

CMU 15-781

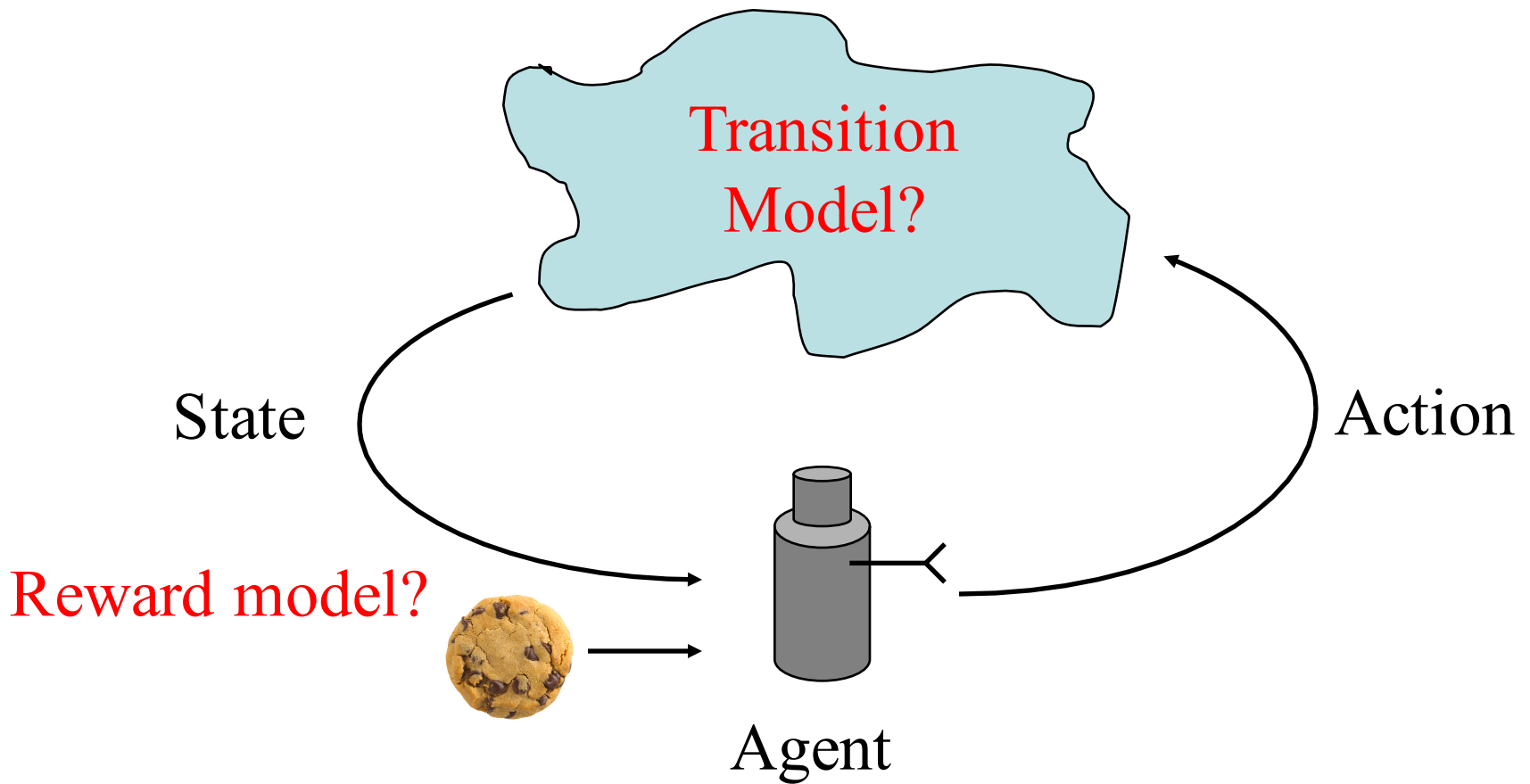
Lecture 12:

Reinforcement Learning

Teacher:

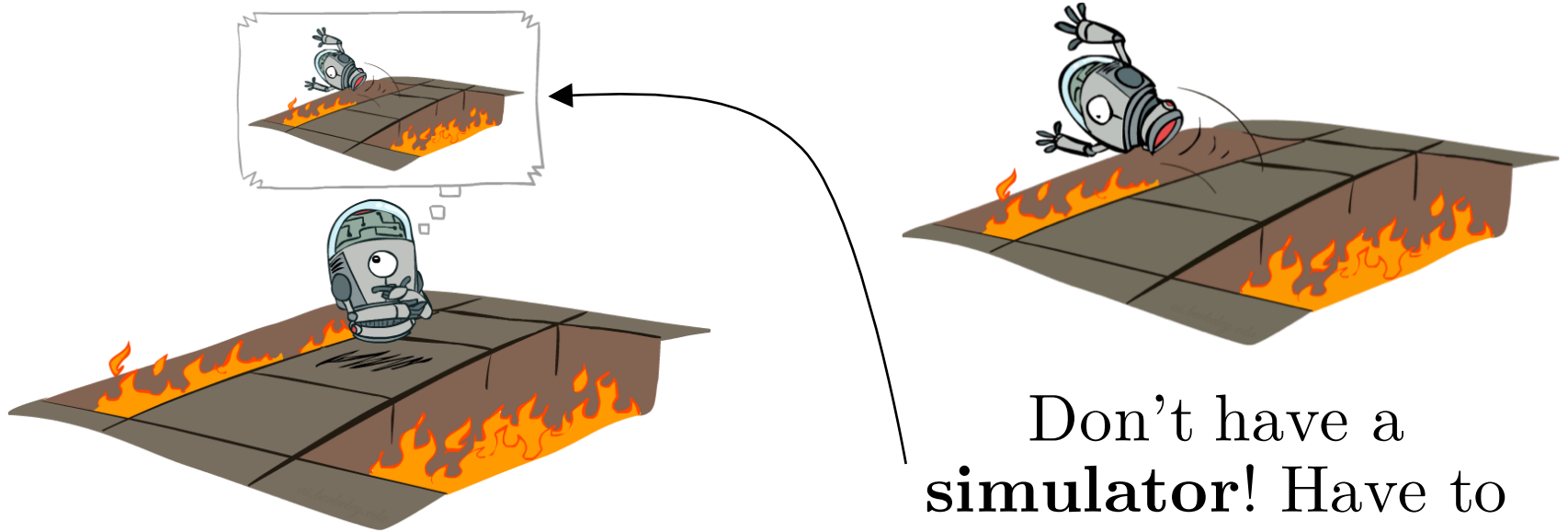
Gianni A. Di Caro

# REINFORCEMENT LEARNING



*Goal: Maximize expected sum of future rewards*

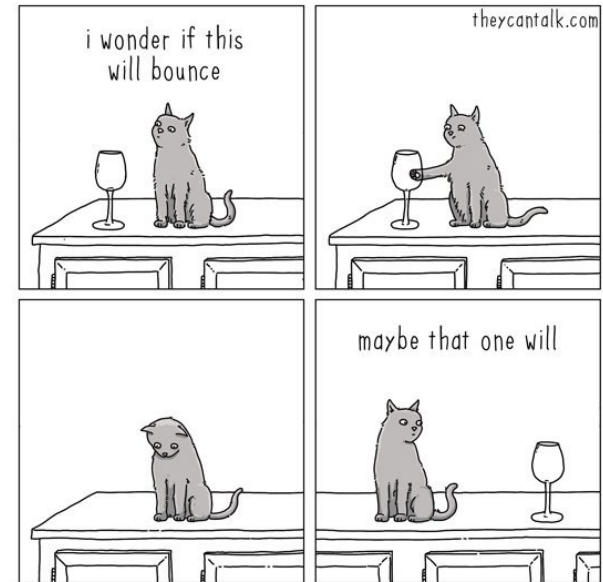
# MDP PLANNING VS REINFORCEMENT LEARNING



Don't have a **simulator!** Have to actually **learn** what happens if take an action in a state

# REINFORCEMENT LEARNING PROBLEM

- ✓ The agents can "sense" the environment (it knows the state) and has goals
- ✓ Learning from interaction with the environment
  - Trial and error search
  - (*Delayed*) Rewards (**Advisory signals  $\neq$  errors signals**)
  - What actions to take?
    - **Exploration- exploitation dilemma**



# PASSIVE REINFORCEMENT LEARNING

- Before figuring out how to act, let's first just try to figure out how good a particular policy is
- **Passive learning:** agent's policy is fixed (i.e., in state  $s$  it always execute action  $\pi(s)$ ) and the task is to policy's value  $\rightarrow$  **Learn state values** ( $V(s)$ ) or state-action values ( $Q(s,a)$ )

Policy evaluation in MDPs  $\sim$  Passive RL

(T,R) Model  
Bellman eqs.

