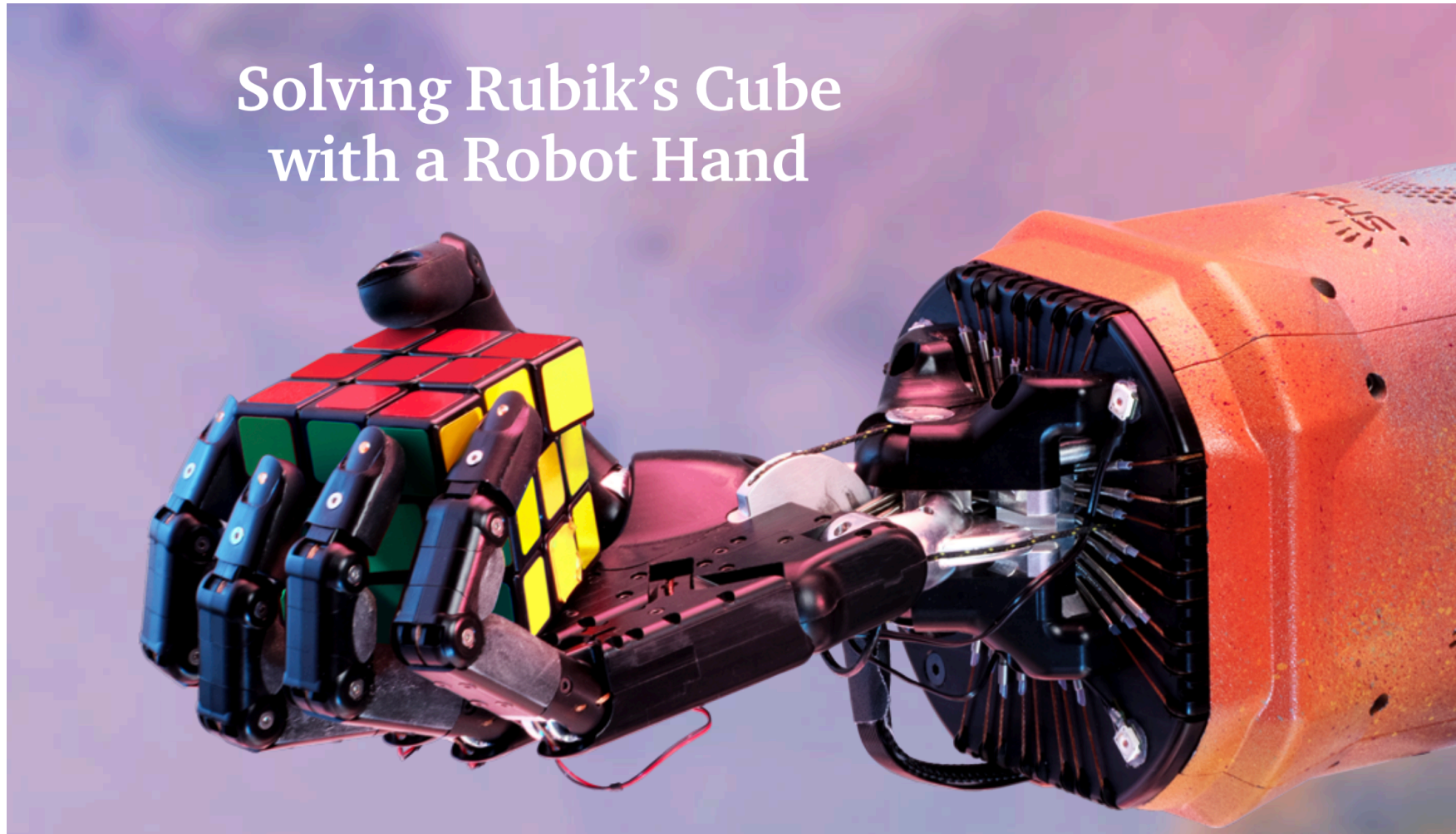
---



# Robotics Rubik Cube

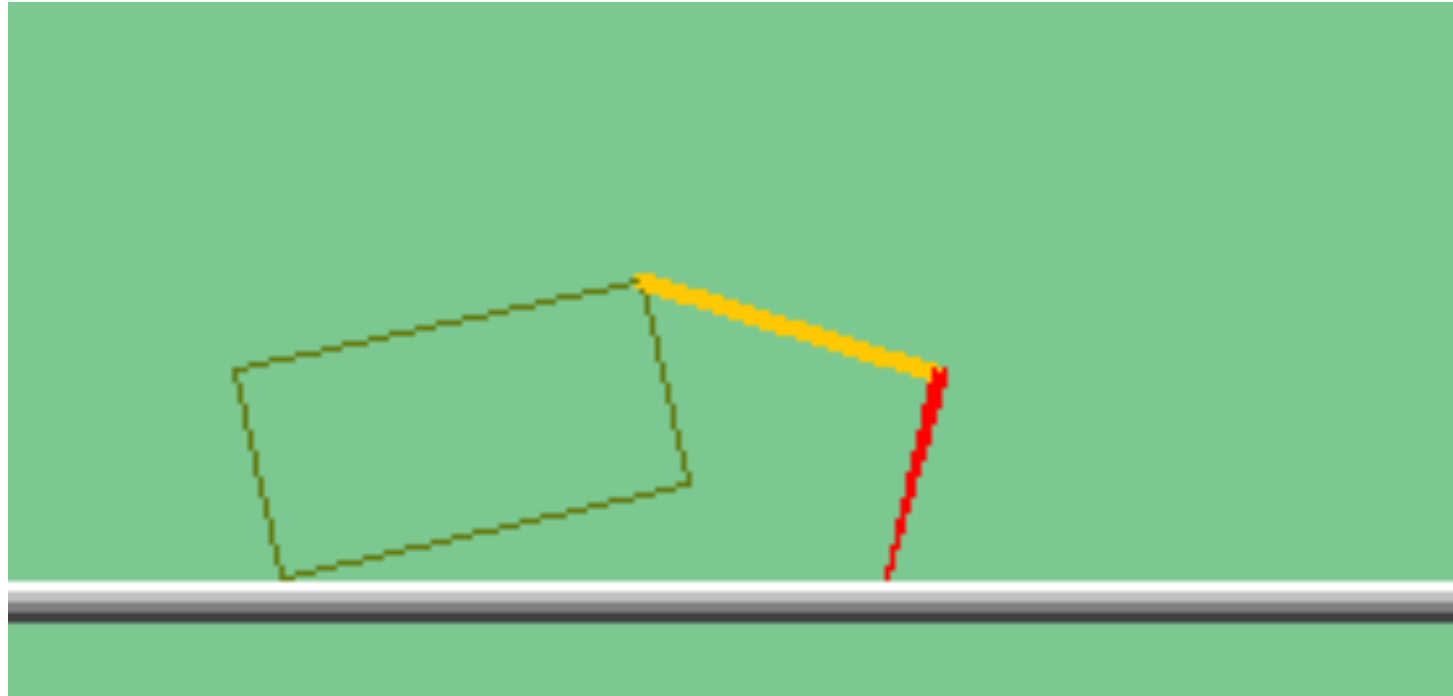
---

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



# The Crawler!

---



# Video of Demo Crawler Bot

---

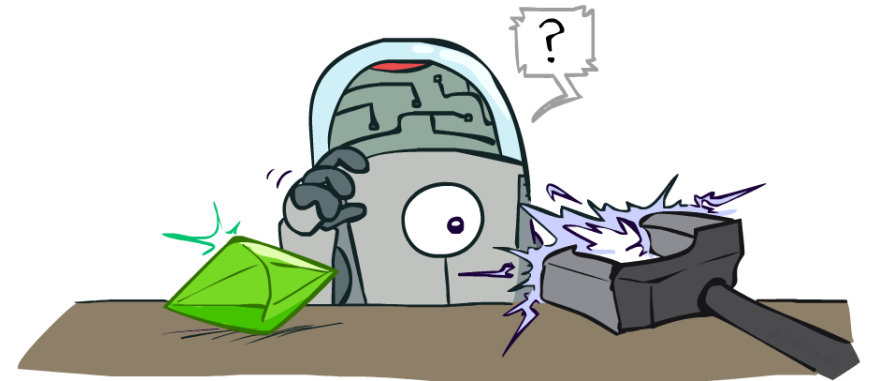


---

# CSE 473: Introduction to Artificial Intelligence

Hanna Hajishirzi  
Reinforcement Learning

slides adapted from  
Dan Klein, Pieter Abbeel [ai.berkeley.edu](http://ai.berkeley.edu)  
And Dan Weld, Luke Zettlemoyer



# Announcements

---

- HW2 grades will be released soon.
- PS3 will be released soon (Due; May 12<sup>th</sup>)
- HW3 will be released on Tue afternoon (Due, May 7<sup>th</sup>)
- Mid-quarter course evaluations:
  - <https://uw.iasystem.org/survey/240219>

# Reinforcement Learning

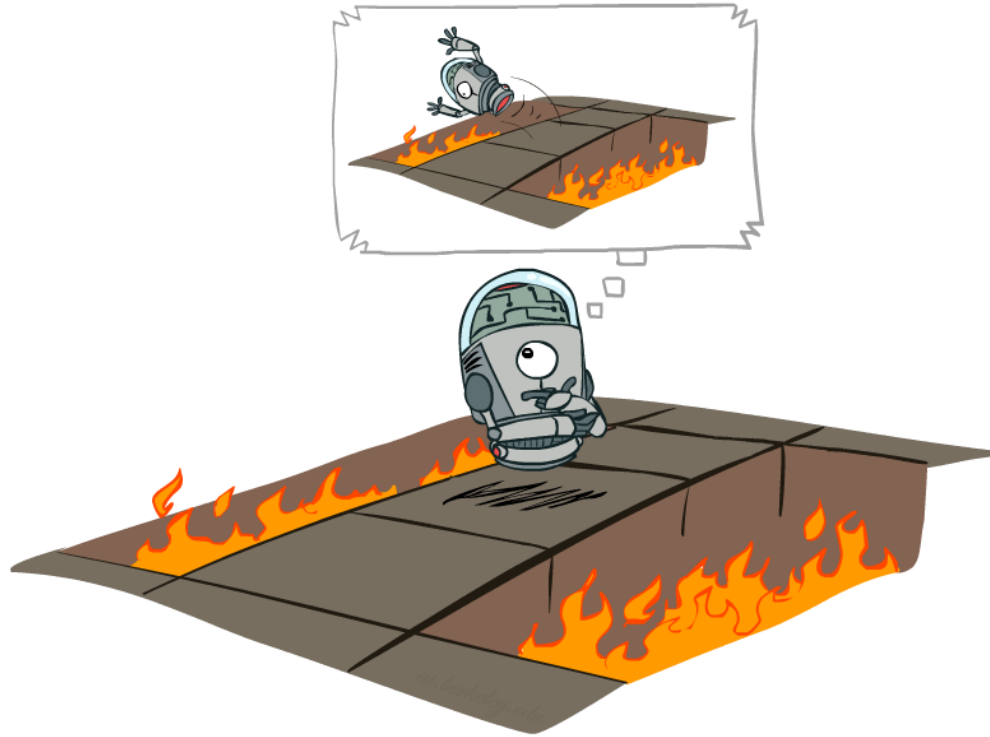
---

- 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

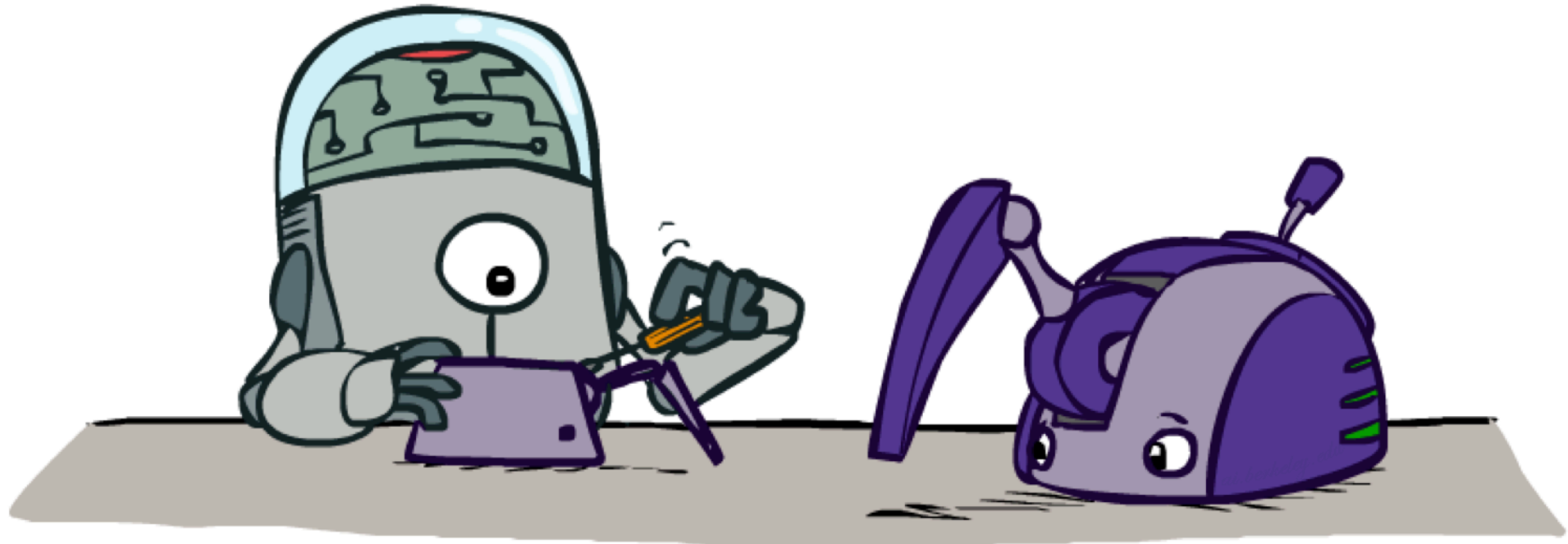


Online Learning



# Model-Based Learning

---



# Model-Based Learning

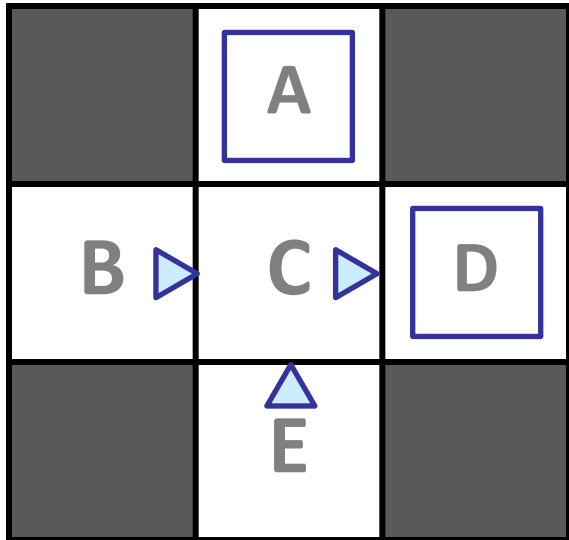
---

- 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  $g\hat{R}(s, a, s')$  estimate of
  - Discover each  $\hat{T}(s, a, s')$  when we experience  $(s, a, s')$
- Step 2: Solve the learned MDP
  - For example, use value iteration, as before



