

Search with Uncertainty



Reinforcement Learning

Review: Solving MDP

MDP

Markov

$$T(s, a \rightarrow s') \Rightarrow P(s'|s, a) : \text{all known}$$

$$R(s) \rightarrow \text{all known}$$

$\{s\}$ - fully obs

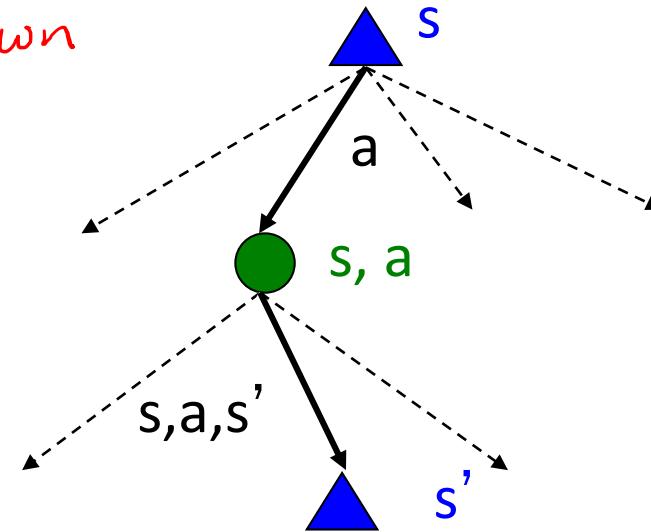
POMDP

finite

discount γ

Policy $\pi > \underline{P(a|s)}$: Seq of actions given s

$$\text{Return } G_t = \sum_t \gamma^t R_t$$



$$\text{Value of } s \quad V^\pi(s) = \mathbb{E}_\pi [G_t | s]$$

$$\text{action value ft} \quad Q^\pi(s, a) = \mathbb{E}_\pi [G_t | s, a]$$

Review: Solving MDP

↳ optimal policy $\pi^*(\text{set of action}) \rightarrow \text{maximize return}$

for each
 $s \rightarrow a^*$

$$V^*(s) = \max_{\pi} V^{\pi}(s) = \max_a \tilde{g}^*(s, a)$$

$s \xrightarrow{a} s' \xrightarrow{a'} s''$

$$\pi^*(s) = \operatorname{argmax}_a \sum_{s'} P(s'|s, a) \underline{V}_{\pi}(s')$$

Bellman Eqn: $V(s) \leftrightarrow V(s')$

$$\underline{V}^{\pi}(s) = R(s) + \gamma \sum_{s'} P(s'|s, \pi(s)) \cdot \underset{\text{future}}{\downarrow} V^{\pi}(s')$$