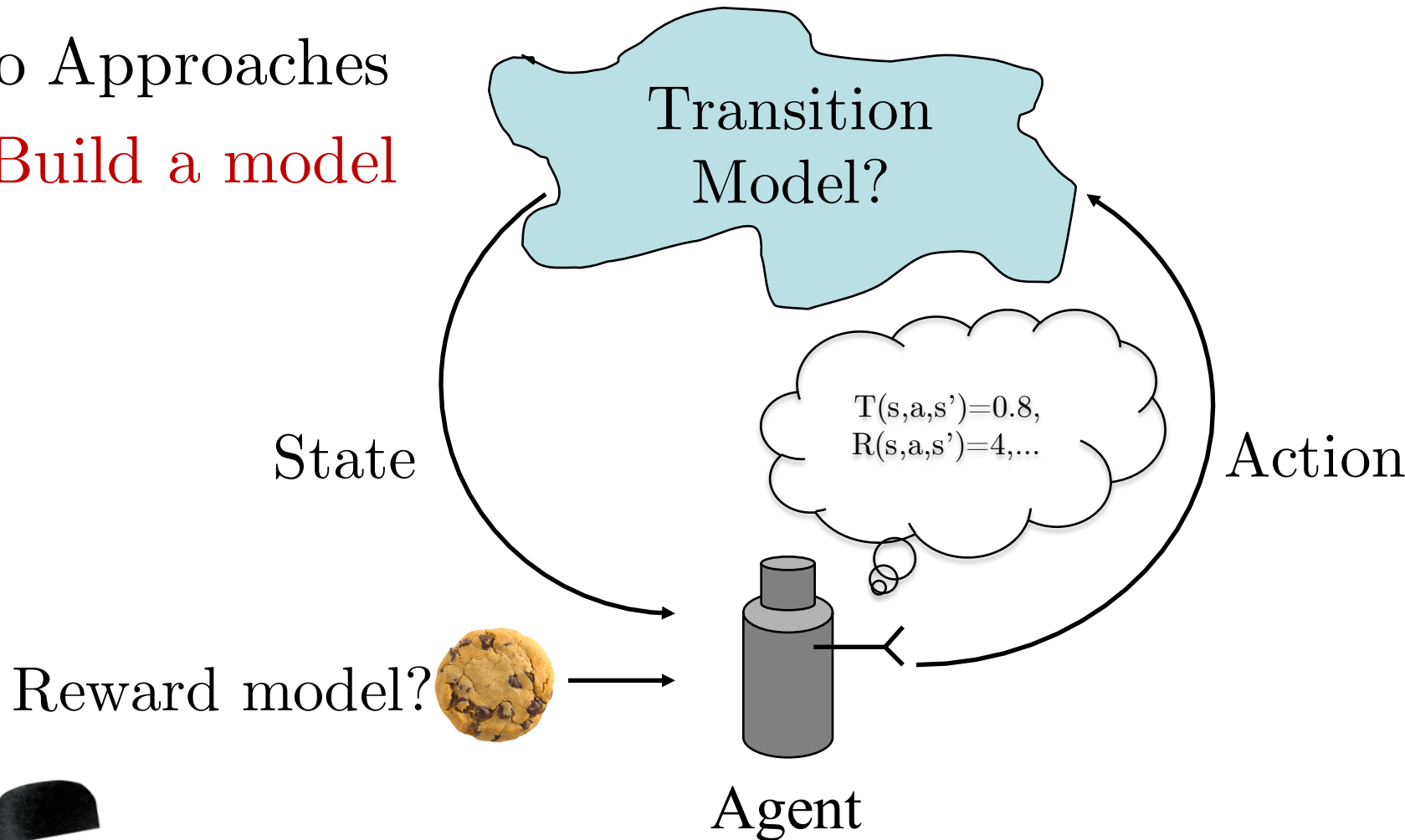
~~(T,R) Model~~  
Learning



# PASSIVE REINFORCEMENT LEARNING

Two Approaches

1. Build a model

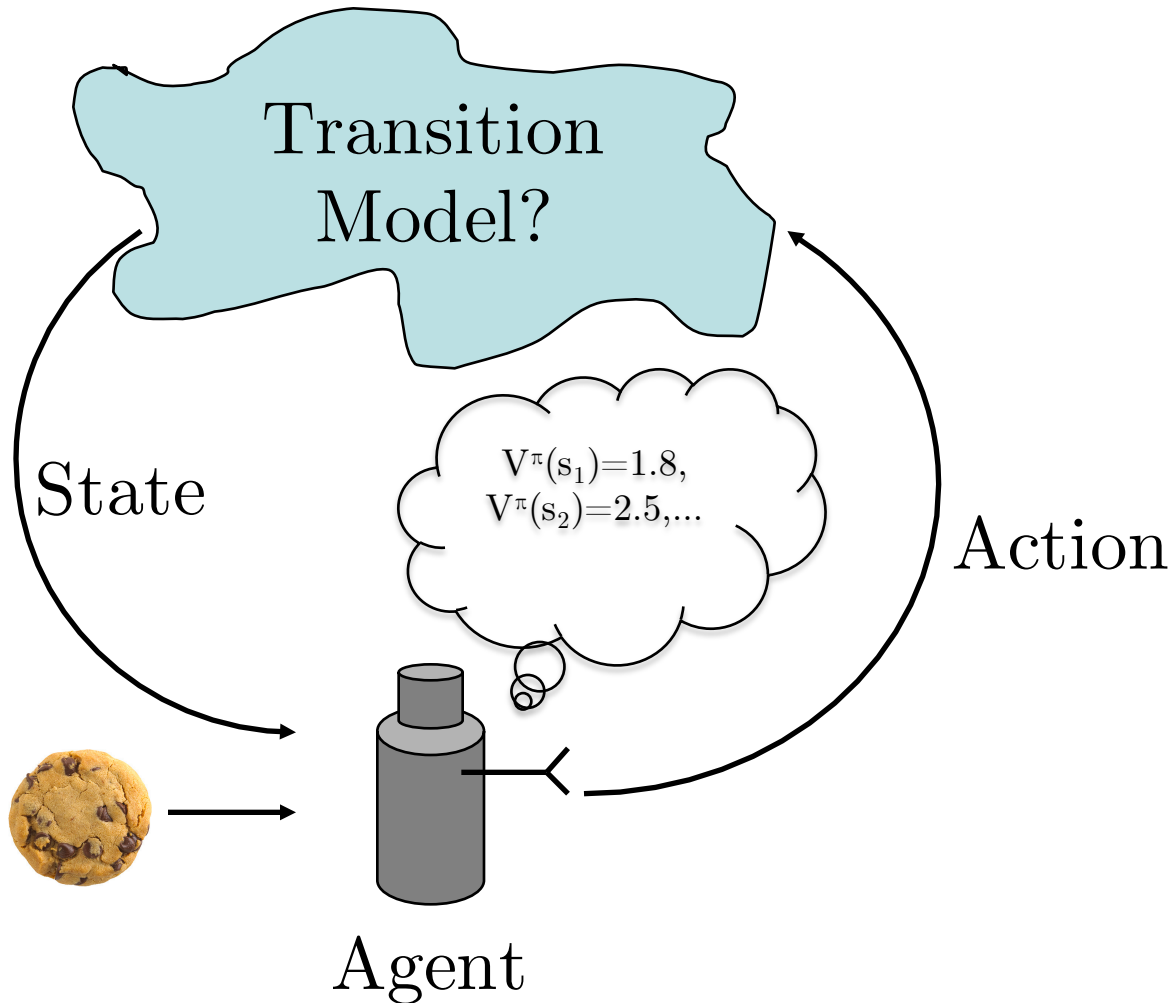


# PASSIVE REINFORCEMENT LEARNING

Two Approaches:

1. Build a model
2. Model-free:  
directly  
estimate  $V^\pi$

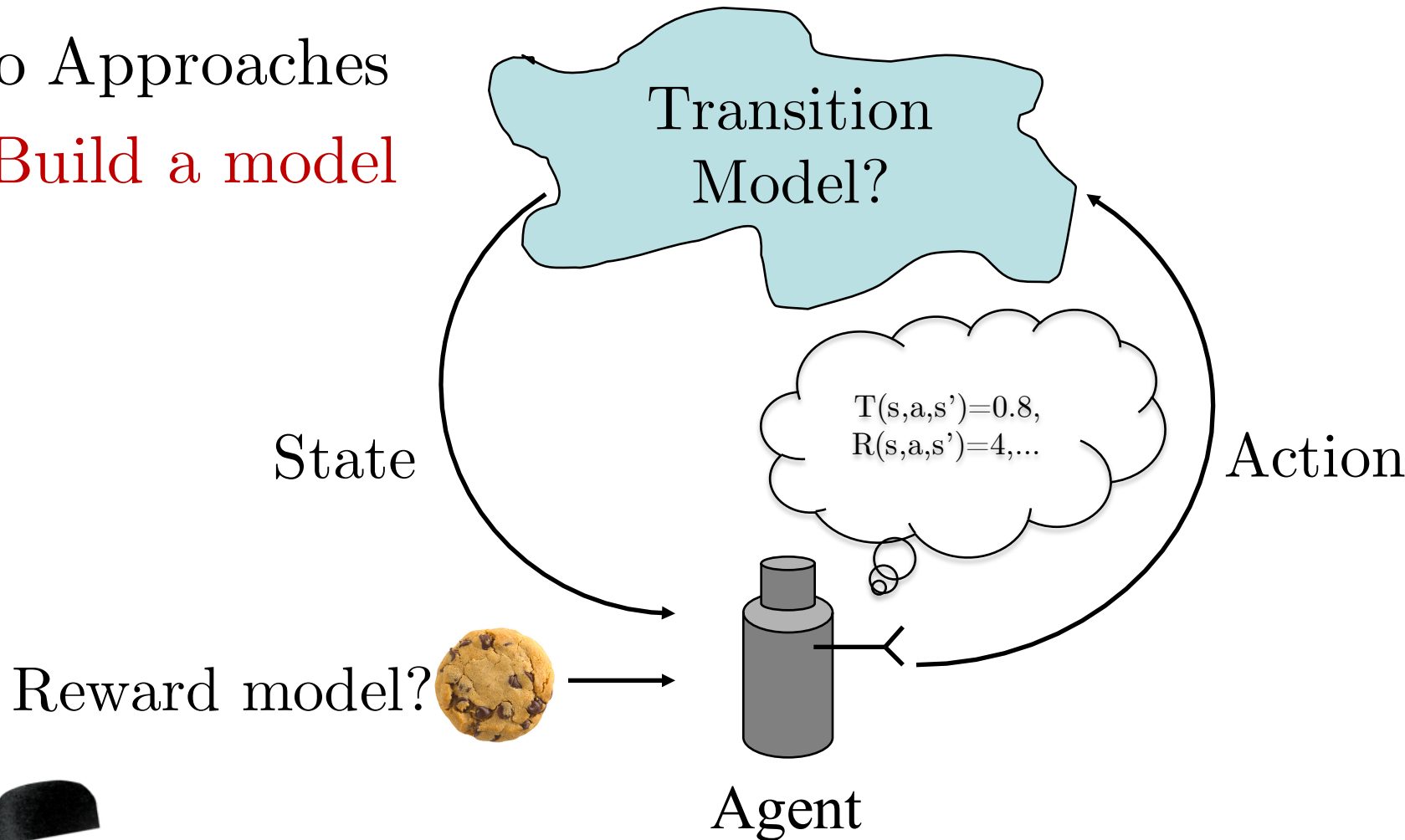
Reward model?



# PASSIVE REINFORCEMENT LEARNING

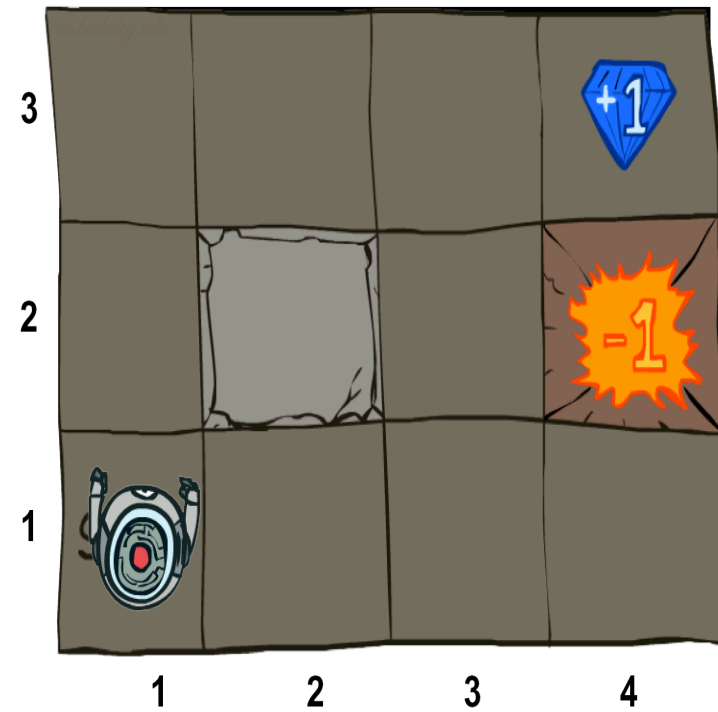
Two Approaches

1. Build a model



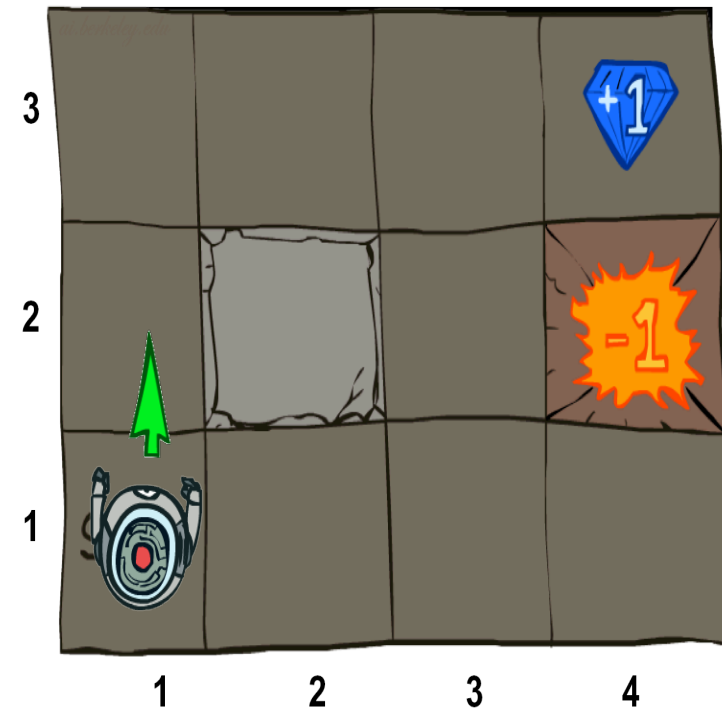


Start at (1,1)



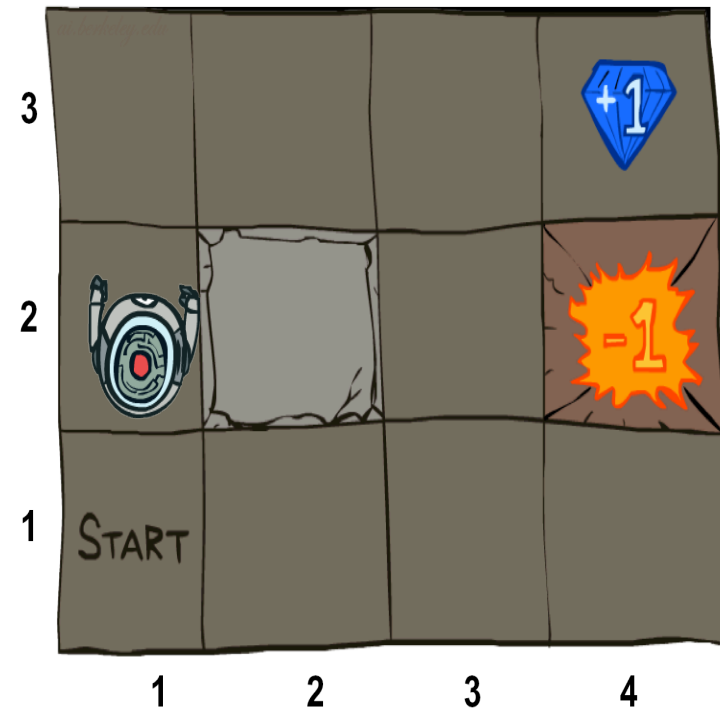
Start at (1,1)

$s=(1,1)$  action= *tup* *try* up



Start at (1,1)

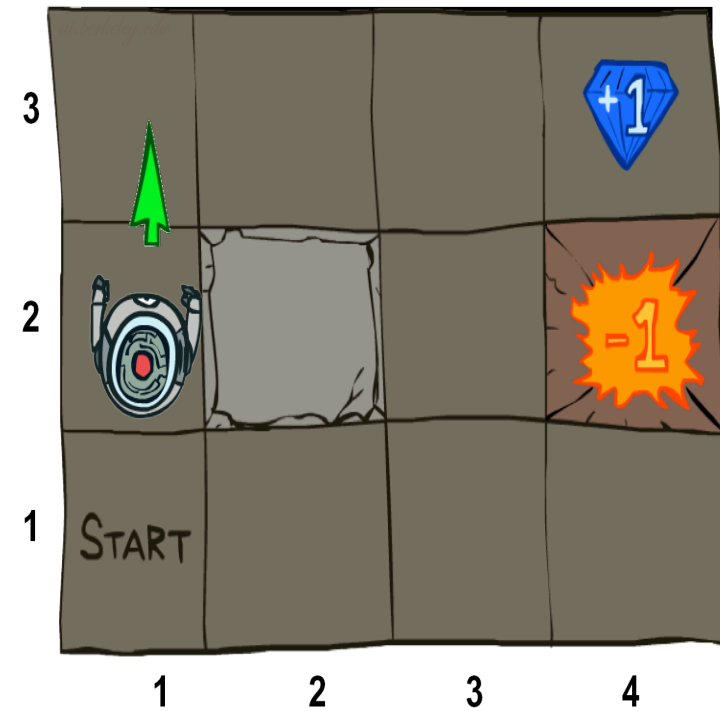
$s=(1,1)$  action= tup,  $s'=(1,2)$ ,  $r = -.01$