# Example: Model-Based Learning

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

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

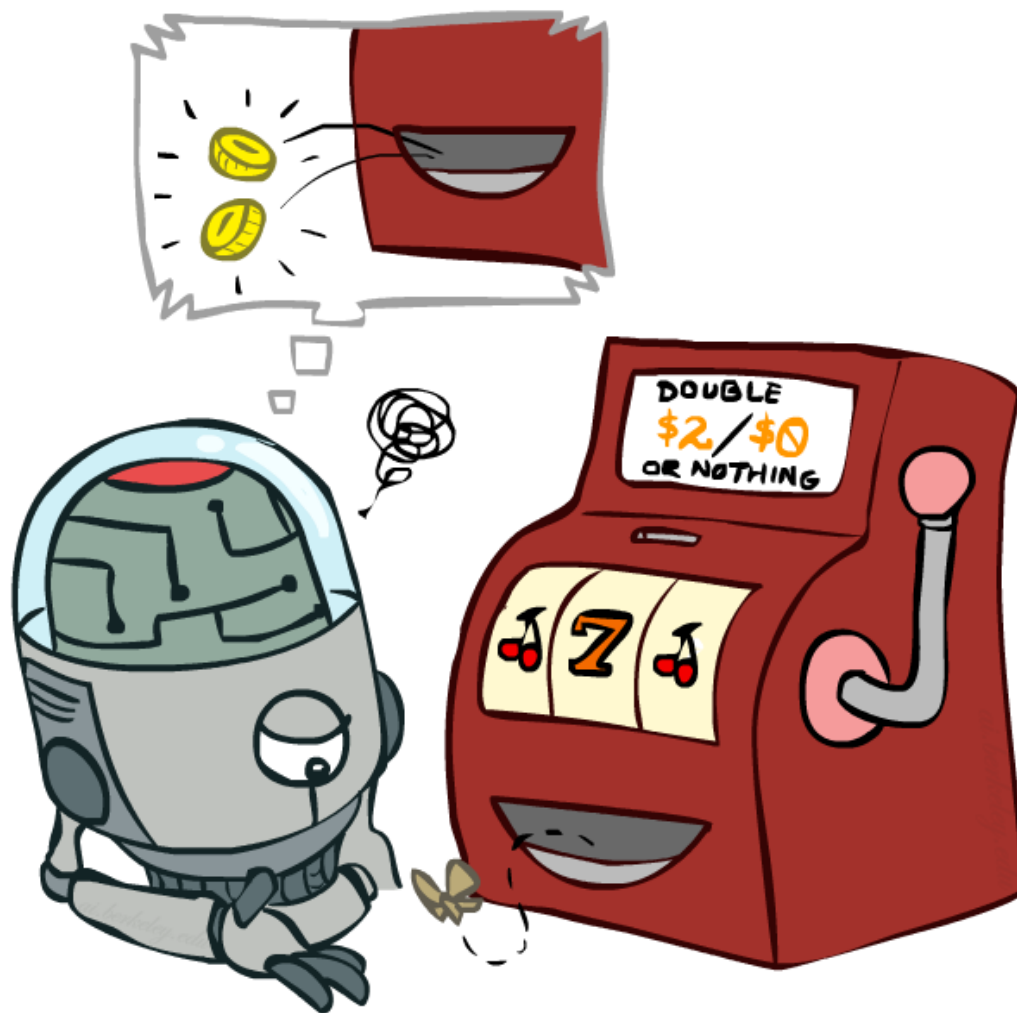
T(B, east, C) = 1.00  
T(C, east, D) = 0.75  
T(C, east, A) = 0.25  
...

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

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

# Model-Free Learning

---



# Direct Evaluation

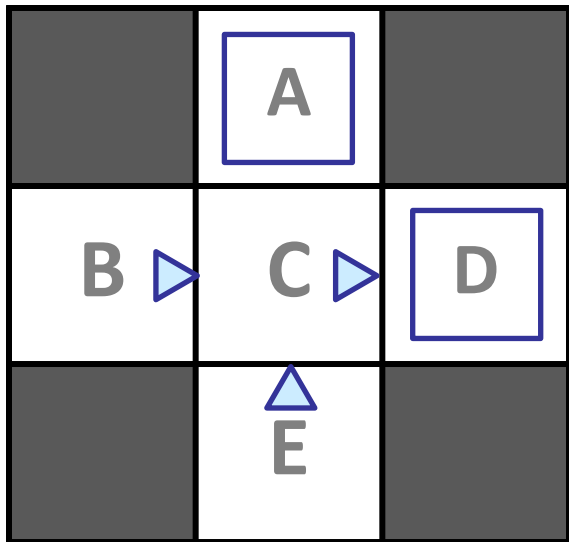
---

- 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



# 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	+4	+10
B	C	D
	-2	
	E	

*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

	-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?*

---

# CSE 473: Introduction to Artificial Intelligence

Hanna Hajishirzi  
Reinforcement Learning

slides adapted from  
Dan Klein, Pieter Abbeel [ai.berkeley.edu](http://ai.berkeley.edu)  
And Dan Weld, Luke Zettlemoyer





# The Story So Far: MDPs and RL

$S, A, T, R$

## Known MDP: Offline Solution

Goal

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

Evaluate a fixed policy  $\pi$

Technique

Value / policy iteration

Policy evaluation

$2R$

## Unknown MDP: Model-Based

Goal

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

Evaluate a fixed policy  $\pi$

Technique

VI/PI on approx. MDP

PE on approx. MDP

## Unknown MDP: (Model-Free)

Goal

Evaluate a fixed policy  $\pi$

Technique

direct evaluation

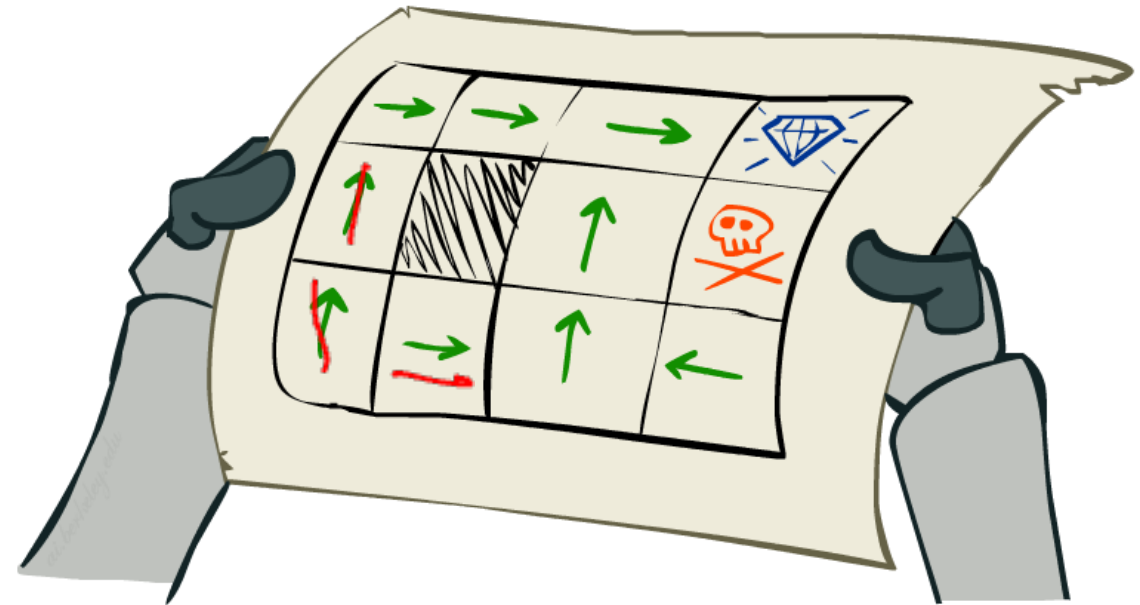
# Passive Reinforcement Learning

- 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:

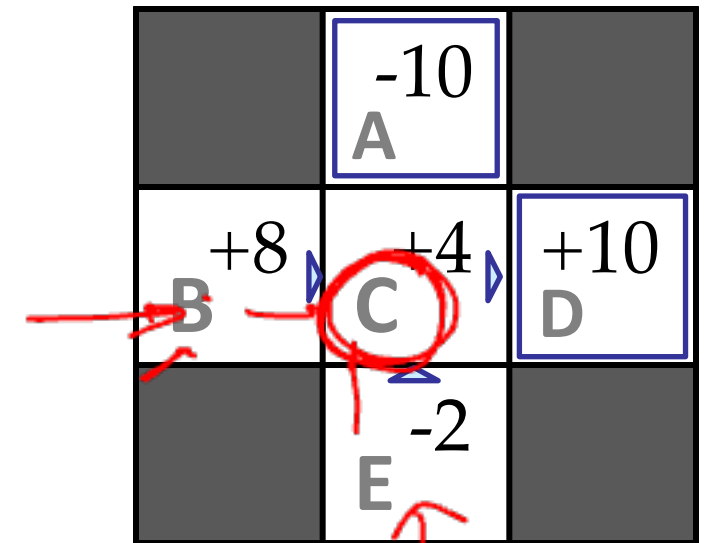
- Learner is "along for the ride"
- No choice about what actions to take
- Just execute the policy and learn from experience
- This is NOT offline planning! You actually take actions in the world.



# 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

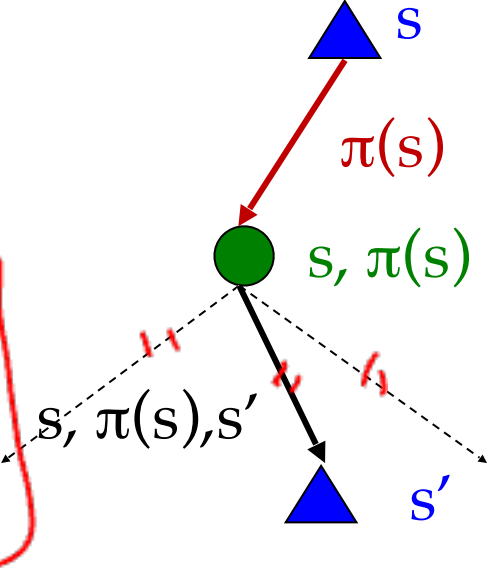
## Output Values



*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  $V$

$$V_0^\pi(s) = 0$$
$$V_{k+1}^\pi(s) \leftarrow \sum_{s'} T(s, \pi(s), s') [R(s, \pi(s), s') + \gamma V_k^\pi(s')]$$


The diagram shows a green circle representing a state-action pair  $(s, \pi(s))$ . A red arrow labeled  $\pi(s)$  points to a blue triangle labeled  $s$ . A black arrow points to a blue triangle labeled  $s'$ . Dashed lines represent transitions to other states, with red arrows indicating the transition probabilities  $T(s, \pi(s), s')$ .

- 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:

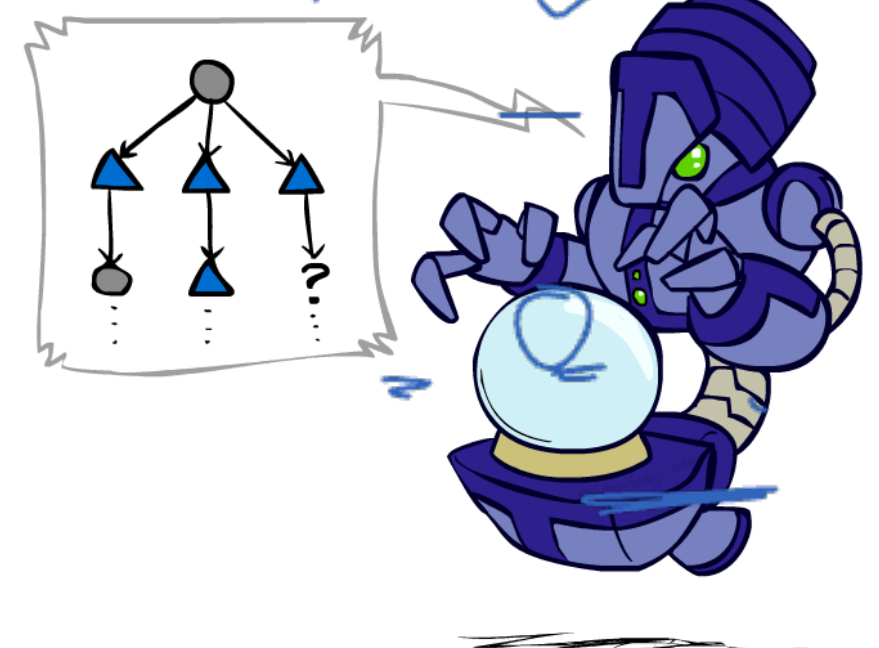
$$V_{k+1}^\pi(s) \leftarrow \sum_{s'} T(s, \pi(s), s') [R(s, \pi(s), s') + \gamma V_k^\pi(s')]$$

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

$$\begin{aligned} \text{sample}_1 &= R(s, \pi(s), s'_1) + \gamma V_k^\pi(s'_1) \\ \text{sample}_2 &= R(s, \pi(s), s'_2) + \gamma V_k^\pi(s'_2) \\ &\dots \\ \text{sample}_n &= R(s, \pi(s), s'_n) + \gamma V_k^\pi(s'_n) \end{aligned}$$

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

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



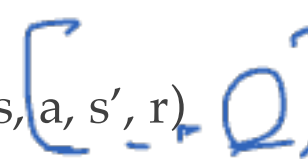
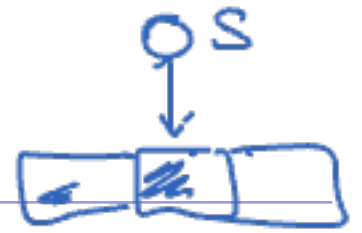
# Temporal Difference Learning

- Big idea: learn from every experience!