Value Iter Policy Iter

Review: Solving MDP

Policy Iter

$$V^{\pi} \rightarrow 0$$

$\pi \rightarrow \text{random}$

1) eval V^{π}

$$V_i^{\pi}(s) = R(s) + \gamma \sum_{s'} P(s'|s, \pi_i(s)) V_i^{\pi}(s')$$

2) Policy update

$$\max_a \sum_{s'} P(s'|s, a) V_i(s')$$

$$\pi_{i+1} \leftarrow a$$

$$Q(s, a) \rightarrow V$$

Value Iter

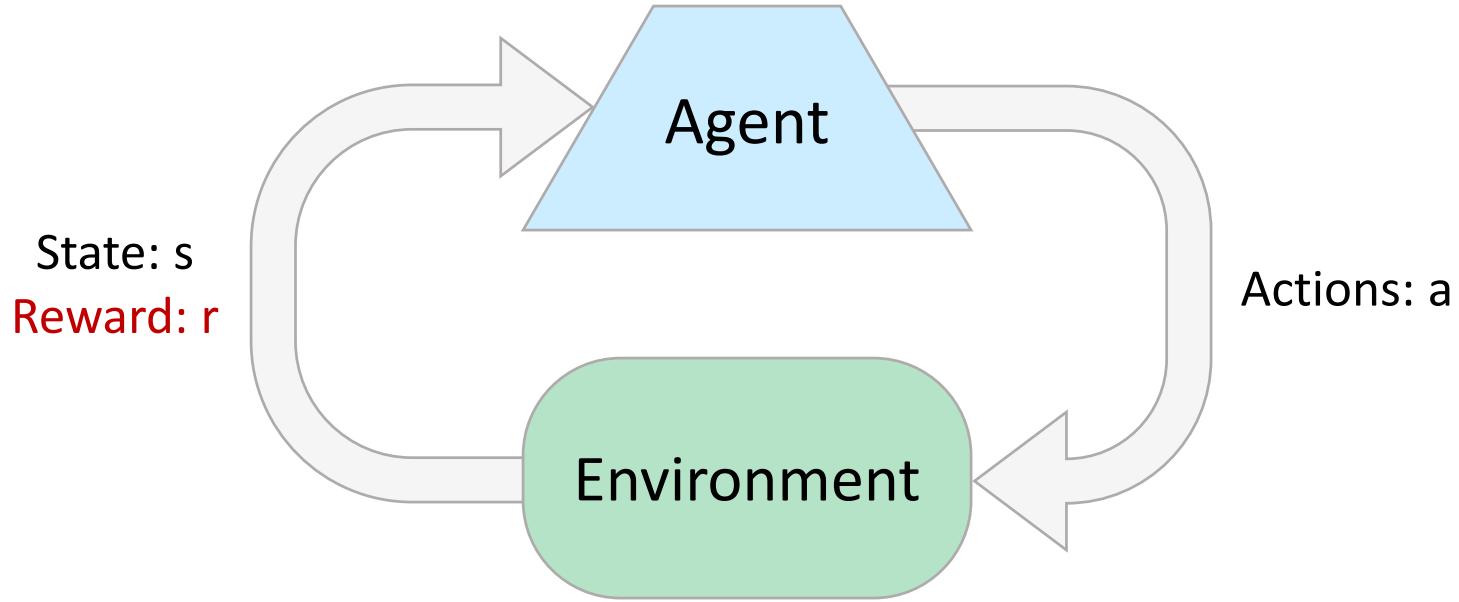
$$V^{\pi}(s) \rightarrow 0$$

$$V_{i+1}(s) \leftarrow$$

$$R(s) + \gamma \max_a \left(\sum_{s'} P(s'|s, a) V_i(s') \right)$$

$$\pi^*, \underset{a}{\operatorname{argmax}} \left(\sum_{s'} P(s'|s, a) V_i(s') \right)$$

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!