Start at (1,1)

$s=(1,1)$  action= tup,  $s'=(1,2)$ ,  $r = -.01$

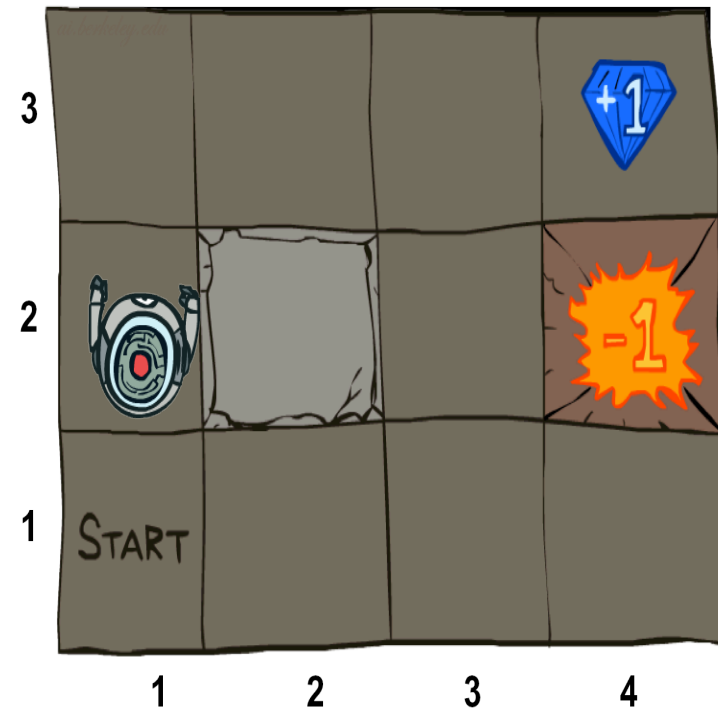
$s=(1,2)$  action= tup



Start at (1,1)

$s=(1,1)$  action= tup,  $s'=(1,2)$ ,  $r = -.01$

$s=(1,2)$  action= tup,  $s'=(1,2)$ ,  $r = -.01$



Start at (1,1)

$s=(1,1)$  action=tup,  $s'=(1,2)$ ,  $r = -.01$

$s=(1,2)$  action=tup,  $s'=(1,2)$ ,  $r = -.01$

$s=(1,2)$  action=tup,  $s'=(1,3)$ ,  $r = -.01$

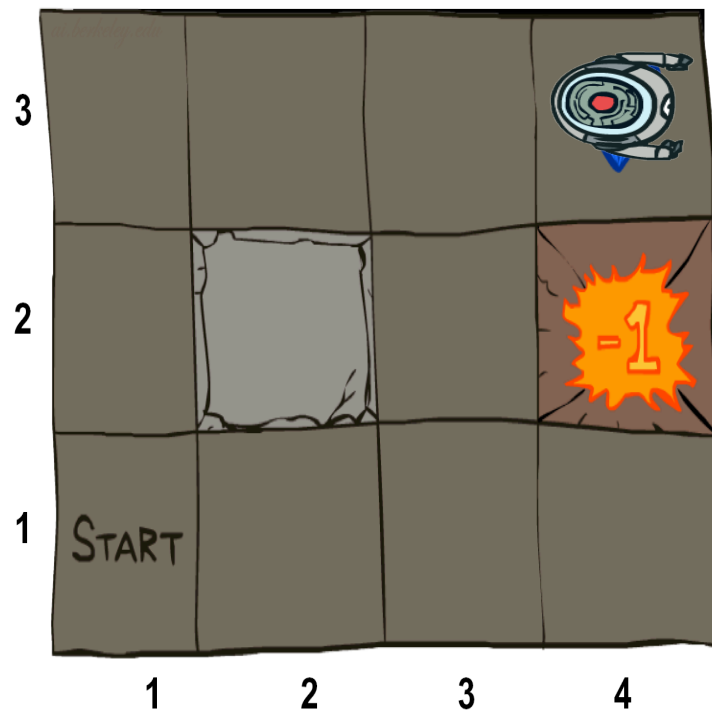
$s=(1,3)$  action=tright,  $s'=(2,3)$ ,  $r = -.01$

$s=(2,3)$  action=tright,  $s'=(3,3)$ ,  $r = -.01$

$s=(3,3)$  action=tright,  $s'=(3,2)$ ,  $r = -.01$

$s=(3,2)$  action=tup,  $s'=(3,3)$ ,  $r = -.01$

$s=(3,3)$  action=tright,  $s'=(4,3)$ ,  $r = 1$



Start at (1,1)

$s=(1,1)$  action= tup,  $s'=(1,2)$ ,  $r = -.01$

$s=(1,2)$  action= tup,  $s'=(1,2)$ ,  $r = -.01$

$s=(1,2)$  action=tup,  $s'=(1,3)$ ,  $r = -.01$

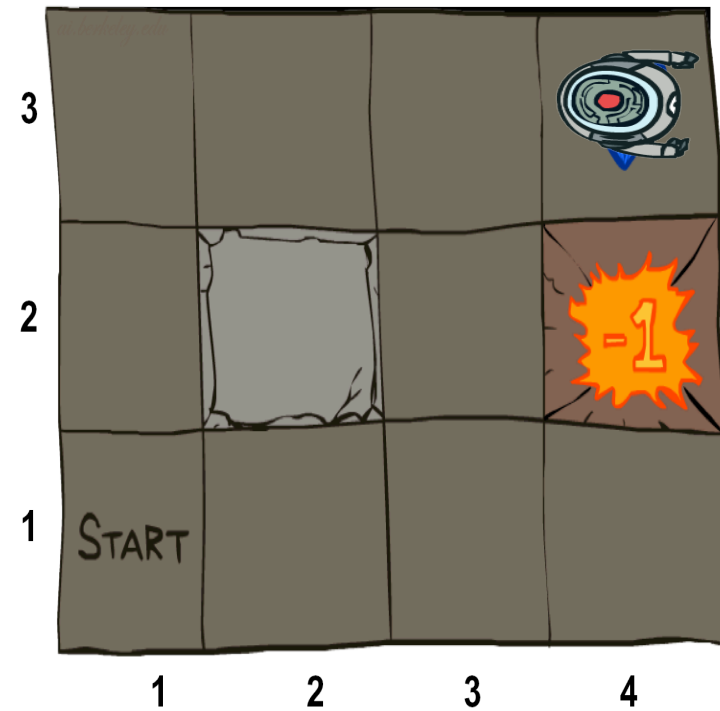
$s=(1,3)$  action=tright,  $s'=(2,3)$ ,  $r = -.01$

$s=(2,3)$  action=tright,  $s'=(3,3)$ ,  $r = -.01$

$s=(3,3)$  action=tright,  $s'=(3,2)$ ,  $r = -.01$

$s=(3,2)$  action=tup,  $s'=(3,3)$ ,  $r = -.01$

$s=(3,3)$  action=tright,  $s'=(4,3)$ ,  $r = 1$



The gathered experience can be used to estimate MDP's  $T$  and  $R$  models



Start at (1,1)

$s=(1,1)$  action=tup,  $s'=(1,2)$ ,  $r = -.01$

$s=(1,2)$  action=tup,  $s'=(1,2)$ ,  $r = -.01$

$s=(1,2)$  action=tup,  $s'=(1,3)$ ,  $r = -.01$

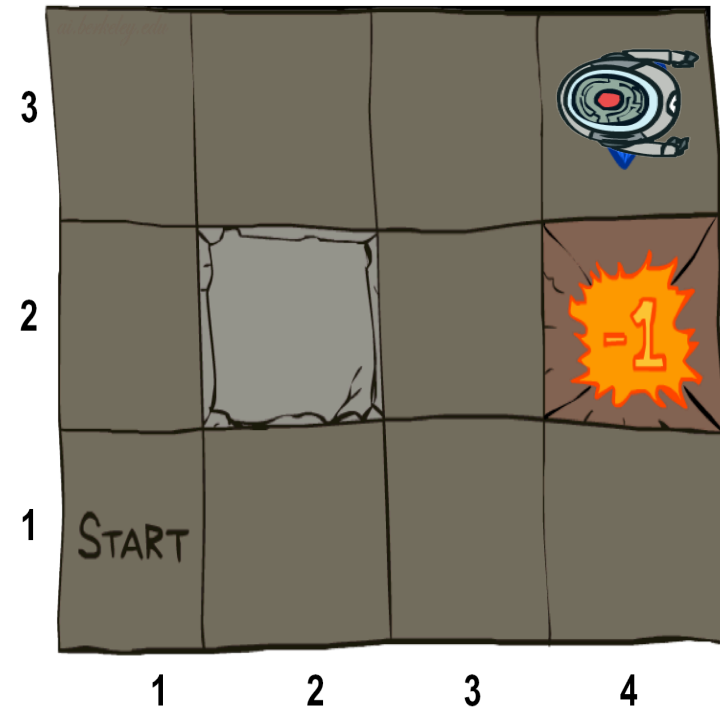
$s=(1,3)$  action=tright,  $s'=(2,3)$ ,  $r = -.01$