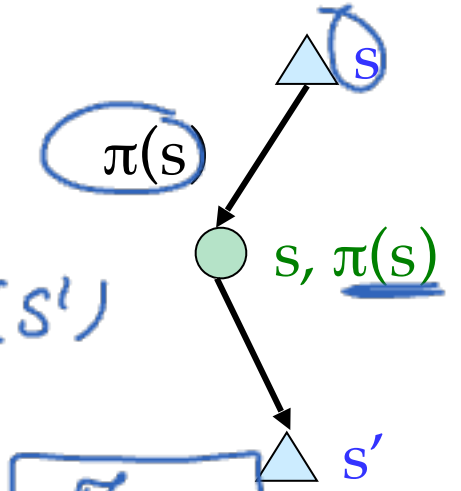
- Update  $V(s)$  each time we experience a transition  $(s, a, s', r)$
- Likely outcomes  $s'$  will contribute updates more often

- Temporal difference learning of values

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



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



$(1 - \alpha) V^{\pi}(s) + \alpha \text{Sample}$

$\alpha = 0.1$

Sample of  $V(s)$ :

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

Update to  $V(s)$ :

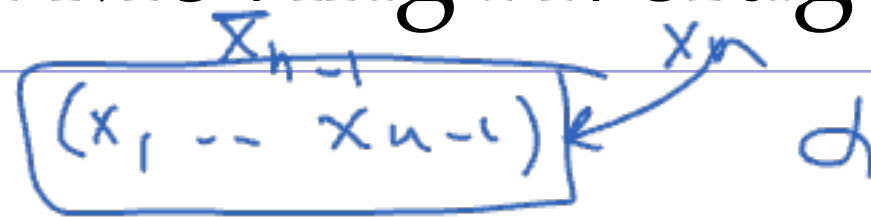
$V^{\pi}(s) \leftarrow (1 - \alpha) V^{\pi}(s) + (\alpha) sample$

Same update:

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

# Exponential Moving Average

$$\frac{x_1 + \dots + x_n}{n}$$



- 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

$$\frac{\alpha \cdot x_n + \alpha(1-\alpha) \cdot x_{n-1} + \dots}{\alpha + \alpha(1-\alpha) + \dots}$$

# Example: Temporal Difference Learning

Sample =  $R + \gamma \cdot V^\pi(s')$

$-2 + 1 \times 8 = 6$   
 $(1 - \alpha) \cdot 0 + \alpha \cdot 6 = 3$

States

$V^\pi(s) = (1 - \alpha)V^\pi(s) + \alpha \cdot \text{Sample}$

Observed Transitions

B, east, C, -2

C, east, D, -2

	A	
B	C	D
	E	

	0	
0	0	8
	0	

	0	
-1	0	8
	0	

	0	
-1	3	8
	0	

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

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

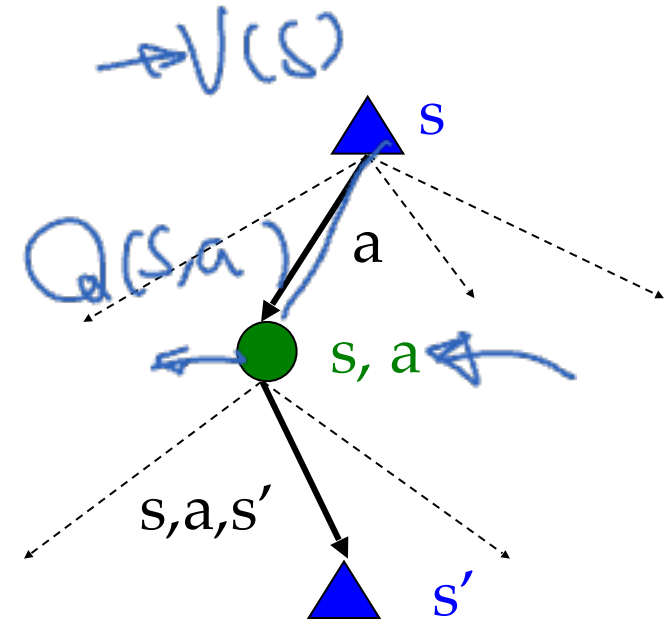


# Problems with TD Value Learning

- 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:

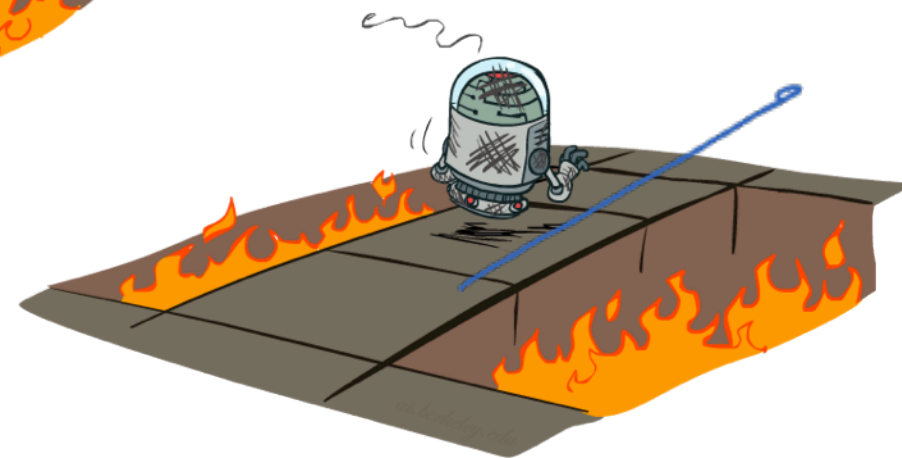
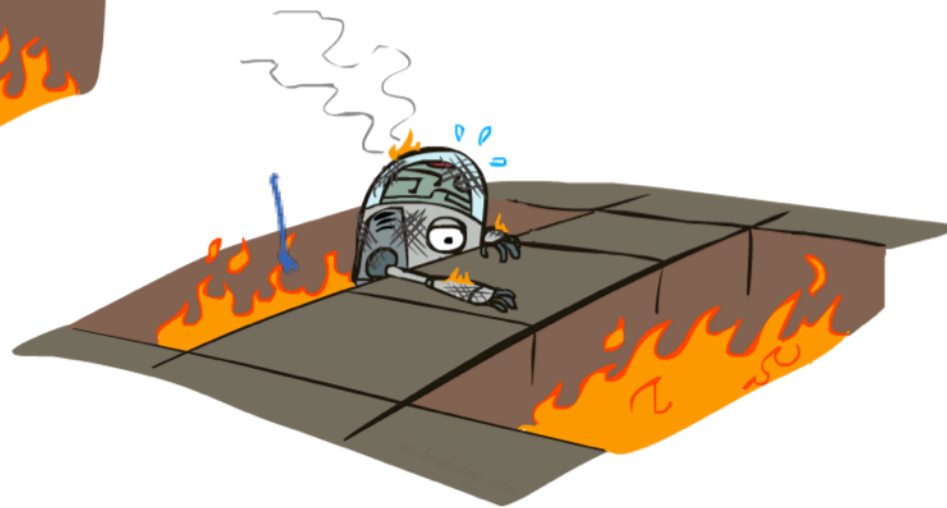
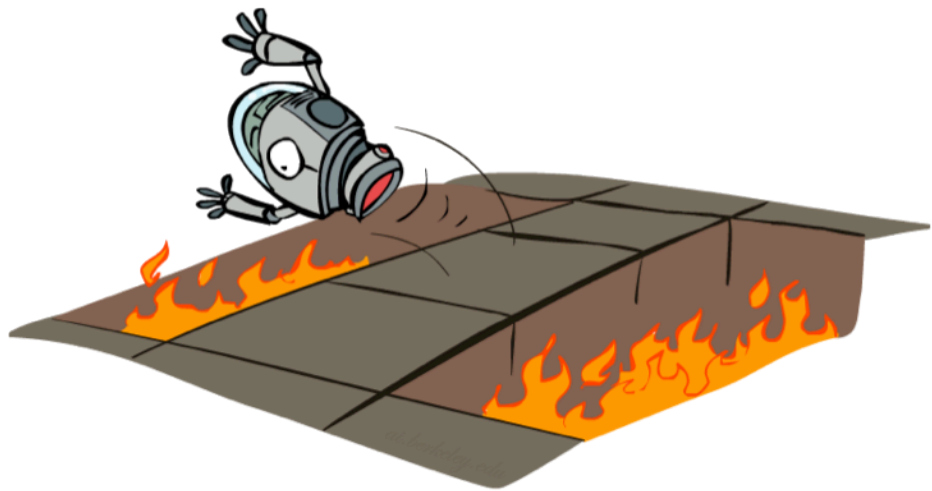
$$\pi(s) = \arg \max_a Q(s, a)$$
$$Q(s, a) = \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V(s')]$$

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



# Active Reinforcement Learning

---



---

# CSE 473: Introduction to Artificial Intelligence

Hanna Hajishirzi  
Reinforcement Learning

slides adapted from  
Dan Klein, Pieter Abbeel [ai.berkeley.edu](http://ai.berkeley.edu)  
And Dan Weld, Luke Zettlemoyer



# The Story So Far: MDPs and RL

## Known MDP: Offline Solution

$\mathcal{T}, R$

Goal

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

Evaluate a fixed policy  $\pi$

Technique

Value / policy iteration

Policy evaluation

## $\mathcal{T}, R$ Unknown MDP: Model-Based

Goal

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

Evaluate a fixed policy  $\pi$

Technique

VI/PI on approx. MDP

PE on approx. MDP

## Unknown MDP: Model-Free

Goal

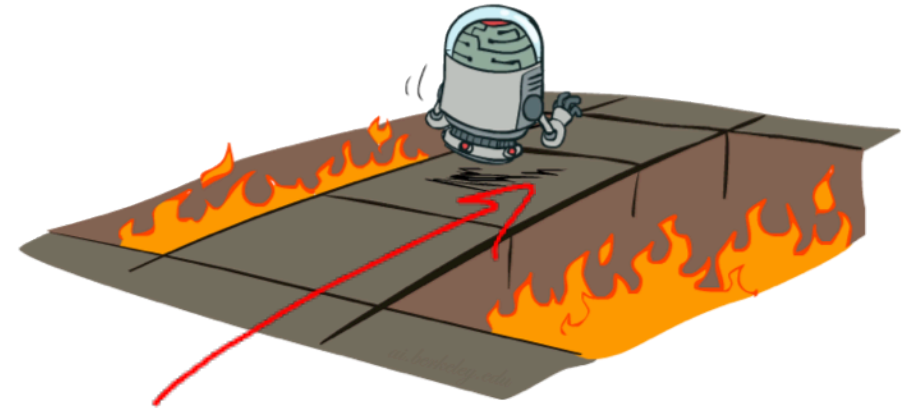
Evaluate a fixed policy  $\pi$

Technique

[direct evaluation]

# (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 annotations: A red arrow points from the label  $V_k$  on the left to the  $V_k$  term in the equation. A red arrow points from the  $\max_a$  term to the  $a$  variable. A red arrow points from the  $\sum_{s'}$  term to the  $s'$  variable. A red arrow points from the  $T(s, a, s')$  term to the  $s, a, s'$  variables. A red arrow points from the  $R(s, a, s')$  term to the  $s, a, s'$  variables. A red arrow points from the  $\gamma V_k(s')$  term to the  $V_k$  term.*

- But Q-values are more useful, so compute them instead

- Start with  $Q_0(s, a) = 0$ , which we know is right

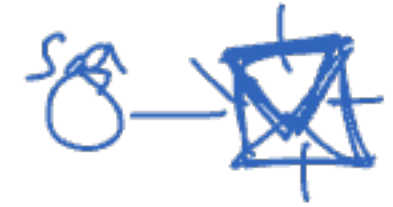
- Given  $Q_k$ , calculate the depth  $k+1$  values for all  $s, a$  states:

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

*Handwritten annotations: A red arrow points from the label  $Q_{k+1}(s, a)$  to the  $s, a$  variables. A red arrow points from the  $\sum_{s'}$  term to the  $s'$  variable. A red arrow points from the  $T(s, a, s')$  term to the  $s, a, s'$  variables. A red arrow points from the  $R(s, a, s')$  term to the  $s, a, s'$  variables. A red arrow points from the  $\max_{a'}$  term to the  $a'$  variable. A red arrow points from the  $Q_k(s', a')$  term to the  $s', a'$  variables.*