Cheetah



OpenAI,

Atari



Example: Learning to Walk



Before training

Example: Learning to Walk

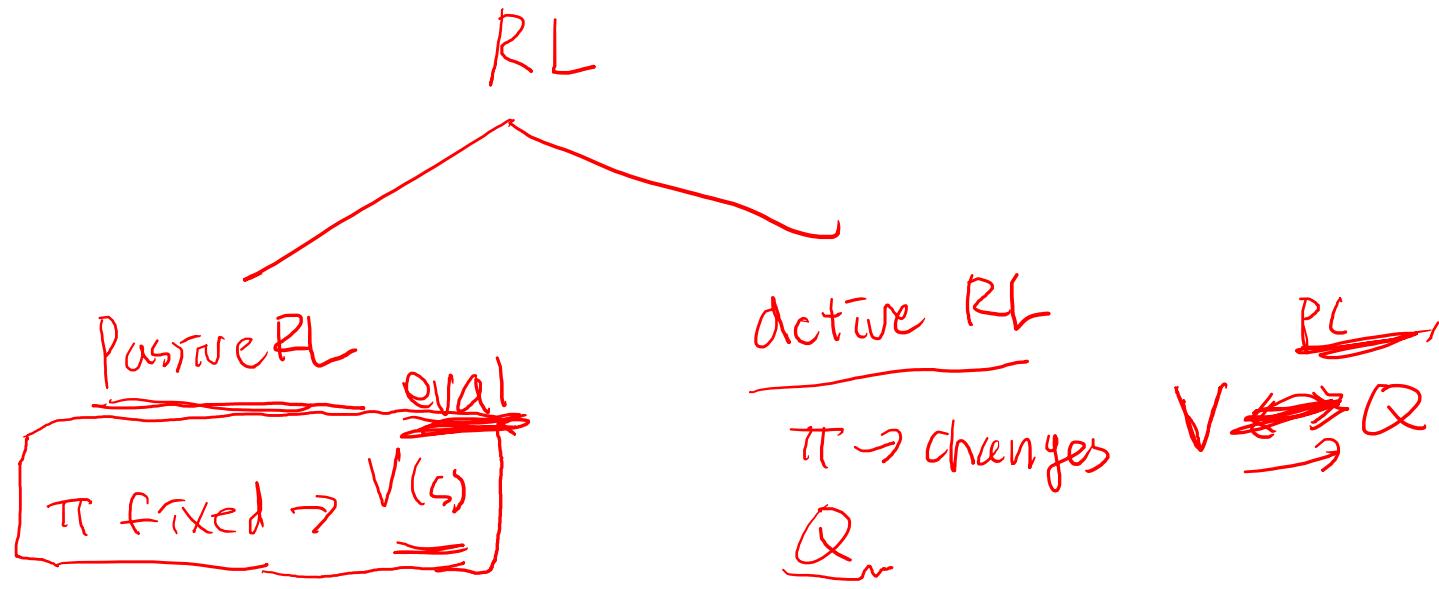
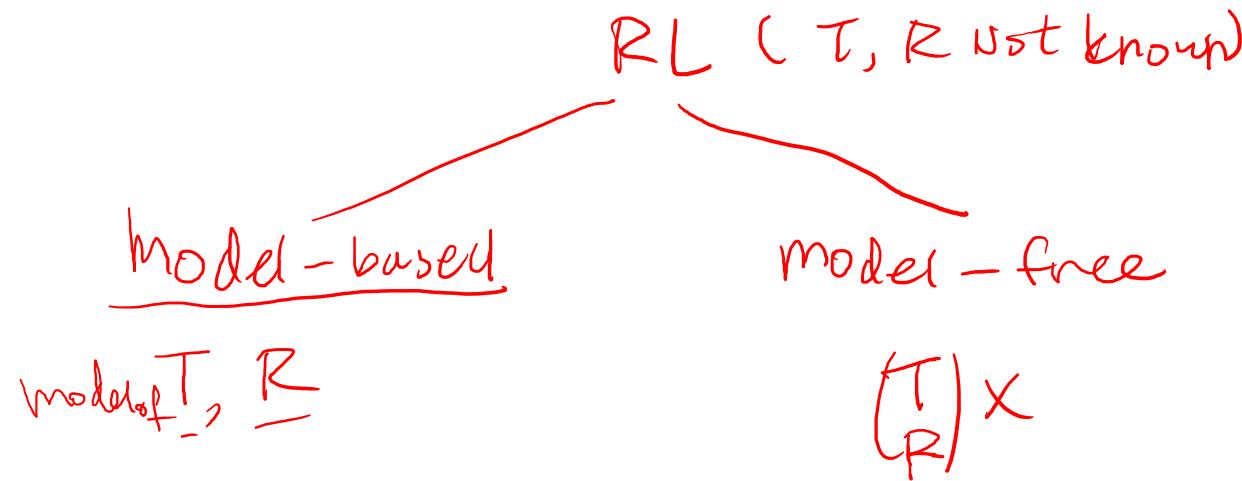


After training

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 out actions and states to learn

Reinforcement Learning approaches



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

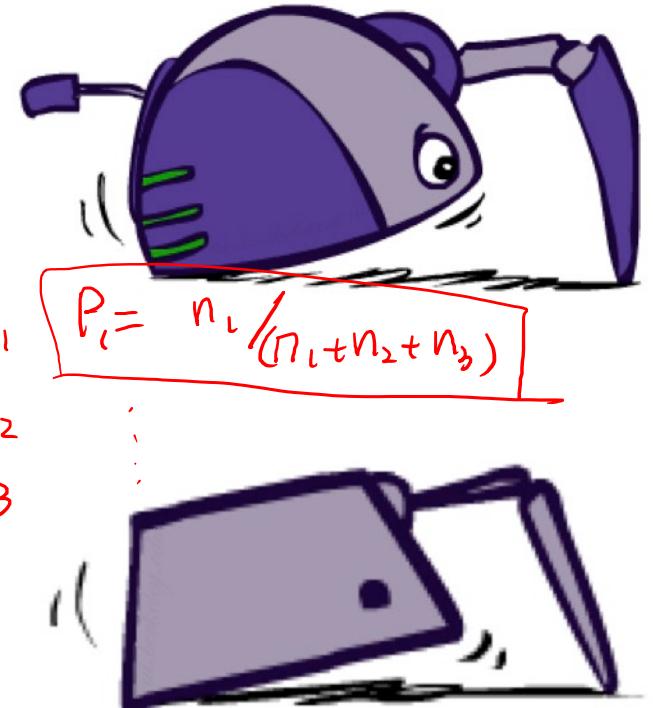
$$\begin{array}{c} \xrightarrow{s, a} s' \\ \hat{T}(s, a, s') \\ \hat{R}(s) \end{array}$$