$s=(2,3)$  action=tright,  $s'=(3,3)$ ,  $r = -.01$

$s=(3,3)$  action=tright,  $s'=(3,2)$ ,  $r = -.01$

$s=(3,2)$  action=tup,  $s'=(3,3)$ ,  $r = -.01$

$s=(3,3)$  action=tright,  $s'=(4,3)$ ,  $r = 1$



The gathered experience can be used to estimate MDP's  $T$  and  $R$  models

Estimate of  $T(\langle 1,2 \rangle, \text{tup}, \langle 1,3 \rangle) = 1/2$



# MODEL-BASED PASSIVE REINFORCEMENT LEARNING

1. Follow policy  $\pi$
2. Estimate MDP model parameters  $T$  and  $R$  given observed transitions and rewards
  1. If finite set of states and actions, can just count and average counts
3. Use estimated MDP to do policy evaluation of  $\pi$

Does this give us all the parameters for an MDP?



Start at (1,1)

$s=(1,1)$  action= tup,  $s'=(1,2)$ ,  $r = -.01$

$s=(1,2)$  action= tup,  $s'=(1,2)$ ,  $r = -.01$

$s=(1,2)$  action=tup,  $s'=(1,3)$ ,  $r = -.01$

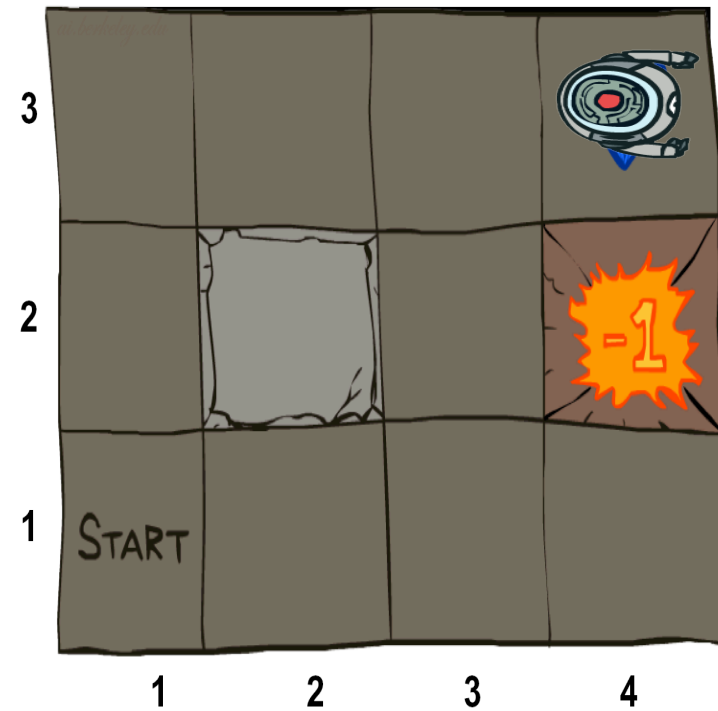
$s=(1,3)$  action=tright,  $s'=(2,3)$ ,  $r = -.01$

$s=(2,3)$  action=tright,  $s'=(3,3)$ ,  $r = -.01$

$s=(3,3)$  action=tright,  $s'=(3,2)$ ,  $r = -.01$

$s=(3,2)$  action=tup,  $s'=(3,3)$ ,  $r = -.01$

$s=(3,3)$  action=tright,  $s'=(4,3)$ ,  $r = 1$



Estimate of  
 $T(\langle 1,2 \rangle, \text{tright}, \langle 1,3 \rangle)$ ?  
No idea! Never tried this  
action...



DOES THIS GIVE US ALL THE PARAMETERS FOR A MDP?

No.

BUT DOES THAT MATTER FOR COMPUTING POLICY VALUE?

No.

HAVE ALL PARAMETERS WE NEED!

(AFTER WE VISIT ALL STATES AT LEAST ONCE)

- Follow policy  $\pi$
- Estimate MDP model parameters given observed transitions and rewards
  - If finite set of states and actions, can just count and average counts
- Use estimated MDP to do policy evaluation of  $\pi$



Start at (1,1)

$s=(1,1)$  action= tup,  $s'=(1,2)$ ,  $r = -.01$

$s=(1,2)$  action= tup,  $s'=(1,2)$ ,  $r = -.01$

$s=(1,2)$  action=tup,  $s'=(1,3)$ ,  $r = -.01$

$s=(1,3)$  action=tright,  $s'=(2,3)$ ,  $r = -.01$

$s=(2,3)$  action=tright,  $s'=(3,3)$ ,  $r = -.01$

$s=(3,3)$  action=tright,  $s'=(3,2)$ ,  $r = -.01$

$s=(3,2)$  action=tup,  $s'=(3,3)$ ,  $r = -.01$

$s=(3,3)$  action=tright,  $s'=(4,3)$ ,  $r = 1$

$s=(1,1)$  action= tup,  $s'=(2,1)$ ,  $r = -.01$

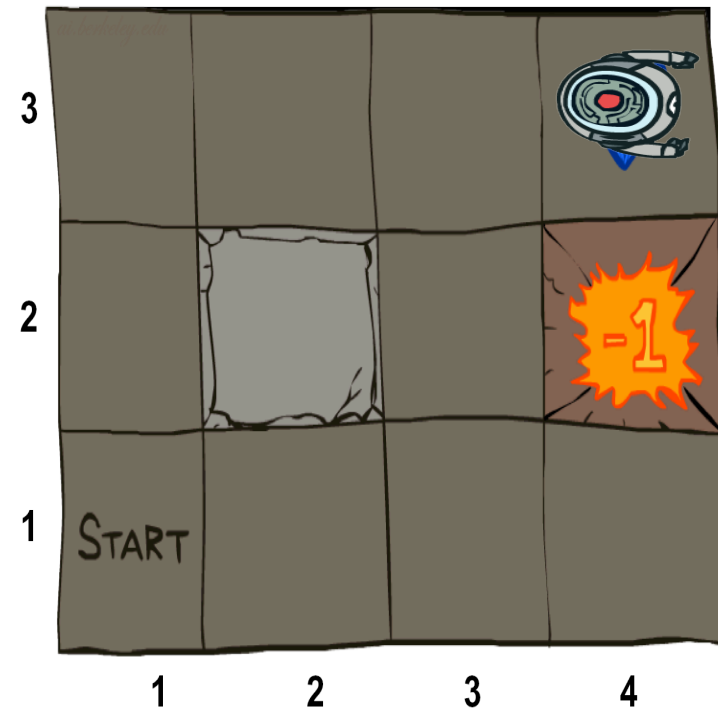
$s=(2,1)$  action= tright,  $s'=(3,1)$ ,  $r = -.01$

$s=(3,1)$  action= tup,  $s'=(4,1)$ ,  $r = -.01$

$s=(4,1)$  action= tleft,  $s'=(3,1)$ ,  $r = -.01$

$s=(3,1)$  action= tup,  $s'=(3,2)$ ,  $r = -.01$

$s=(3,2)$  action= tup,  $s'=(4,2)$ ,  $r = -1$



2 episodes of experience in MDP.  
Use to estimate MDP parameters  
& evaluate  $\pi$

Is the computed policy value  
likely to be correct?