# 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

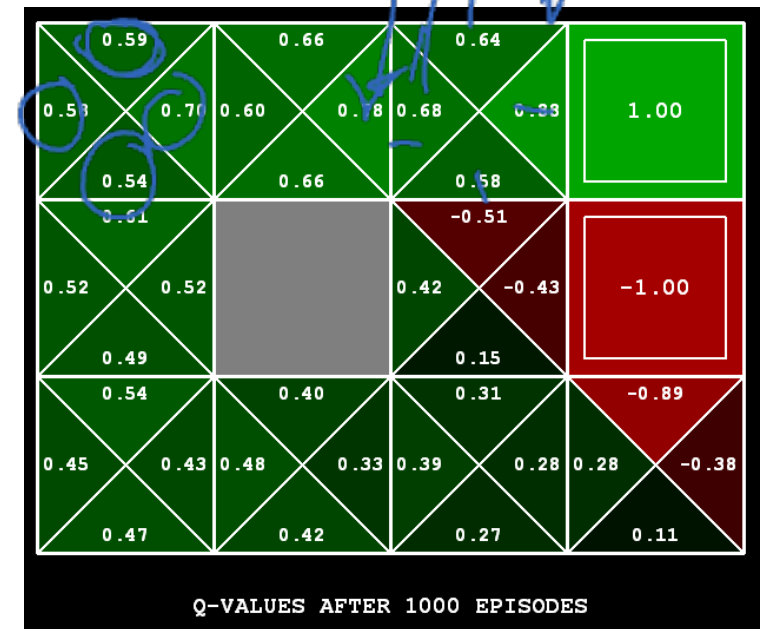
- Receive a sample  $(s, a, s', r)$
- Consider your old estimate:  $Q(s, a)$
- Consider your new sample estimate:

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

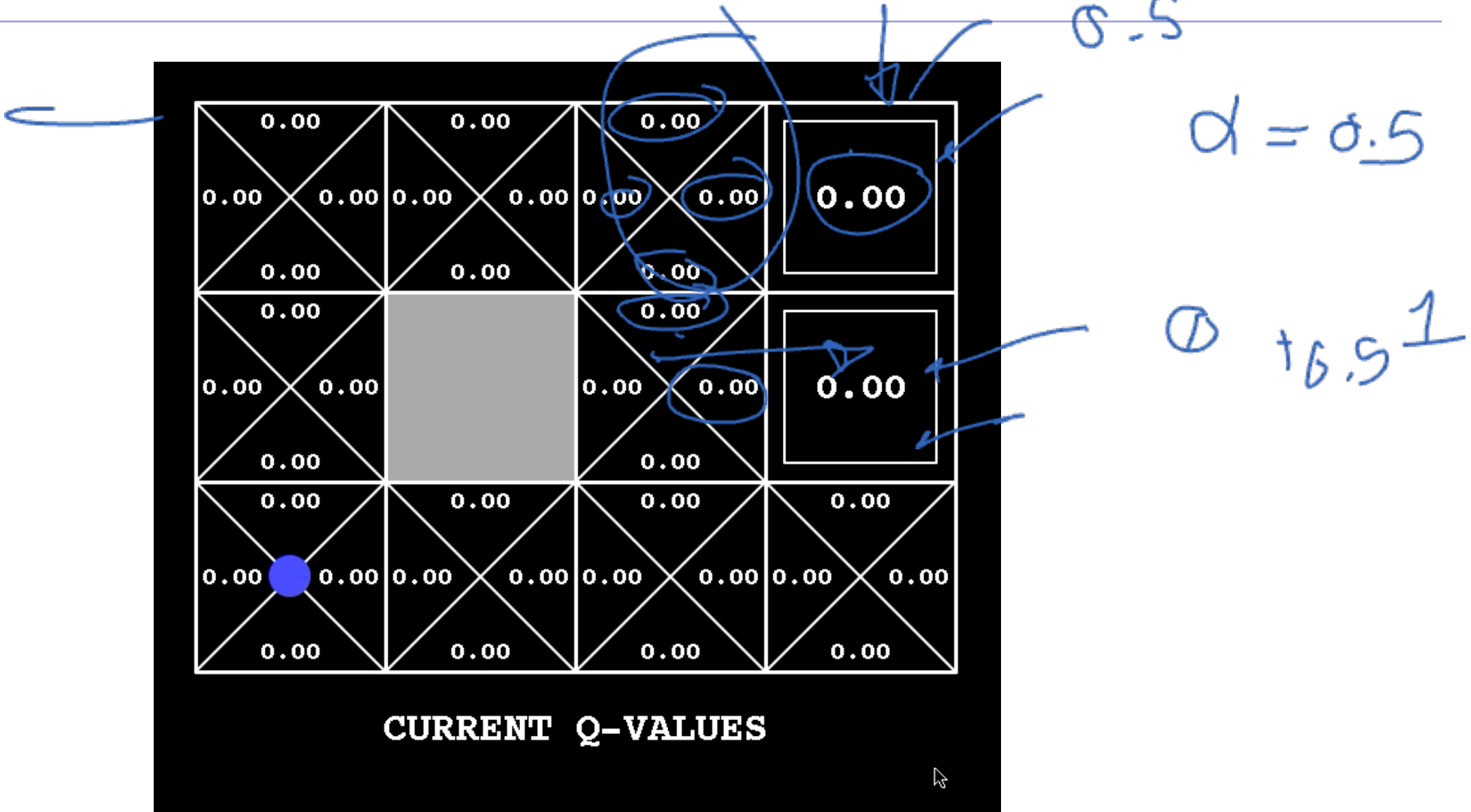
no longer policy evaluation!

- Incorporate the new estimate into a running average:

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



# Q-Learning Demo





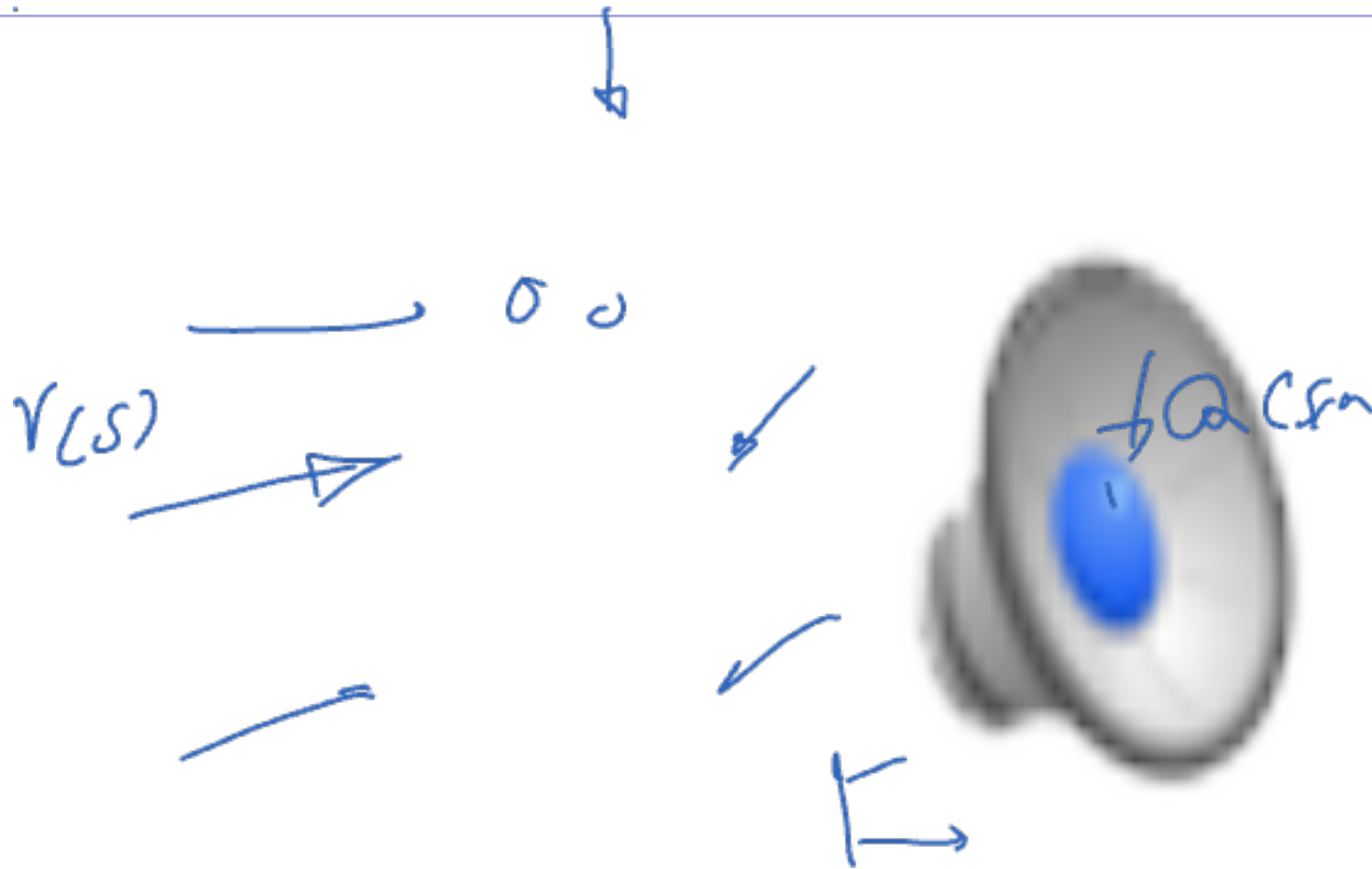
# Video of Demo Q-Learning -- Gridworld

---



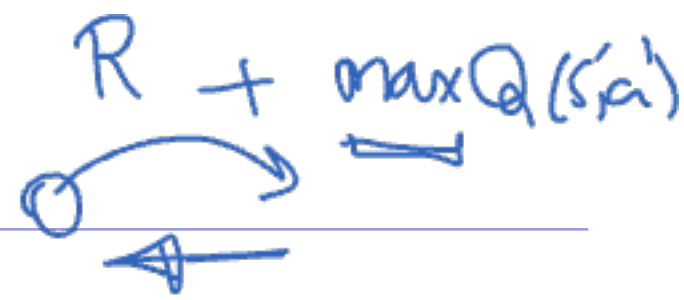
# Video of Demo Q-Learning -- Crawler

---



$\frac{\alpha}{\alpha+1}$

# Q-Learning Properties

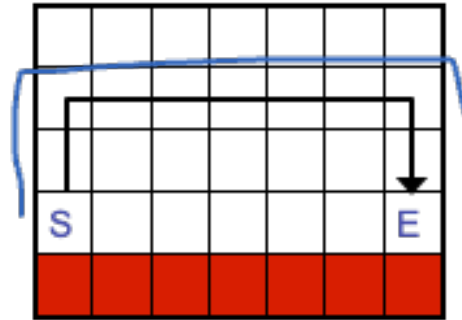
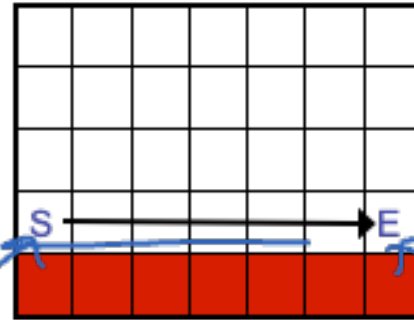


- Amazing result: Q-learning converges to optimal policy -- even if you're acting suboptimally!

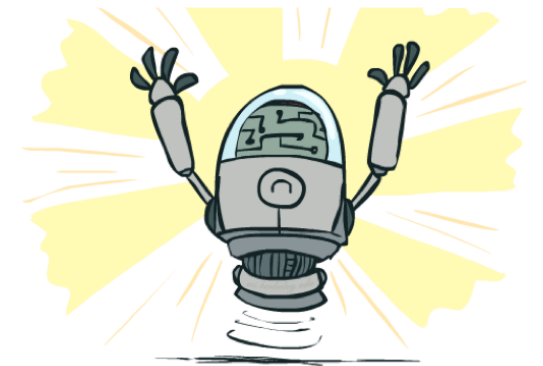
- This is called **off-policy learning**

- 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 (!)