$$\begin{array}{c} s \rightarrow s'_1 \\ \downarrow \\ s'_2 \\ \downarrow \\ s'_3 \end{array}$$

$$P_i = \frac{n_i}{(n_1 + n_2 + n_3)}$$

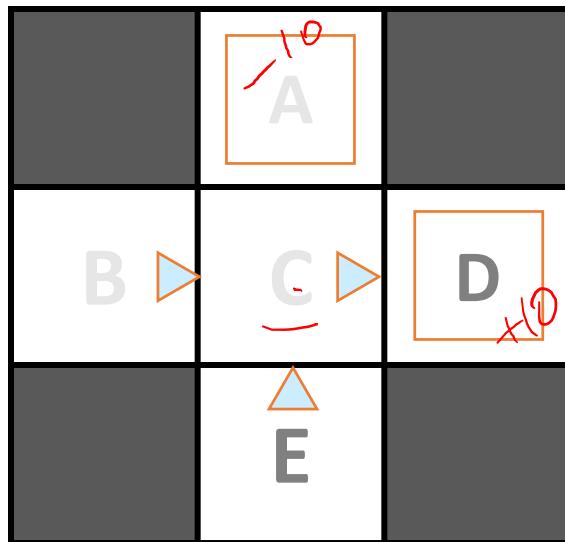
n_1
 n_2
 n_3

$$\underline{\hat{T}, \hat{R}}$$



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

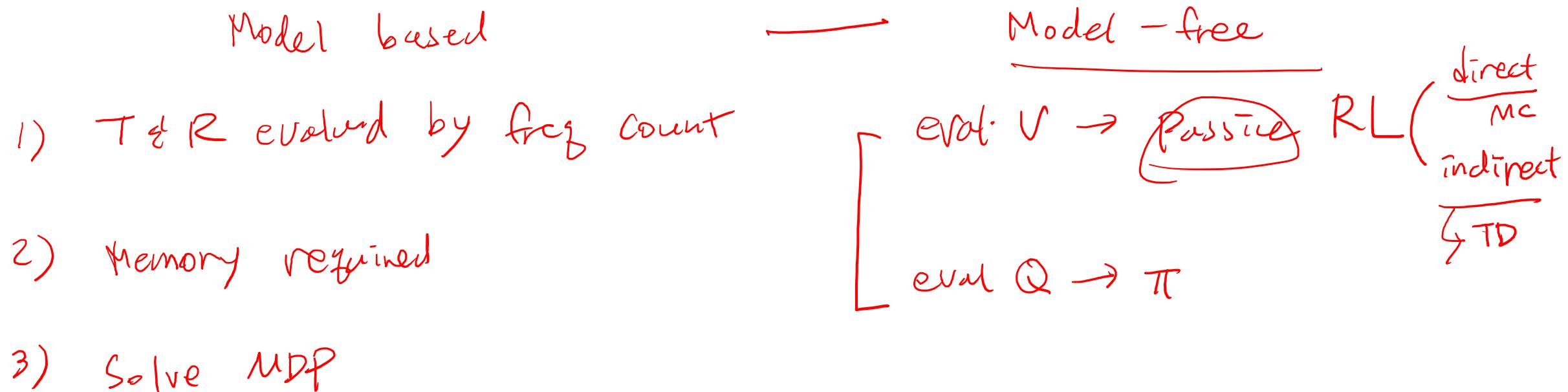
$$\hat{T}(s, a, s')$$

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-Free Learning



Example: Expected Age

Goal: Compute expected age of students in a class

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.