(1) Yes (2) No (3) Not sure

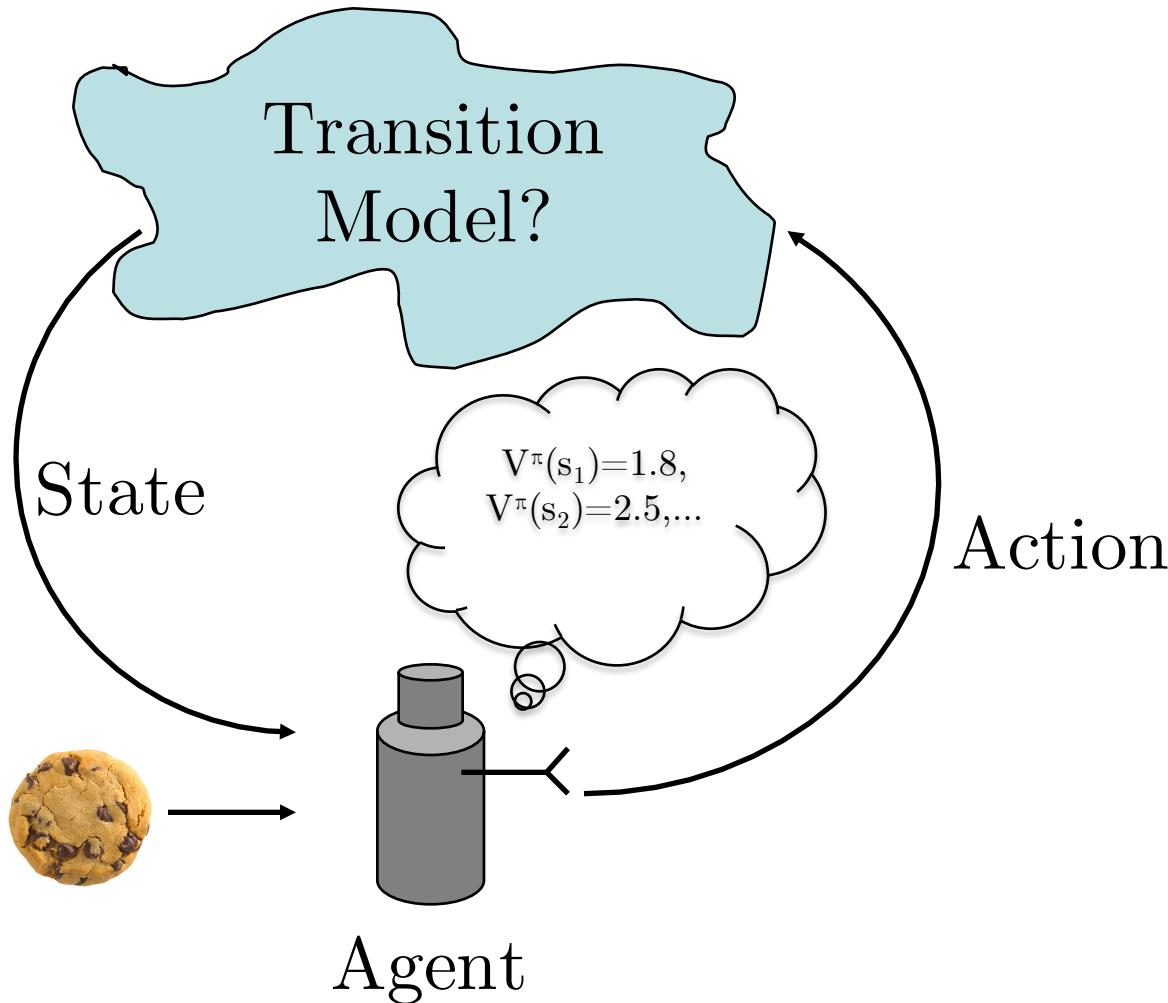


# PASSIVE REINFORCEMENT LEARNING

Two Approaches:

1. Build a model
2. Model-free:  
directly  
estimate  $V^\pi$

Reward model?



# DIRECT UTILITY ESTIMATION: MONTE CARLO POLICY EVALUATION

## First-visit MC policy evaluation (returns $V \approx v_\pi$ )

Initialize:

$\pi \leftarrow$  policy to be evaluated

$V \leftarrow$  an arbitrary state-value function

$Returns(s) \leftarrow$  an empty list, for all  $s \in \mathcal{S}$

Repeat forever:

Generate an episode using  $\pi$

For each state  $s$  appearing in the episode:

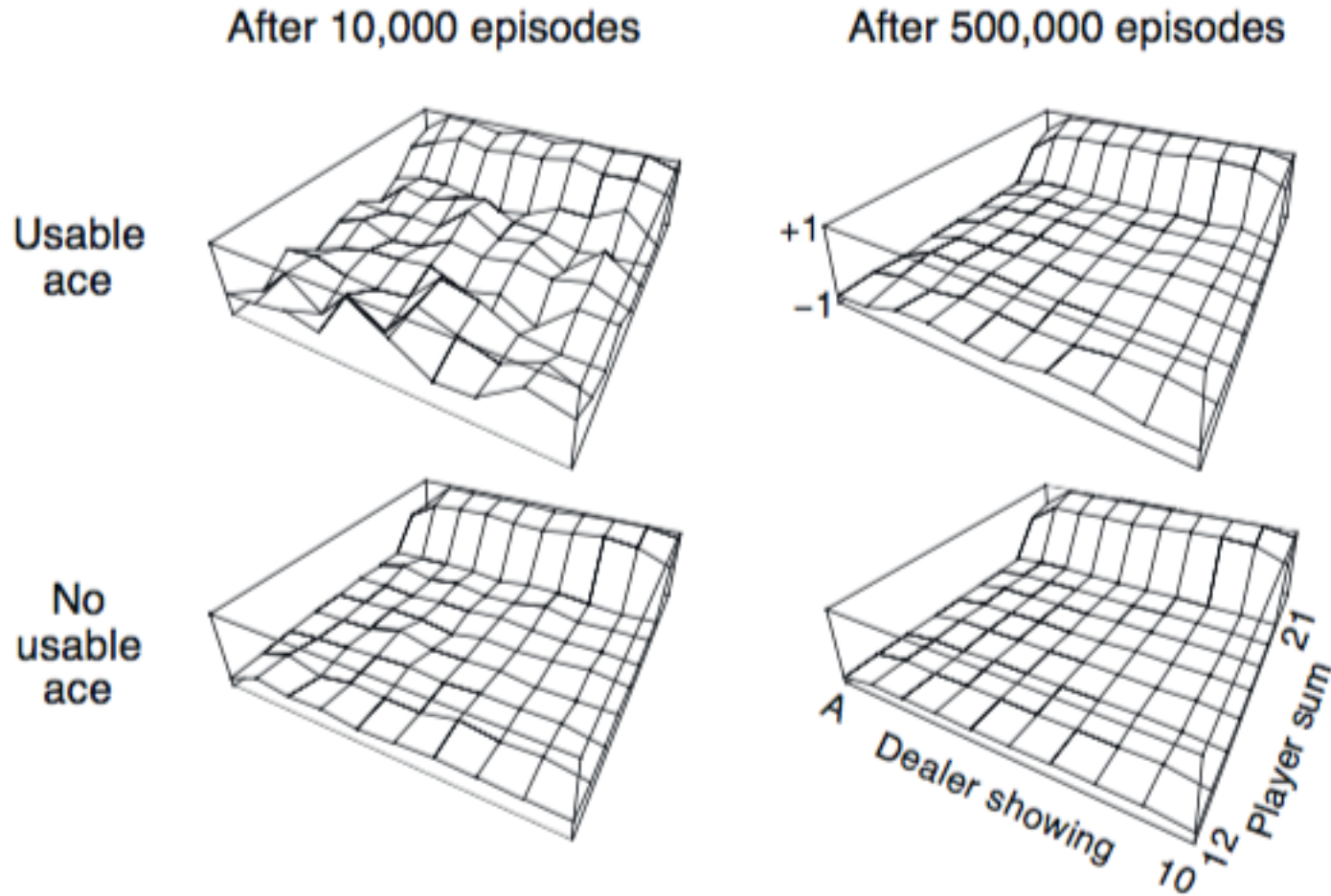
$G \leftarrow$  return following the first occurrence of  $s$

Append  $G$  to  $Returns(s)$

$V(s) \leftarrow$  average( $Returns(s)$ )



# MONTE CARLO POLICY EVALUATION



# MONTE CARLO POLICY EVALUATION

What information is *missing* in the way we are proceeding performing policy evaluation for our MDP? (it's an MDP, we just do not know its parameters ...)

**The recursive Bellman relations among states!**

$$V^\pi(s) = \sum_{s' \in S} p(s' | s, \pi(s)) \left[ R(s, \pi(s), s') + \gamma V^\pi(s') \right] \quad \forall s \in S$$

We should exploit this structure in the MDP...





# TEMPORAL DIFFERENCES (TD) POLICY EVALUATION

## Tabular TD(0) for estimating $v_\pi$

Input: the policy  $\pi$  to be evaluated  
Initialize  $V(s)$  arbitrarily (e.g.,  $V(s) = 0, \forall s \in \mathcal{S}^+$ )  
Repeat (for each episode):  
  Initialize  $S$   
  Repeat (for each step of episode):  
     $A \leftarrow$  action given by  $\pi$  for  $S$   
    Take action  $A$ , observe  $R, S'$   
     $V(S) \leftarrow V(S) + \alpha [R + \gamma V(S') - V(S)]$   
     $S \leftarrow S'$   
  until  $S$  is terminal

Update  
at each step!