$(1-\alpha)Q + \alpha \text{ Sample}$



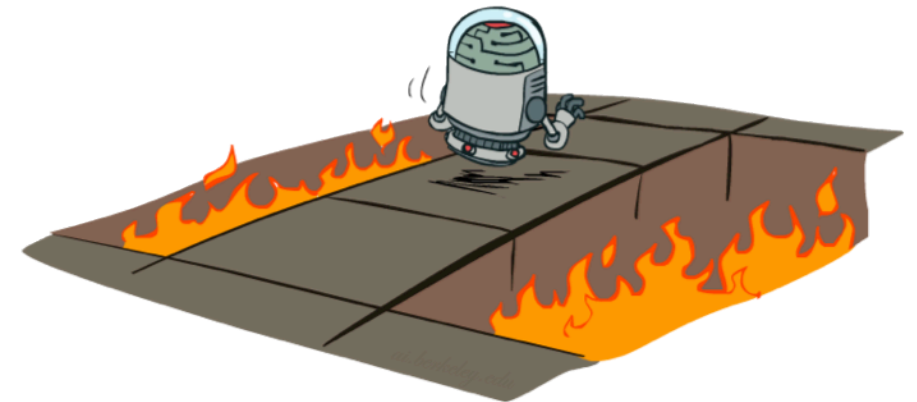
# Discussion: Model-Based vs Model-Free RL

$T, R \rightarrow V^*$   $V^*, \pi^*$

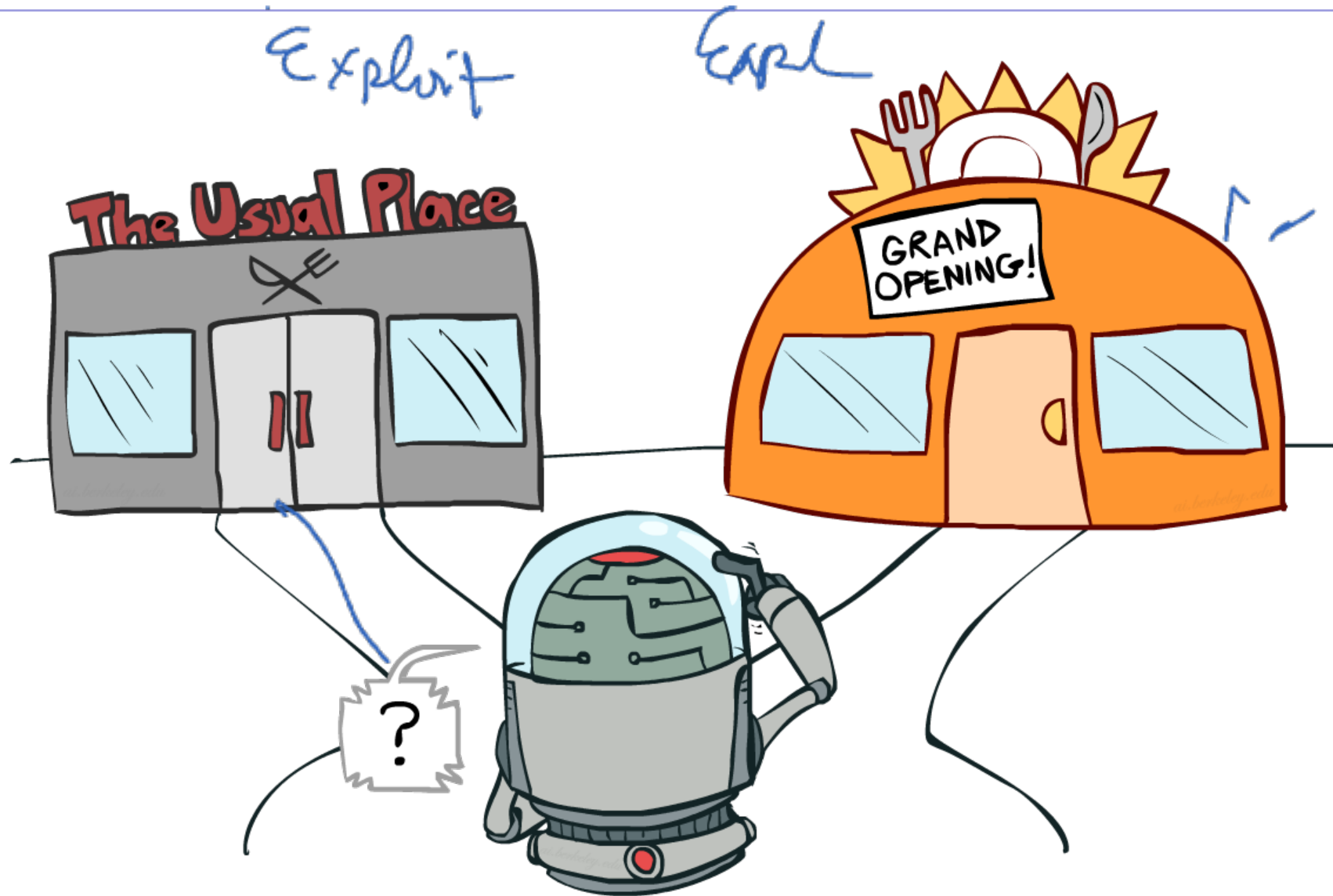
- Model-Based vs. Model Free

- Active vs. Passive

- act according to current optimal (based on Q-Values)
- but also explore...



# Exploration vs. Exploitation



$\epsilon$

# How to Explore?

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



# 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.



$$f(u, n) = u + \frac{k}{n+1}$$

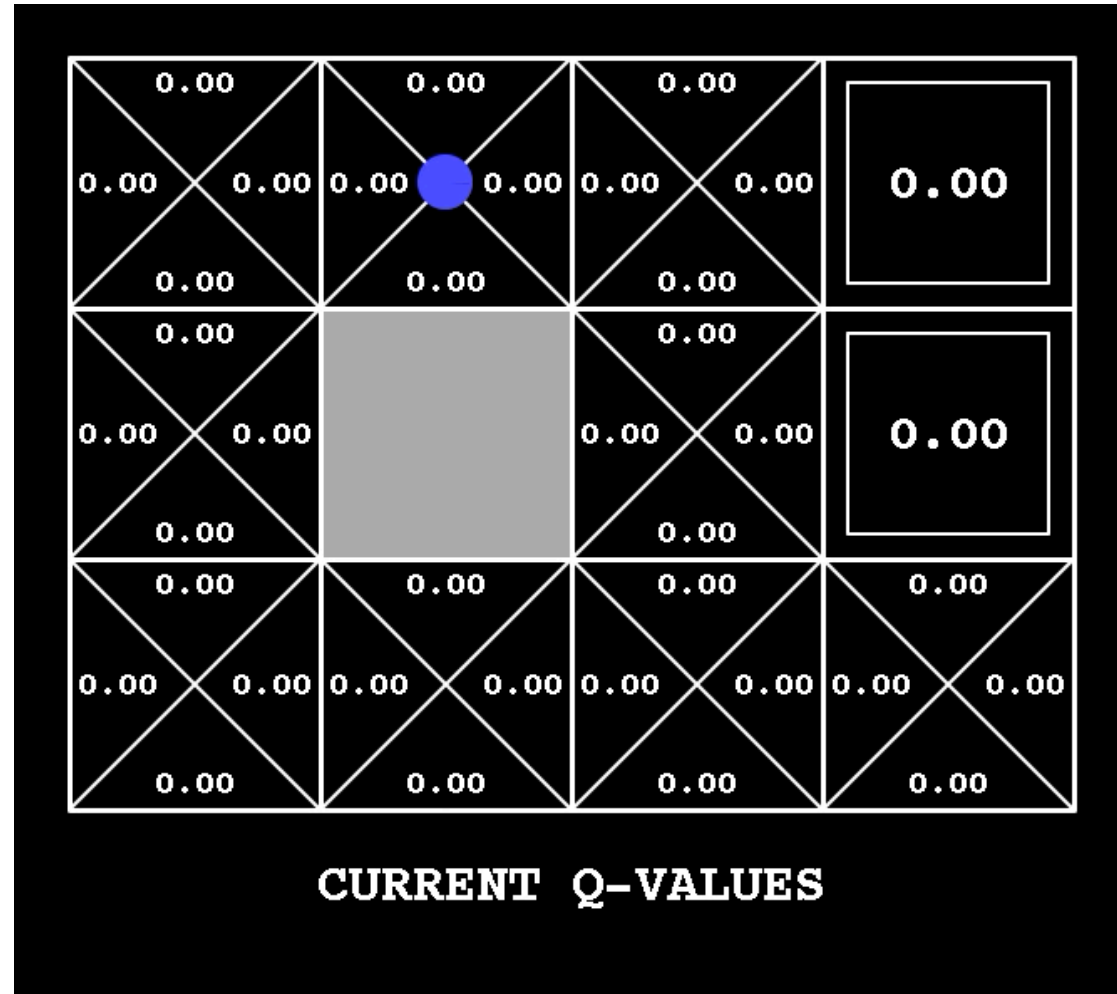
$u + k$   
 $u + \frac{k}{15}$

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

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

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

# Q-Learn Epsilon Greedy



$Q(s,a)$

Sample: 1 ✓

$$Q(s,a) = 0.5$$

Sample: 1

$$0.5 \times 0.5 + 0.5 \times 1$$

0.75



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

---

0,



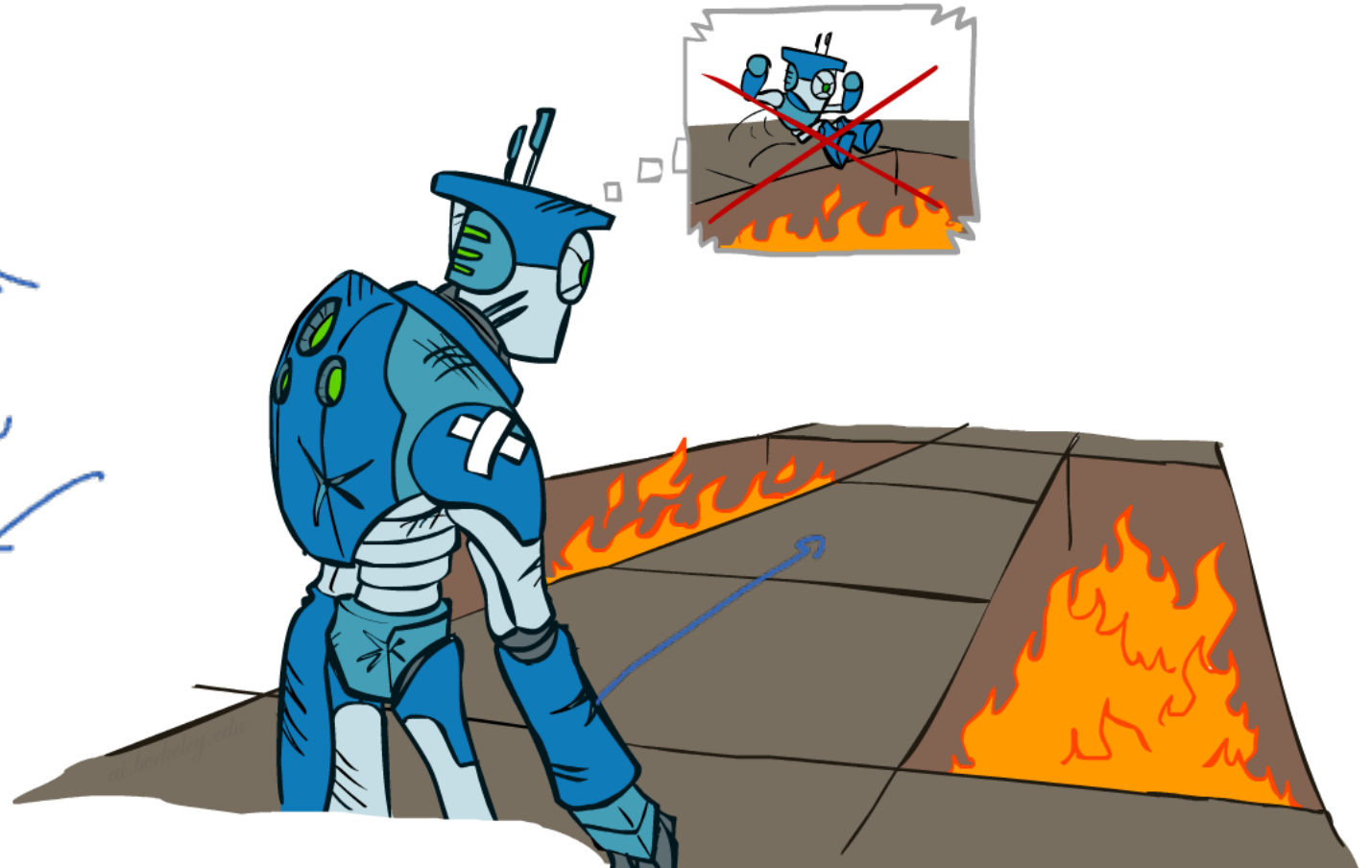
# Video of Demo Q-learning – Exploration Function – Crawler