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

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

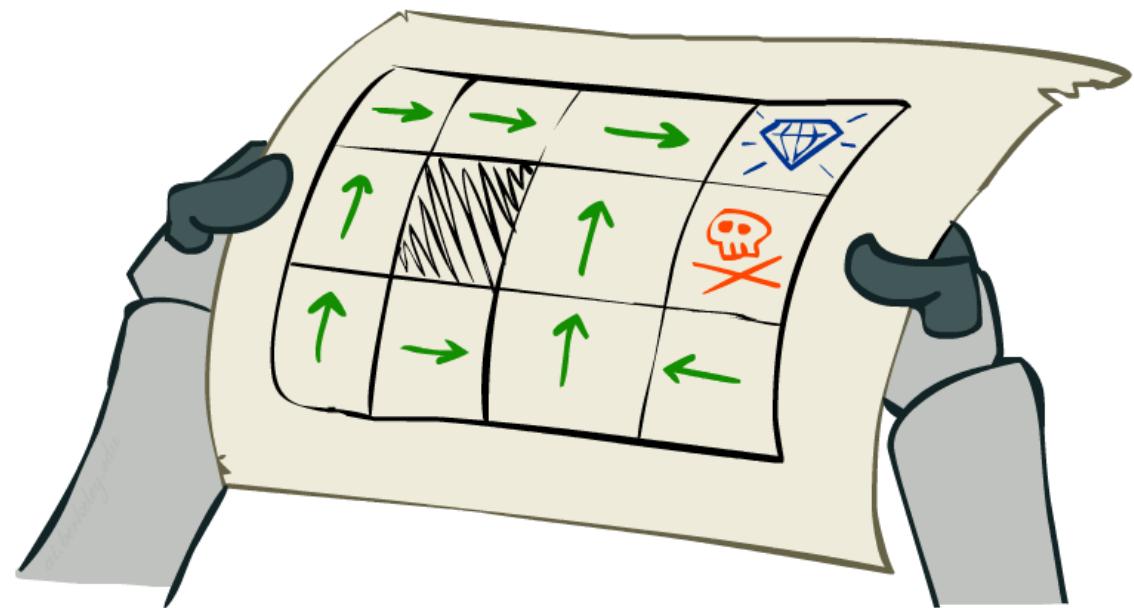
Unknown $P(A)$: “Model Free”

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

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

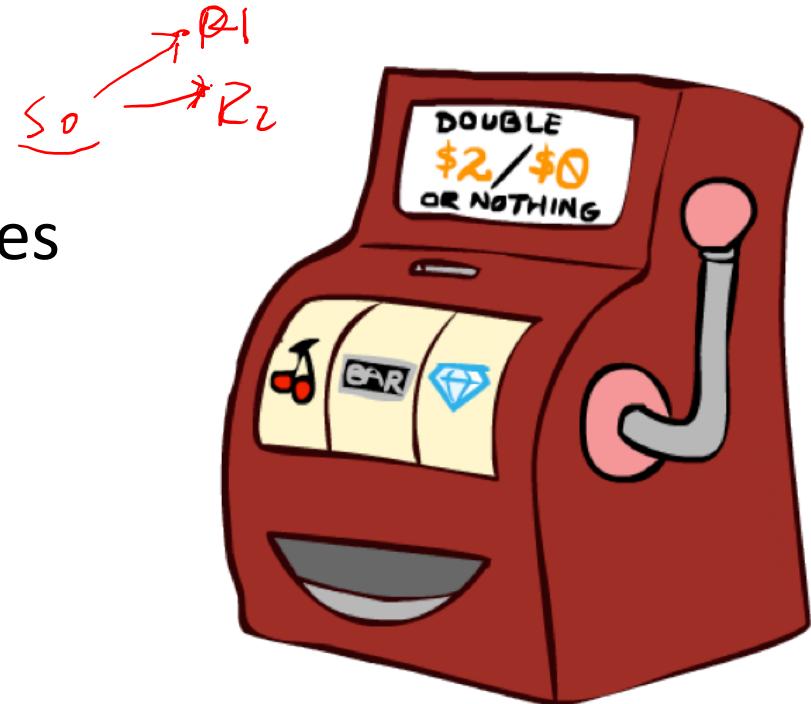
Passive Reinforcement Learning

- Simplified task: policy evaluation
 - Input: a fixed policy $\pi(s)$ given
 - 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.