Sample +  
Bellman

TD *bootstraps* on available state information

# TEMPORAL DIFFERENCE LEARNING

- No explicit model of  $T$  or  $R$ !
- Estimate  $V$  and expectation **through samples**
- Update from **every experience**
  - Update  $V(s)$  after each state 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 sample of  $V$ : **running average**

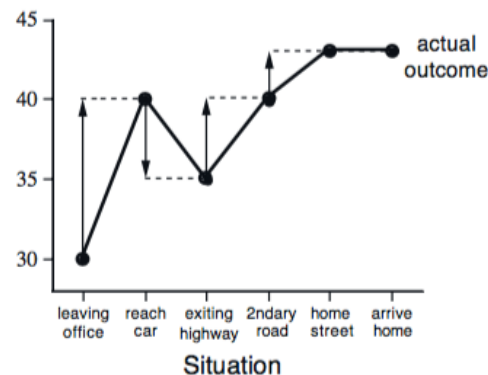
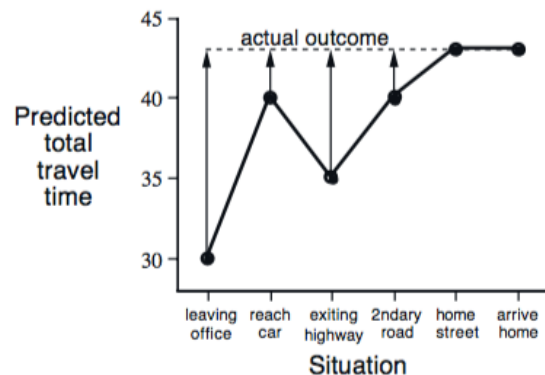
Sample of  $V(s)$ :  $sample = R(s, \pi(s), s') + \gamma V^\pi(s')$

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



# TD EXAMPLE: GOING HOME

<i>State</i>	<i>Elapsed Time (minutes)</i>	<i>Predicted Time to Go</i>	<i>Predicted Total Time</i>
leaving office, friday at 6	0	30	30
reach car, raining	5	35	40
exiting highway	20	15	35
2ndary road, behind truck	30	10	40
entering home street	40	3	43
arrive home	43	0	43



# EXPONENTIAL MOVING AVG

- 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:

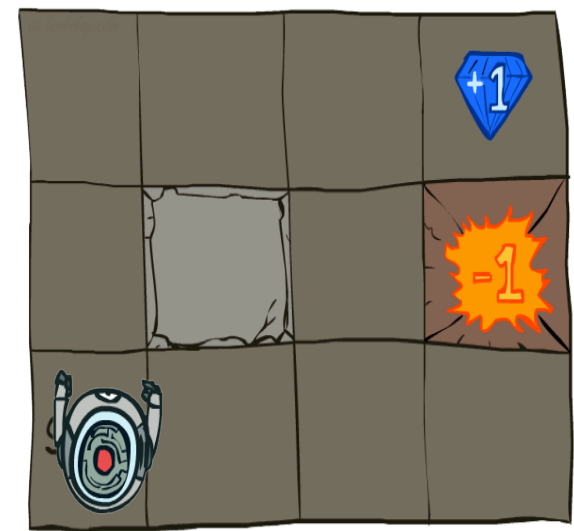
$$\bar{x}_n = \frac{x_n + (1 - \alpha) \cdot x_{n-1} + (1 - \alpha)^2 \cdot x_{n-2} + \dots}{1 + (1 - \alpha) + (1 - \alpha)^2 + \dots}$$

- Forgets about the past

- Decreasing learning rate ( $\alpha$ ) can give converging avgs



# TD Learning Example

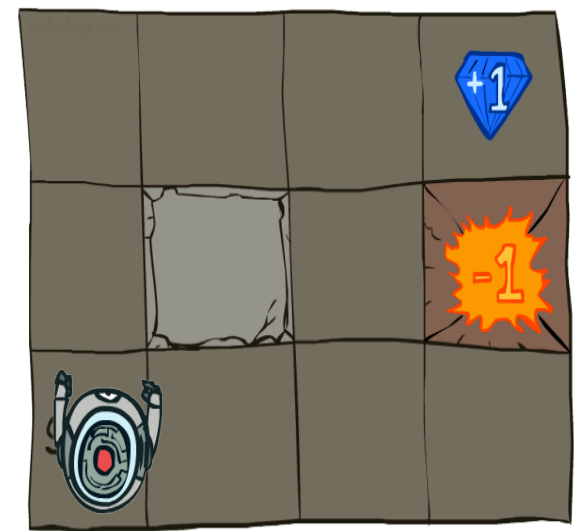


Initialize all  $V^\pi(s)$  values:  $V^\pi(s) = 0$

State	$V^\pi(s)$
(1,1)	0
(1,2)	0
(1,3)	0
(4,1)	0
(2,1)	0
(2,3)	0
(3,1)	0
(3,2)	0
(3,3)	0

## TD Learning Example $\alpha=0.1, \gamma=1$

$s=(1,1)$  action= *tup*,  $s'=(1,2)$ ,  $r = -0.01$



Update  $V^\pi((1,1))$

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

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

$$\text{sample} = -0.01 + V^\pi((1,2)) = -0.01$$

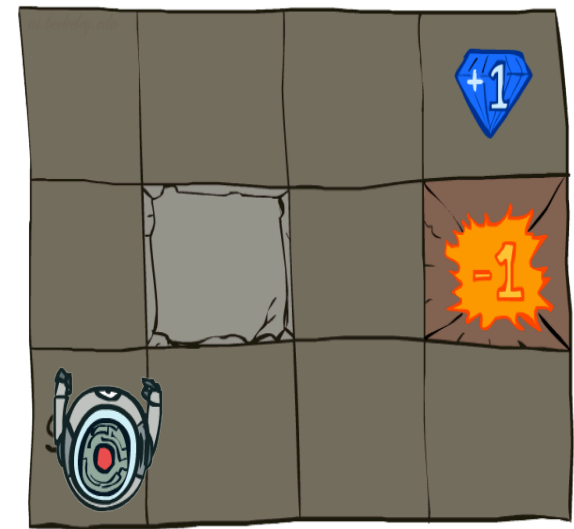
$$V^\pi((1,1)) = (1-\alpha)V^\pi((1,1)) + \alpha*\text{sample}$$

$$= .9*0 + 0.1*-0.01 = -0.001$$

State	$V^\pi(s)$
(1,1)	0
(1,2)	0
(1,3)	0
(4,1)	0
(2,1)	0
(2,3)	0
(3,1)	0
(3,2)	0
(3,3)	0

## TD Learning Example $\alpha=0.1, \gamma=1$

$s=(1,1)$  action= *tup*,  $s'=(1,2)$ ,  $r = -0.01$



Update  $V^\pi((1,1))$

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

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

$$\text{sample} = -0.01 + V^\pi((1,2)) = -0.01$$

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

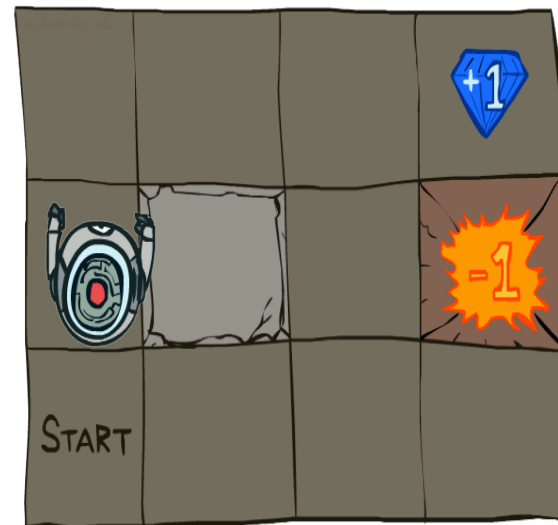
$$= .9 \cdot 0 + 0.1 \cdot -0.01 = -0.001$$

State	$V^\pi(s)$
(1,1)	-0.001
(1,2)	0
(1,3)	0
(4,1)	0
(2,1)	0
(2,3)	0
(3,1)	0
(3,2)	0
(3,3)	0

## TD Learning Example $\alpha=0.1, \gamma=1$

$s=(1,1)$  action= tup,  $s'=(1,2)$ ,  $r = -.01$

$s=(1,2)$  action= tup,  $s'=(1,2)$ ,  $r = -.01$



State	$V^\pi(s)$
(1,1)	-0.001
(1,2)	0
(1,3)	0
(4,1)	0
(2,1)	0
(2,3)	0
(3,1)	0
(3,2)	0
(3,3)	0





# PROBLEMS WITH PASSIVE LEARNING

- Agent wants to ultimately learn to *act* to gather high reward in the environment.
- Using a deterministic policy, gives no experience for other actions (not included in the policy)

**Active reinforcement learning:** the agent decides what action to take with the goal of learning an optimal policy



# ACTIVE RL: EXPLORATION ISSUES

- How the actions are selected for the purpose of efficient learning?
- Consider acting **randomly** in the world
- Can such experience allow the agent to learn the optimal values and policy?



# RECALL MODEL-BASED PASSIVE REINFORCEMENT LEARNING

- Follow policy  $\pi$
- Estimate MDP model parameters given observed transitions and rewards
  - If finite set of states and actions, can just count and average counts
- Use estimated MDP to do policy evaluation of  $\pi$



# MODEL-BASED ACTIVE RL W/RANDOM ACTIONS

- Choose actions randomly
- Estimate MDP model parameters given observed transitions and rewards
  - If finite set of states and actions, can just count and average counts
- Use estimated MDP to compute estimate of optimal values and policy

Will the computed values and policy converge to the true optimal values and policy in the limit of infinite data?