---



# 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 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



# Recap: Q-Learning

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

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

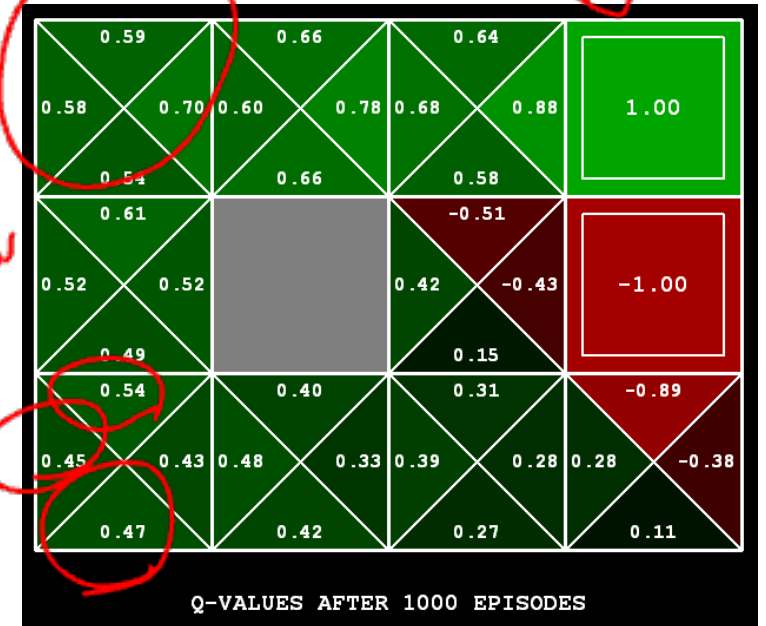
- Learn  $Q(s, a)$  values as you go

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

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

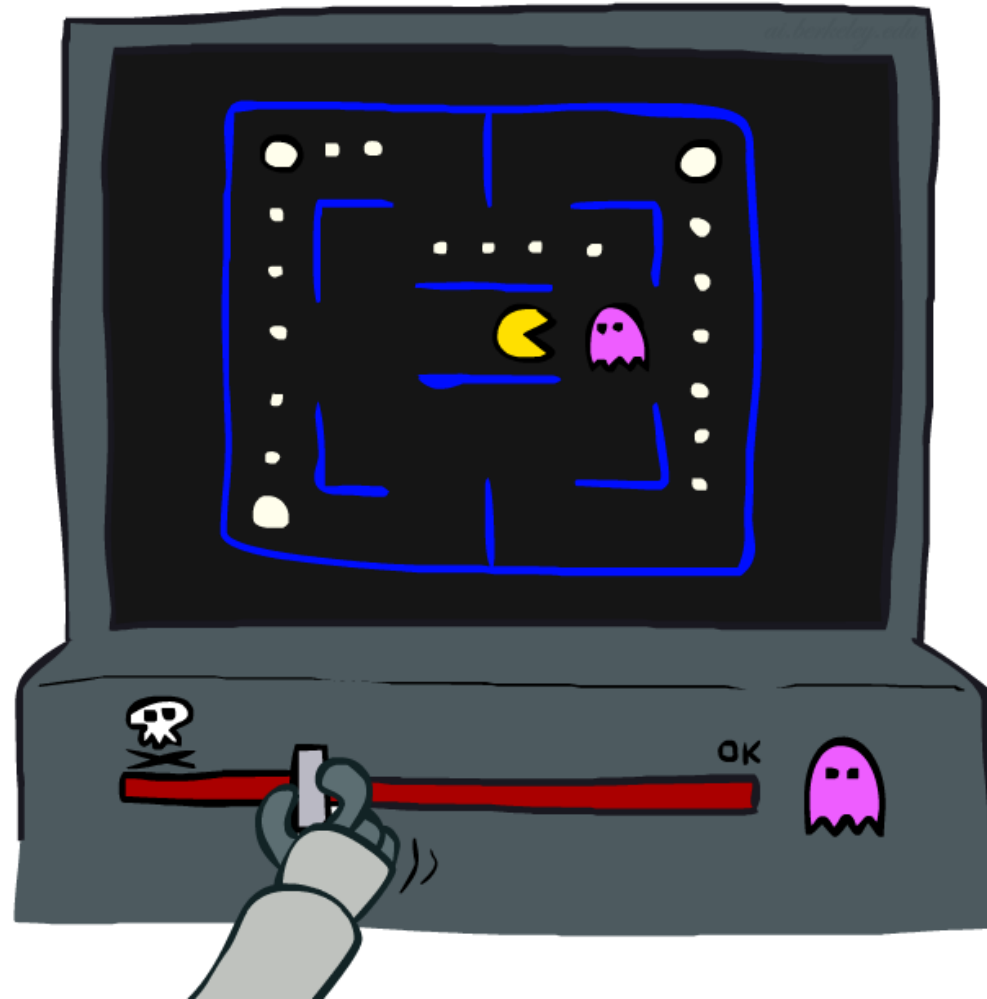
- Incorporate the new estimate into a running average:

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



# Approximate Q-Learning

---

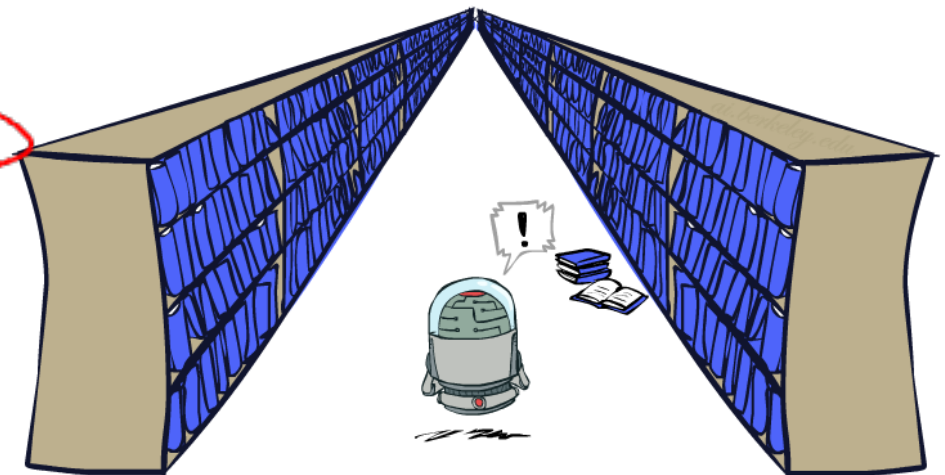
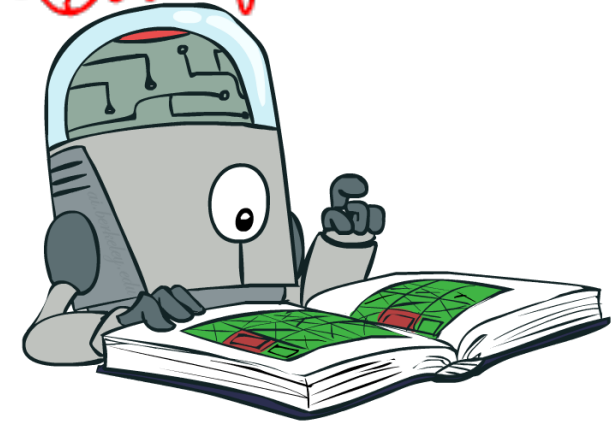


# Generalizing Across States

*tabular form*

- 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

*30  
2*



# Video of Demo Q-Learning Pacman – Tiny – Watch All

---



# Video of Demo Q-Learning Pacman – Tiny – Silent Train

---





# Video of Demo Q-Learning Pacman – Tricky – Watch All

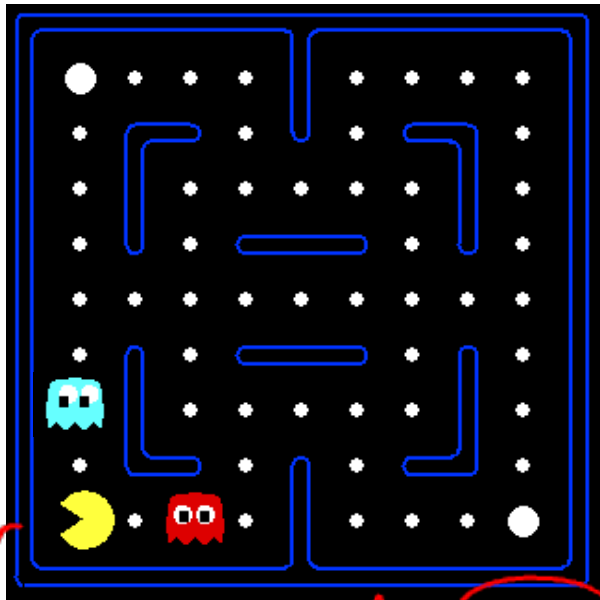
---



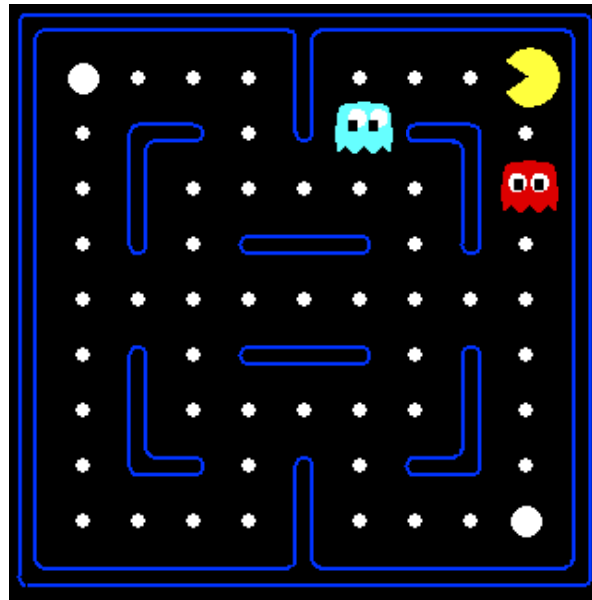
# Example: Pacman

$Q_{\pi}(s_2, a)$

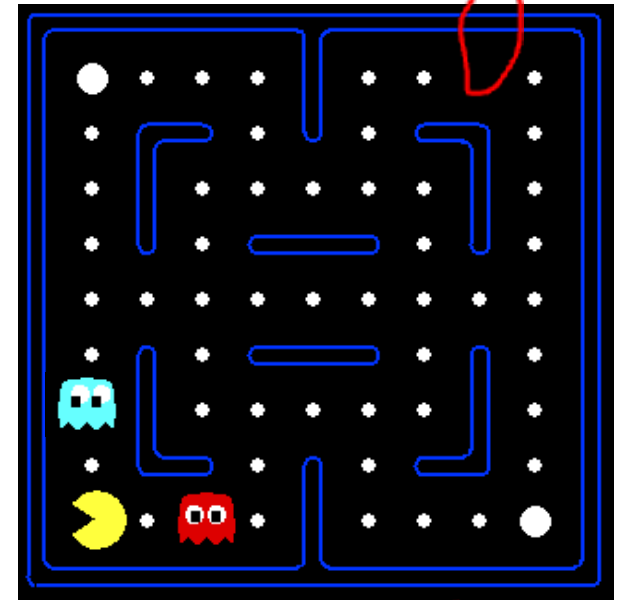
Let's say we discover through experience that this state is bad:



In naïve q-learning, we know nothing about this state:

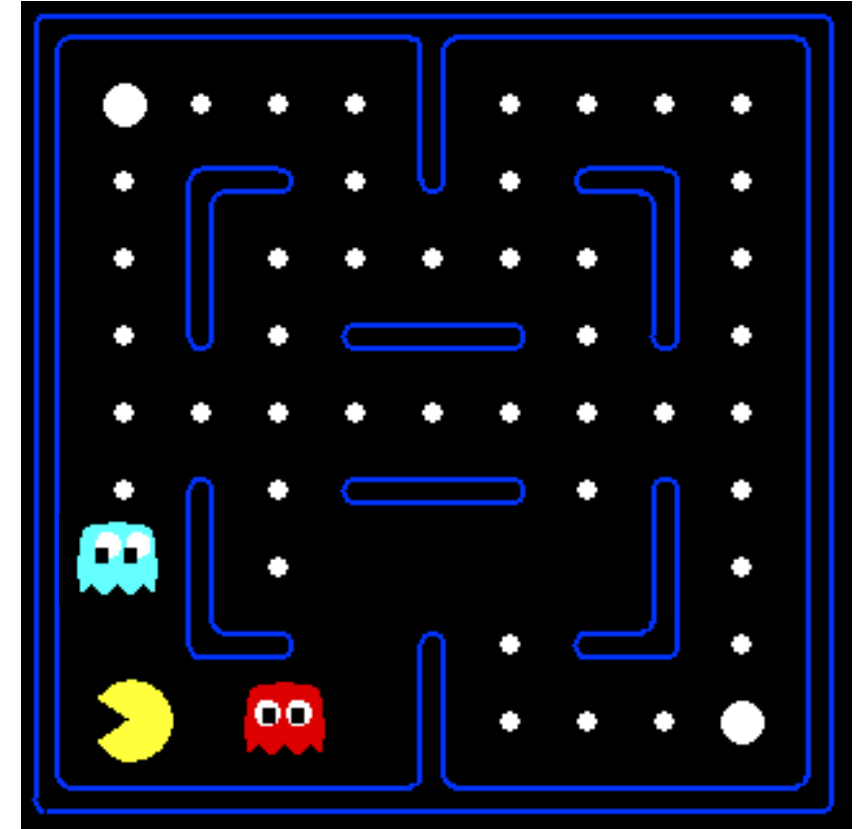


Or even this one!



# 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
    - $1 / (\text{dist to dot})^2$
    - 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:

$$V(s) = w_1 f_1(s) + w_2 f_2(s) + \dots + w_n f_n(s)$$

$$Q(s, a) = w_1 f_1(s, a) + w_2 f_2(s, a) + \dots + w_n f_n(s, a)$$

- Advantage: our experience is summed up in a few powerful numbers
- Disadvantage: states may share features but actually be very different in value!

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

$Q(s, a) \leftarrow Q(s, a) + \alpha [\text{difference}]$

$w_i \leftarrow w_i + \alpha [\text{difference}] f_i(s, a)$

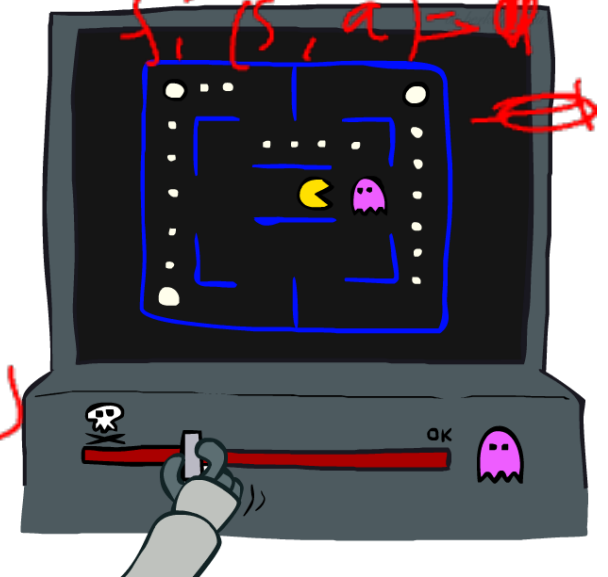
- Intuitive interpretation:

- 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

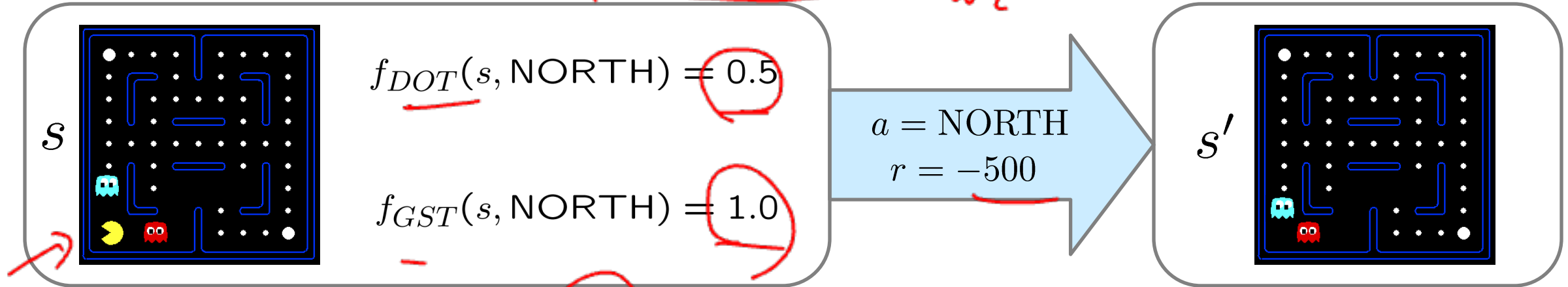
Exact Q's

Approximate Q's



# Example: Q-Pacman

~~$$Q(s, a) = 4.0 f_{DOT}(s, a) - 1.0 f_{GST}(s, a)$$~~



~~$$Q(s, \text{NORTH}) = +1$$

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

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

difference =  $-501$  →

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

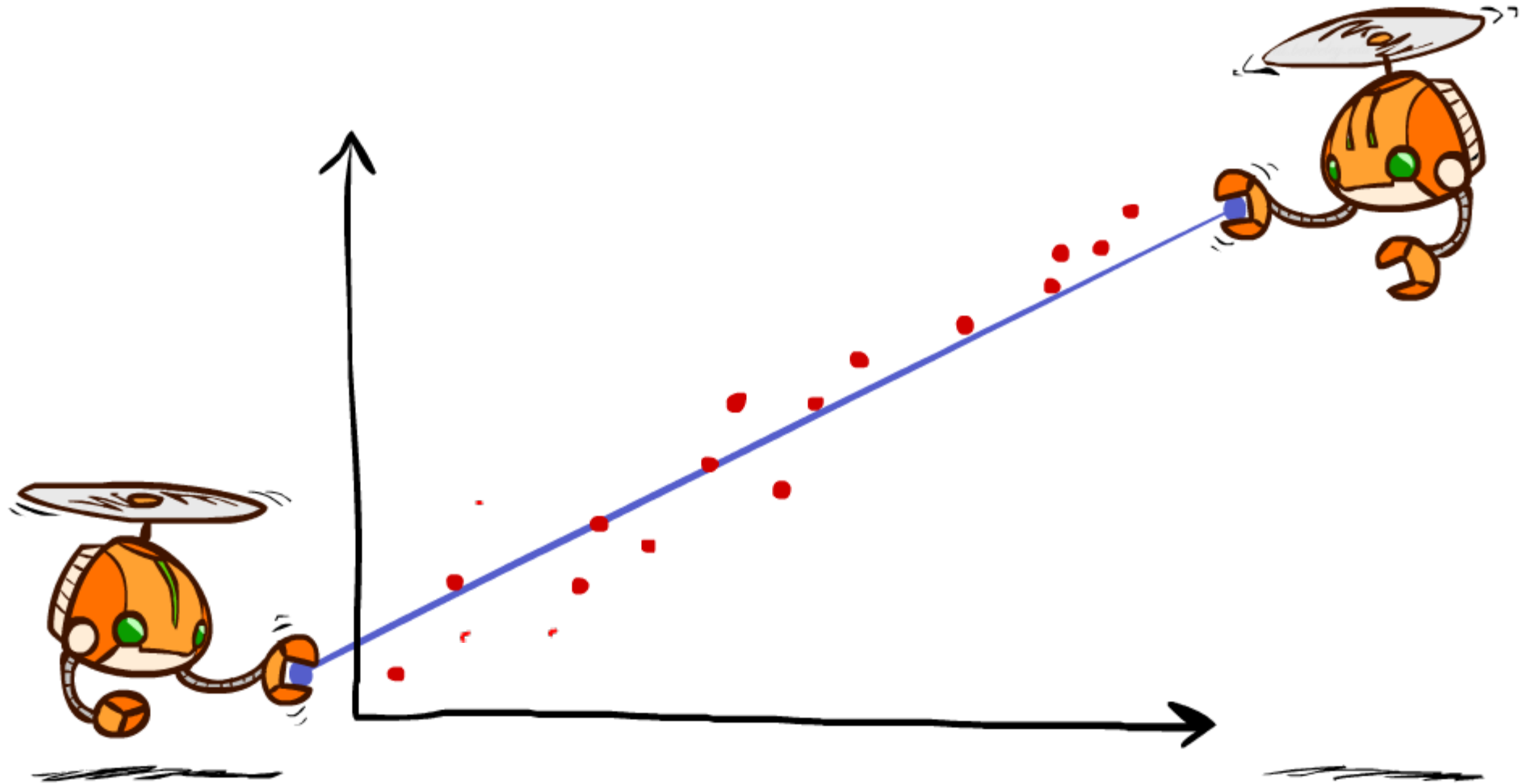
$$w_{GST} \leftarrow -1.0 + \alpha [-501] 1.0$$

# Video of Demo Approximate Q-Learning -- Pacman

---

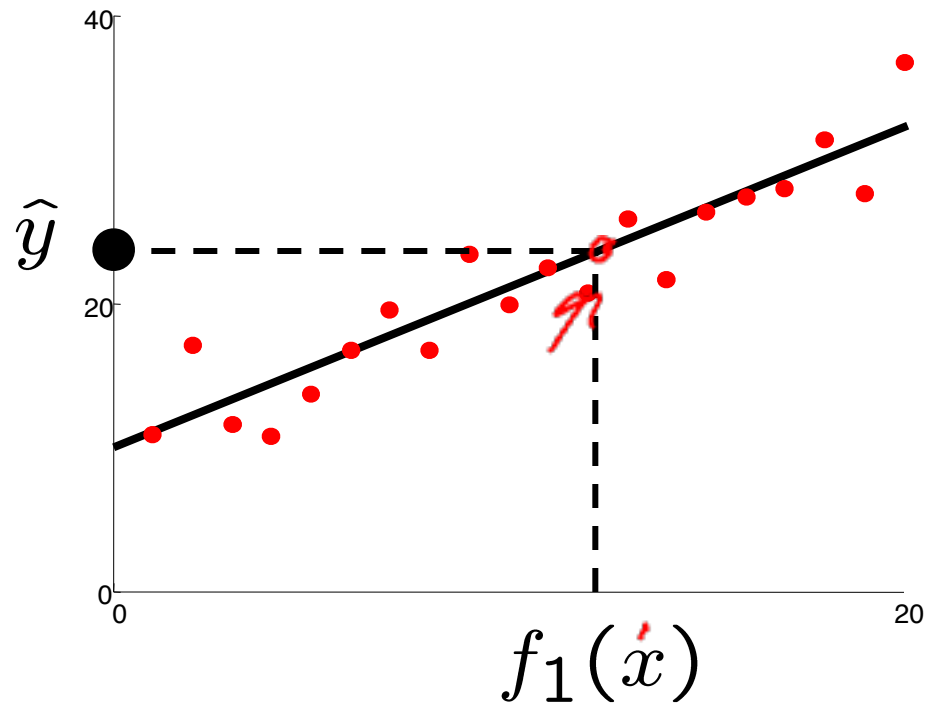


# Q-Learning and Least Squares



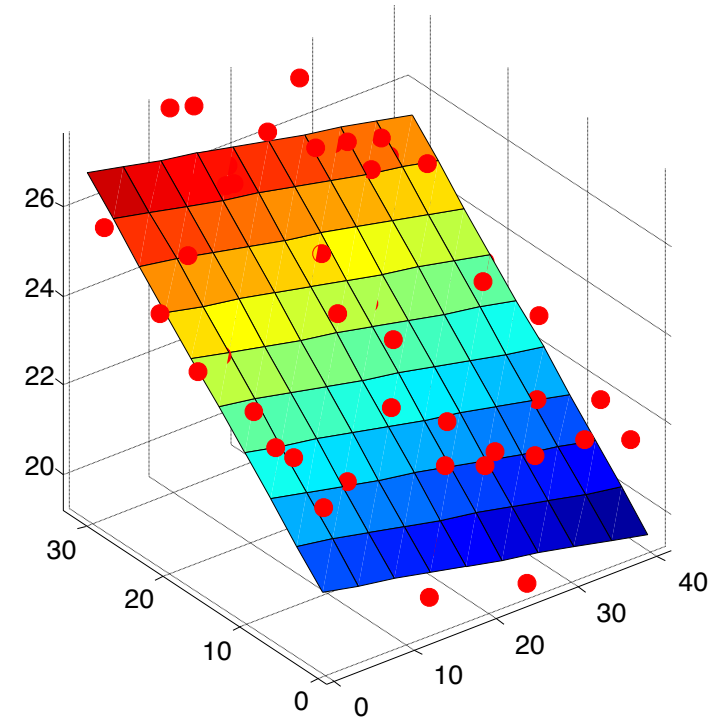


# Linear Approximation: Regression



Prediction:

$$\hat{y} = w_0 + w_1 f_1(x)$$



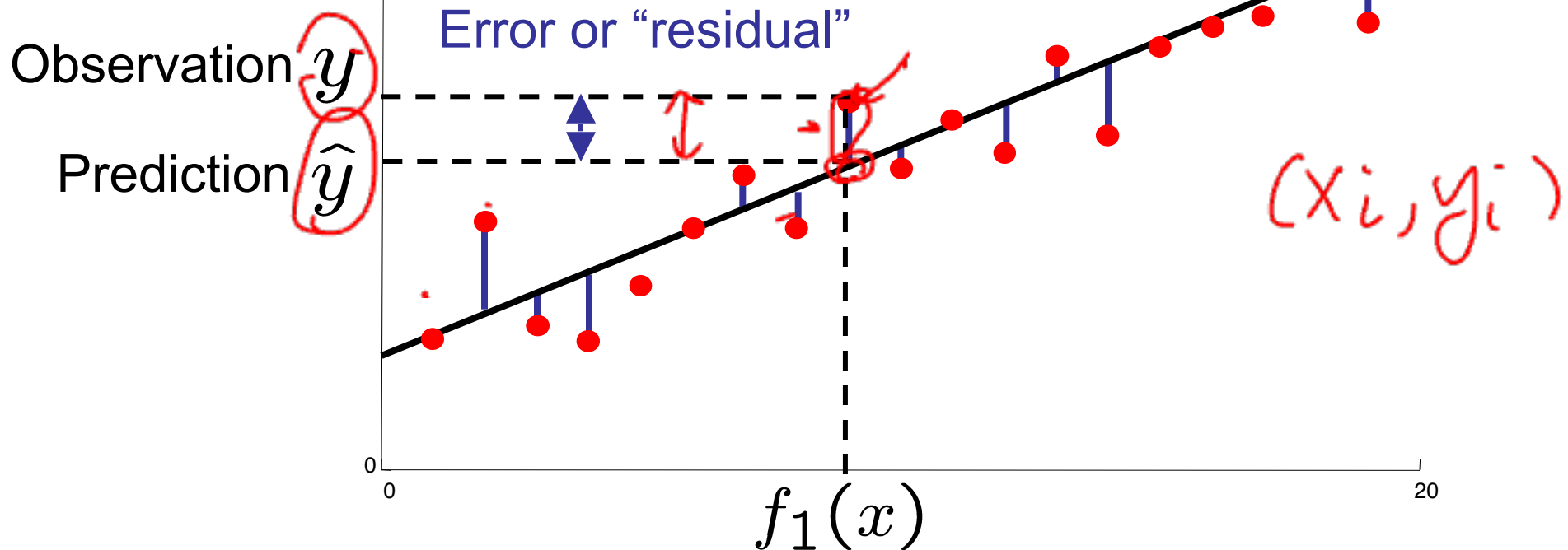
Prediction:

$$\hat{y}_i = w_0 + w_1 f_1(x) + w_2 f_2(x)$$

# Optimization: Least Squares

total error =

$$\sum_i (y_i - \hat{y}_i)^2 = \sum_i \left( y_i - \sum_k w_k f_k(x_i) \right)^2$$



$w_1$   $(w_m)$   $w_m$

# Minimizing Error



Imagine we had only one point  $x$ , with features  $f(x)$ , target value  $y$ , and weights  $w$ :

★  $\text{error}(w) = \frac{1}{2} \left( y - \sum_k w_k f_k(x) \right)^2$

$\frac{\partial \text{error}(w)}{\partial w_m} = - \left( y - \sum_k w_k f_k(x) \right) f_m(x)$

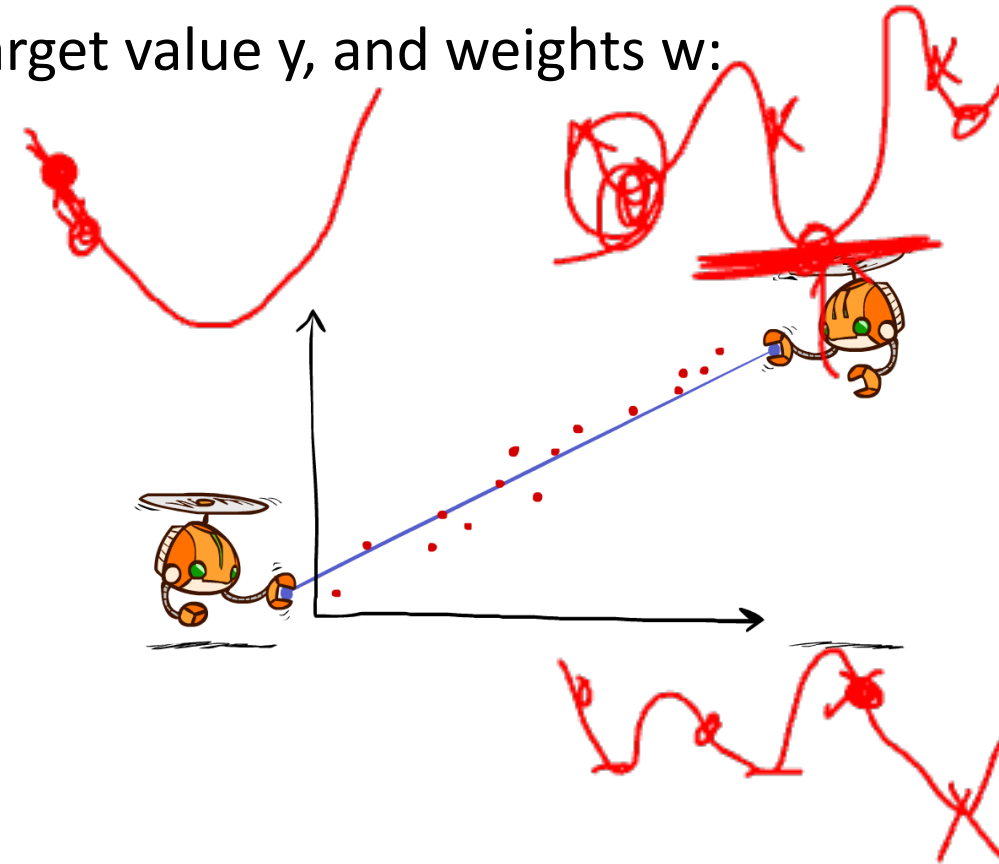
★  $w_m \leftarrow w_m + \alpha \left( y - \sum_k w_k f_k(x) \right) f_m(x)$