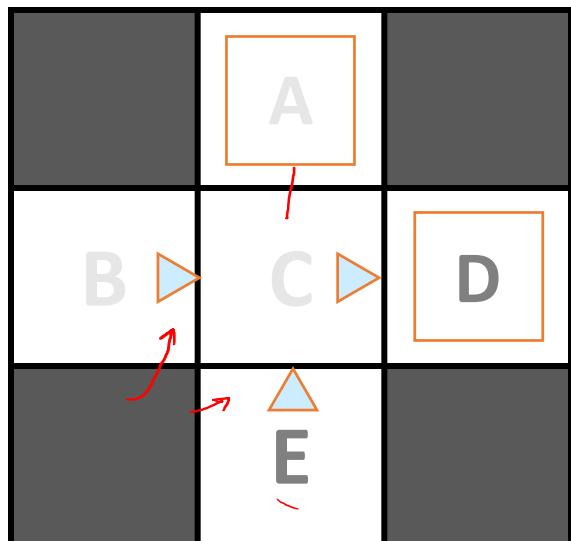
Direct Evaluation

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

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

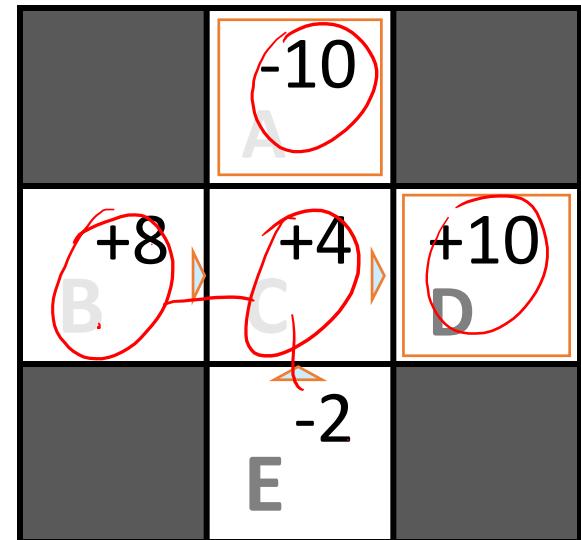
Sum of discounted reward
along the way
Output Values

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

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 basis)~~ how can
their values be different?
$$R(s) + \gamma \sum_{s'} P(s'|s) \pi(a|s')$$

$$V^\pi(s)$$

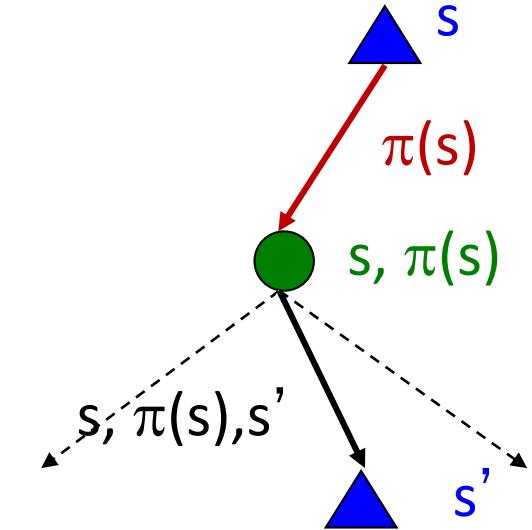
Why Not Use Policy Evaluation?