Approximate q update explained:

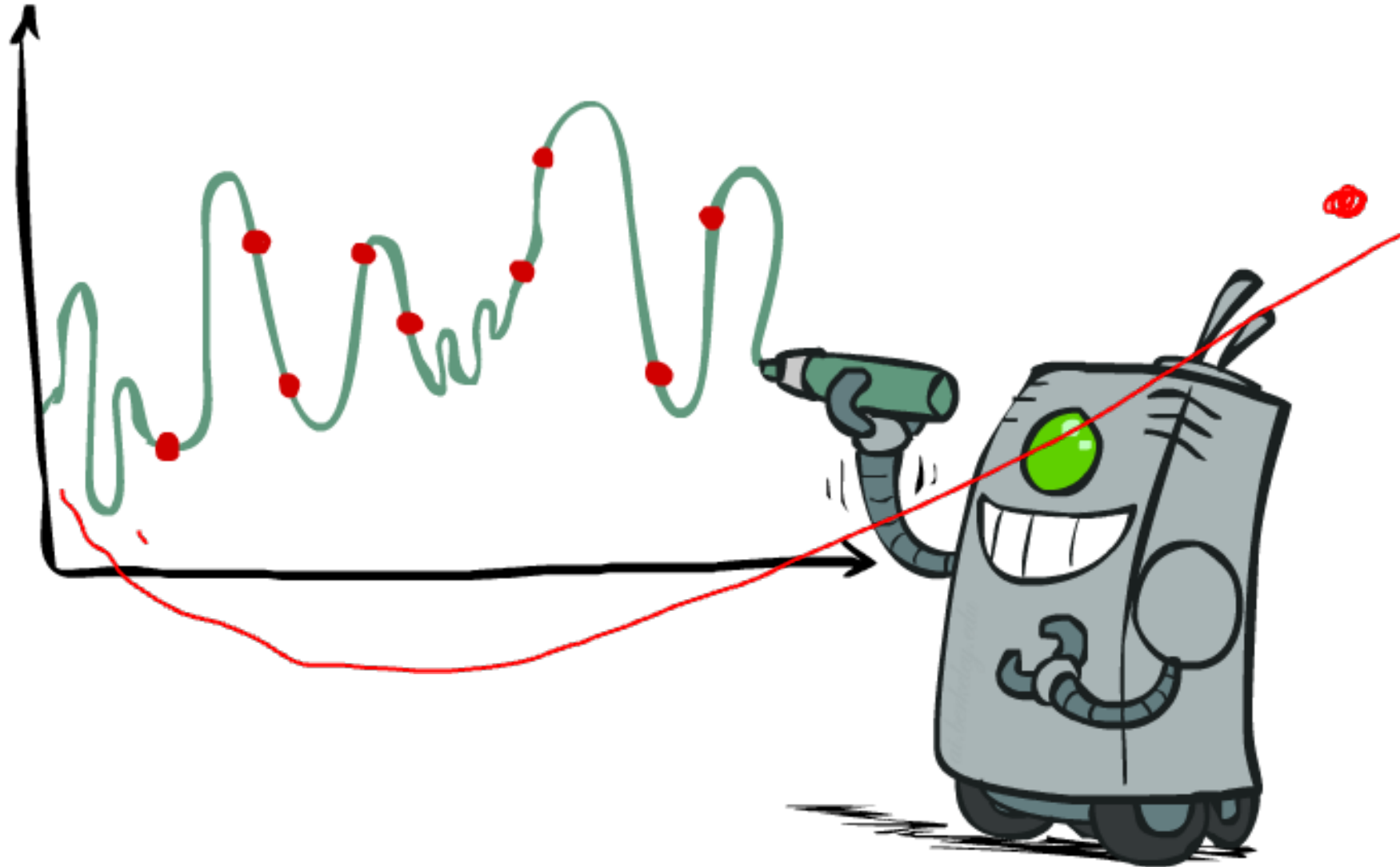
★  $w_m \leftarrow w_m + \alpha \left[ r + \gamma \max_a Q(s', a') - Q(s, a) \right] f_m(s, a)$

“target”

“prediction”



# Overfitting: Why Limiting Capacity Can Help



# Summary: MDPs and RL

## Known MDP: Offline Solution ✓

Goal	Technique
Compute $V^*, Q^*, \pi^*$	Value / policy iteration
Evaluate a fixed policy $\pi$	Policy evaluation

## Unknown MDP: Model-Based <sup>T/IR</sup>

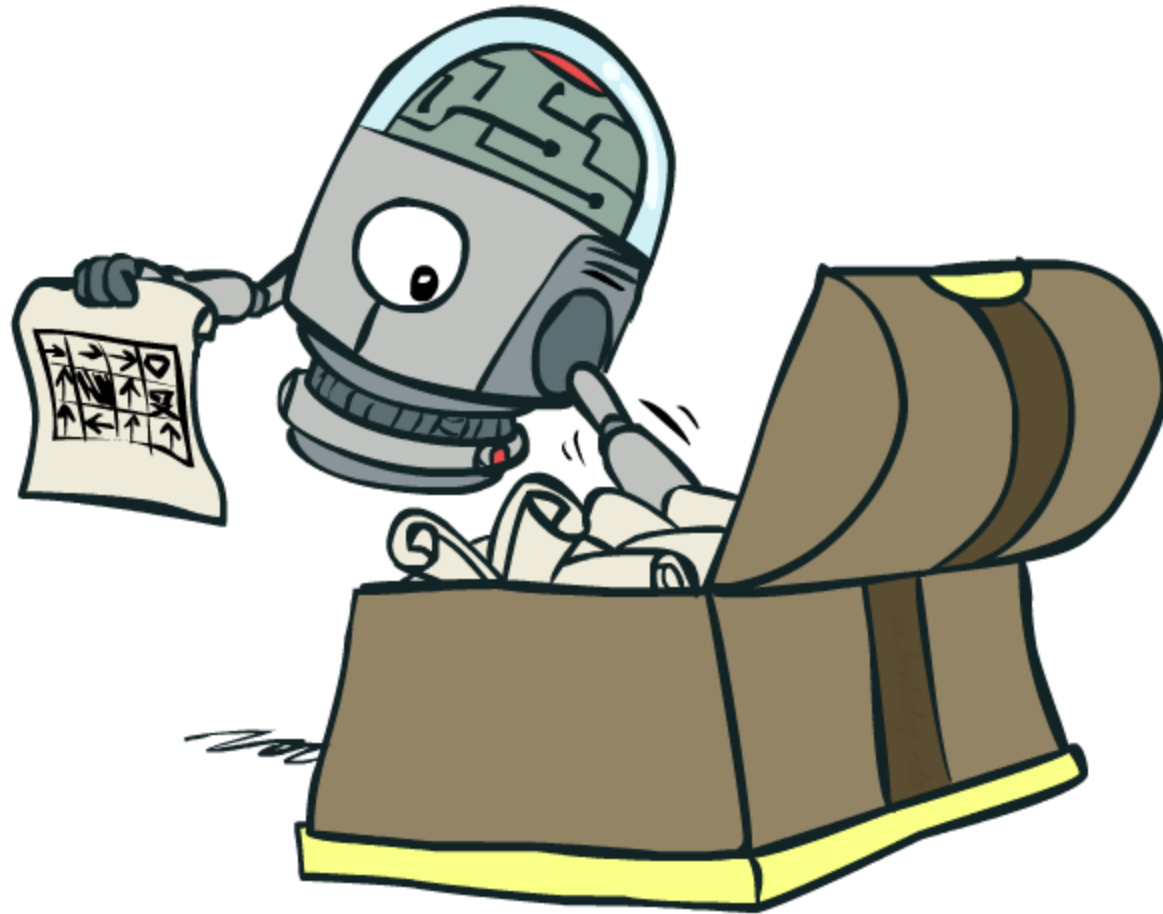
Goal	<i>*use features to generalize</i>	Technique
Compute $V^*, Q^*, \pi^*$		VI/PI on approx. MDP
Evaluate a fixed policy $\pi$		PE on approx. MDP

## Unknown MDP: Model-Free <sup>+</sup>

Goal	<i>*use features to generalize</i>	Technique
Compute $V^*, Q^*, \pi^*$		<u>Q-learning</u>
Evaluate a fixed policy $\pi$		<u>Value Learning</u>

# Policy Search

---



# Policy Search

---

- Problem: often the feature-based policies that work well (win games, maximize utilities) aren't the ones that approximate  $V / Q$  best
  - E.g. your value functions from project 2 were probably horrible estimates of future rewards, but they still produced good decisions
  - Q-learning's priority: get Q-values close (modeling)
  - Action selection priority: get ordering of Q-values right (prediction)
- Solution: learn policies that maximize rewards, not the values that predict them
- Policy search: start with an ok solution (e.g. Q-learning) then fine-tune by nudging each feature weight up and down and see if your policy is better than before



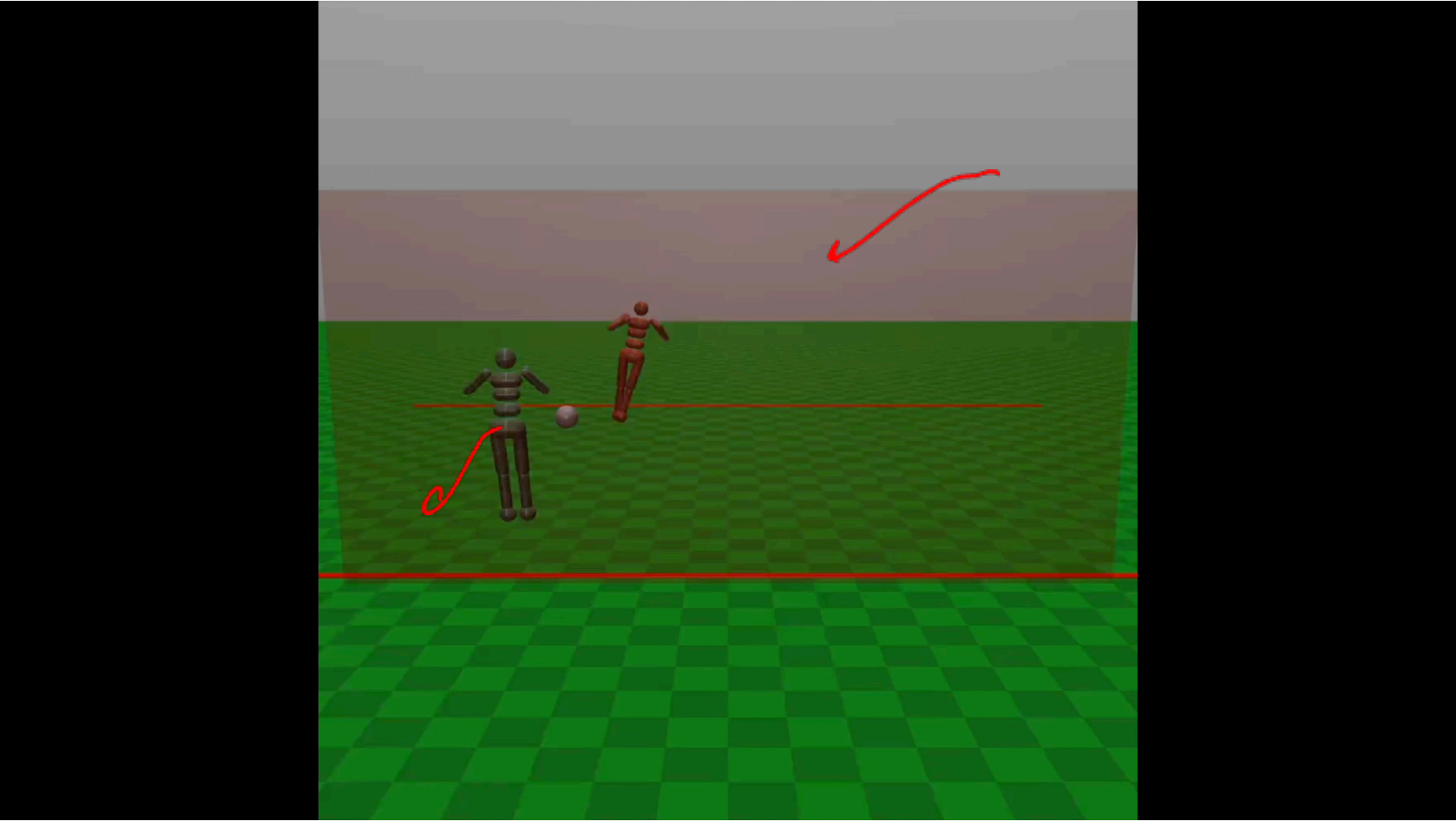
# New in Model-Free RL

## Playing Atari Games

---







# ( $x_i, y_i$ ) Conclusion

- We've seen how AI methods can solve problems in:
  - Search
  - Games
  - Markov Decision Problems
  - Reinforcement Learning
- Next up: Uncertainty and Learning!