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

$$V_0^\pi(s) = 0 \quad \begin{matrix} s \rightarrow s \\ V(s) \end{matrix}$$

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

- 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

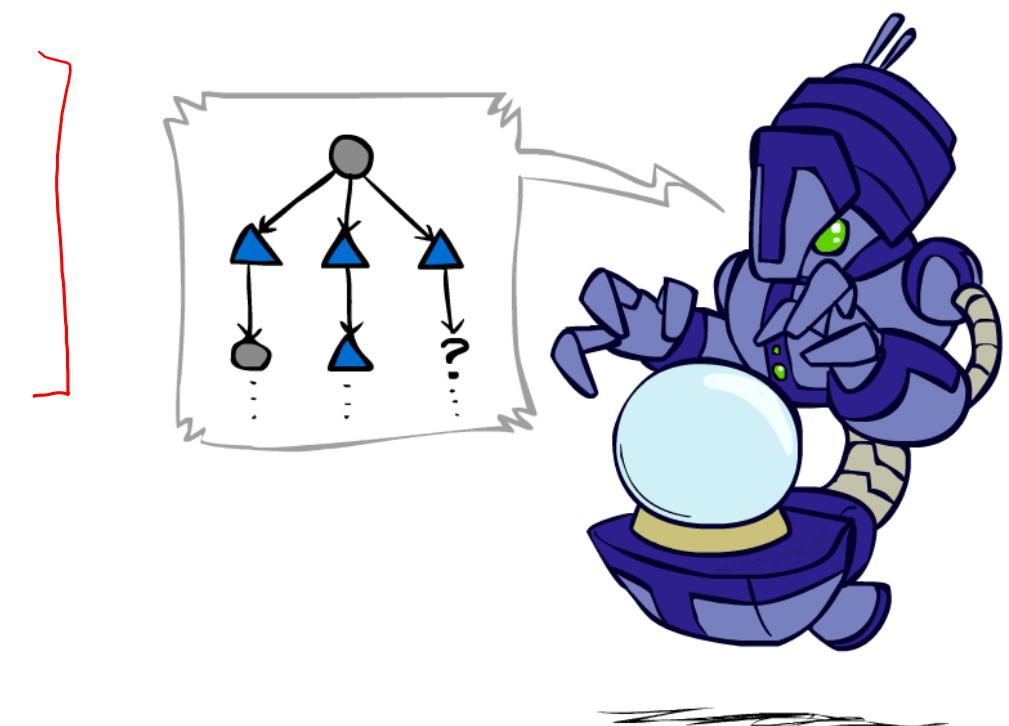
$$sample_1 = R(\cancel{s}, \pi(s), s'_1) + \gamma V_k^{\pi}(s'_1)$$

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

...

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

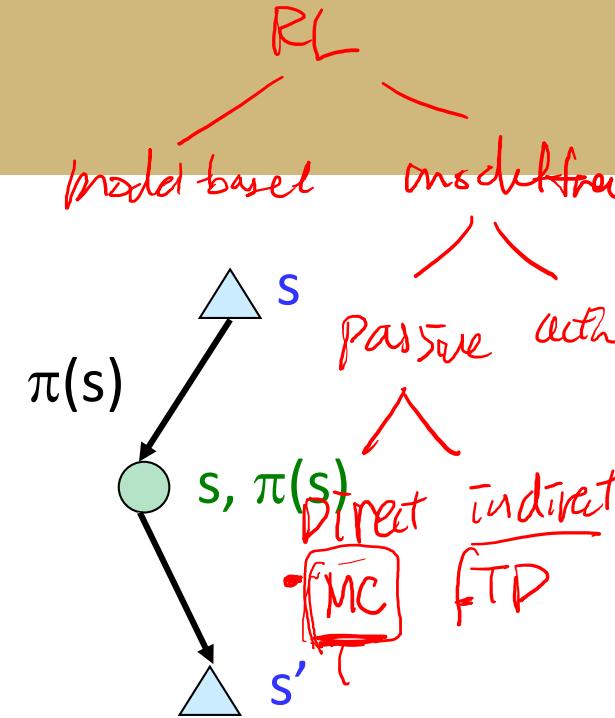
$$V_{k+1}^{\pi}(s) \leftarrow \boxed{\frac{1}{n} \sum_i sample_i}$$



Temporal Difference Learning

(Indirect)

- 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 $V(s)$:

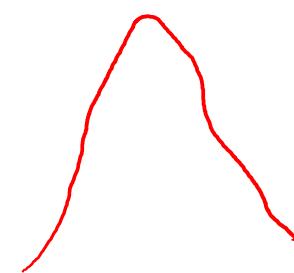
$$sample = R(s, \pi(s), s') + \gamma V^\pi(s')$$

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

- Makes recent samples more important:

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

The equation shows the formula for the Exponential Moving Average. The term x_n is circled in red with an arrow pointing to it. The term $(1 - \alpha)$ is circled in red with an arrow pointing to it. The term x_{n-1} is circled in red with an arrow pointing to it. The term $(1 - \alpha)^2$ is circled in red with an arrow pointing to it. The term x_{n-2} is circled in red with an arrow pointing to it. The ellipsis \dots is circled in red with an arrow pointing to it.

- 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	

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

$s = C$

$s' = D$

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

~~$V^\pi(s)$~~ \cancel{s} $\cancel{V^\pi(s)}$ $\cancel{\gamma}$ $\cancel{V^\pi(s')}$

$3 \quad 16$

-2

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!