# REACHABILITY

- When acting randomly forever, still need to be able to **visit each state and take each action many times**
- **Want all states to be reachable from any other state**
- Quite mild assumption but doesn't always hold

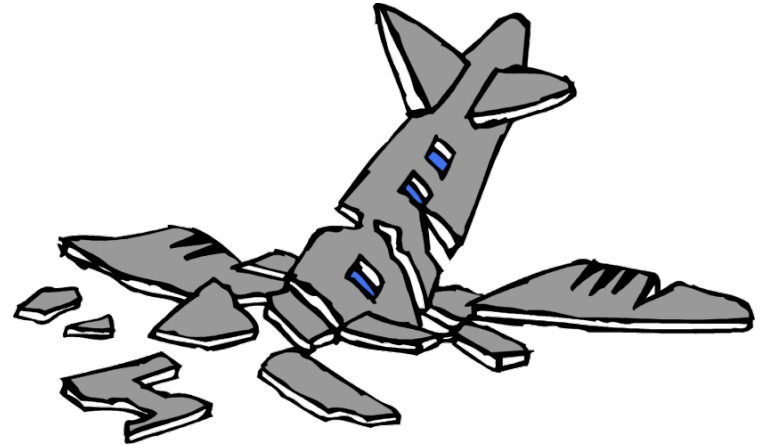


Image source: <http://ancient-heritage.blogspot.com/2014/05/crash-course-on-flying-in-face-of-logic.html>



# ON-POLICY / OFF-POLICY RL LEARNING

An Active RL agent can have two (different) policies:

- **Behavior policy** → Used to generate actions  
(↔ Interact with environment to gather sample data)
- **Learning policy** → Target action policy to learn (the “good”/optimal policy the agent eventually aims to discover through interaction)
- If Behavior policy = Learning policy → **On-policy learning**
- If Behavior policy  $\neq$  Learning policy → **Off-policy learning**



# ON-POLICY / OFF-POLICY

- **On-policy learning:** the behavior policy used to generate samples is *the same* as the target policy → Agent learns the value of the policy being used, including exploration actions.
  - The used policy is usually "soft" and non-deterministic, to ensure there is always exploration.
- **Off-policy learning:** the behavior and the target policy are different. The target policy is learned regardless (independently) of the actions chosen for exploring the environment. The agent follows a policy but learns the value of a different policy.
  - **Easier converge proofs and analysis**
  - In *Q-learning*, the action-value function  $Q(s,a)$  is learned using a **greedy** action policy, such that it directly estimates  $Q^*$



# EXAMPLE OF ON-POLICY LEARNING

- To gather sample data, the agent applies an  *$\varepsilon$ -greedy behavior policy*  $\pi_\varepsilon$ , such that with a probability  $(1-\varepsilon)=0.9$  the **optimal** action (as currently estimated from  $Q$ ) is chosen, and with probability  $\varepsilon=0.1$  a **random** action is chosen.
- The actions defined by the behavior policy are also used to update the **target** policy (i.e., to update  $Q(s,a)$  estimates)
- Since also random actions are being used, the target policy will learn the value of the policy that includes *both* optimal and random actions, selected with a probability  $\varepsilon$
- In SARSA (next slide), each  $Q(s,a)$  is updated with a sample  $Q$  value that can refer to a random action (i.e, the “value” of action  $a$  in state  $s$  is the exponential average of the observed rewards associated to both optimal and random actions)





# SARSA: ON-POLICY TD STATE-ACTION LEARNING

Initialize  $Q(s, a), \forall s \in \mathcal{S}, a \in \mathcal{A}(s)$ , arbitrarily, and  $Q(\text{terminal-state}, \cdot) = 0$

Repeat (for each episode):

Initialize  $S$

Choose  $A$  from  $S$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)

Repeat (for each step of episode):

Take action  $A$ , observe  $R, S'$

Choose  $A'$  from  $S'$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)

$$Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma Q(S', A') - Q(S, A)]$$

$S \leftarrow S'; A \leftarrow A';$

until  $S$  is terminal

or keep acting forever,  
or termination criterion

Update  $Q$  estimate with the  
sample data generated according  
to the behavior policy

# *Q-LEARNING: OFF-POLICY TD* *STATE-ACTION LEARNING*

Any  
exploration  
policy

Initialize  $Q(s, a), \forall s \in \mathcal{S}, a \in \mathcal{A}(s)$ , arbitrarily, and  $Q(\text{terminal-state}, \cdot) = 0$

Repeat (for each episode):

Initialize  $S$

Repeat (for each step of episode):

Choose  $A$  from  $S$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)

Take action  $A$ , observe  $R, S'$

$$Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma \max_a Q(S', a) - Q(S, A)]$$

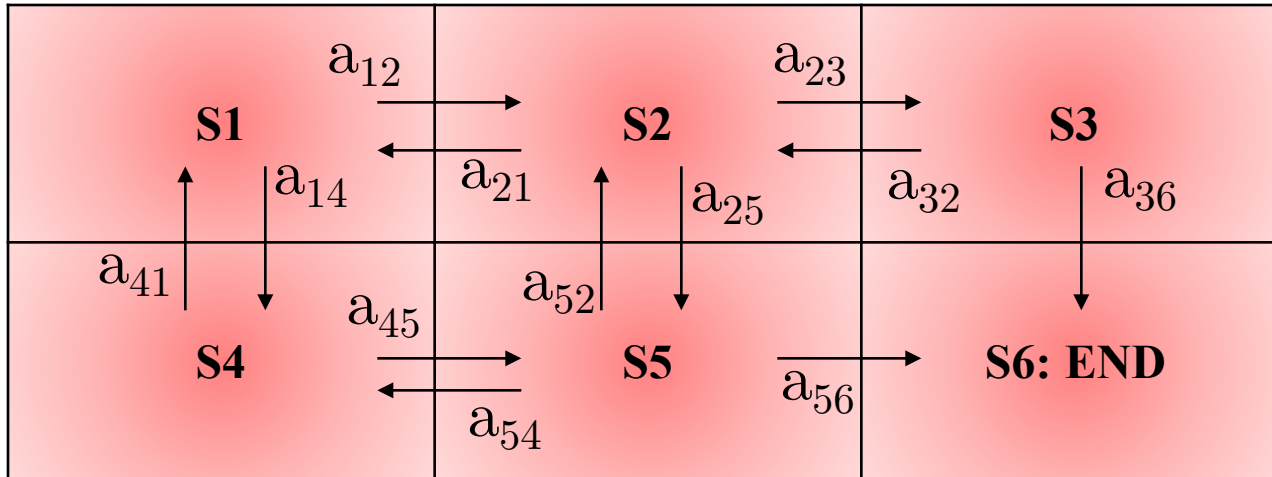
$S \leftarrow S'$

until  $S$  is terminal

or keep acting forever,  
or termination criterion

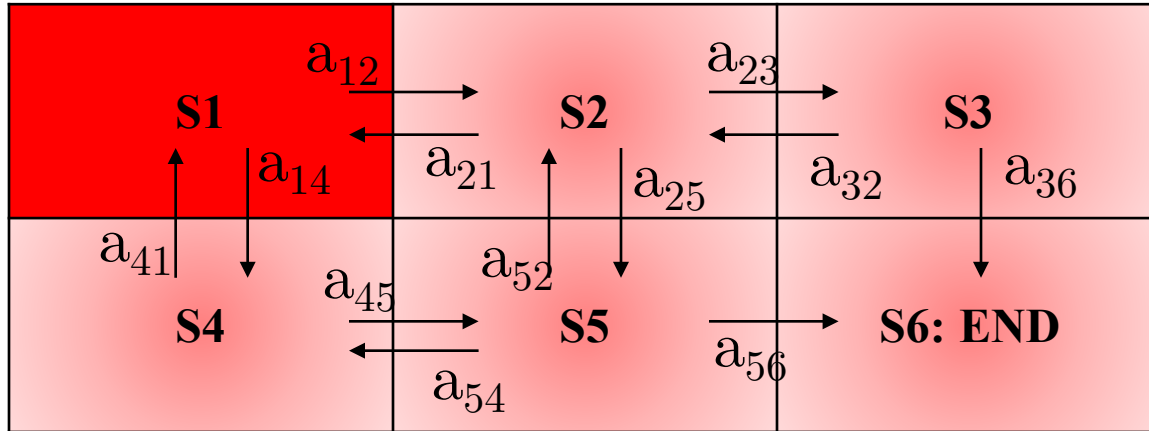
Update  $Q$  estimate with the sample data but according to a greedy policy for action selection (take the max)  $\neq$  from behavior policy

# Q-LEARNING EXAMPLE



- 6 states,  $S1, \dots, S6$
- 12 actions  $a_{ij}$  for state transitions, deterministic
- $R=100$  in  $S6$  (terminal state),  $R=0$  otherwise
- $\gamma=0.5$ ,  $\alpha = 1$
- Random behavior policy

# INITIAL STATE

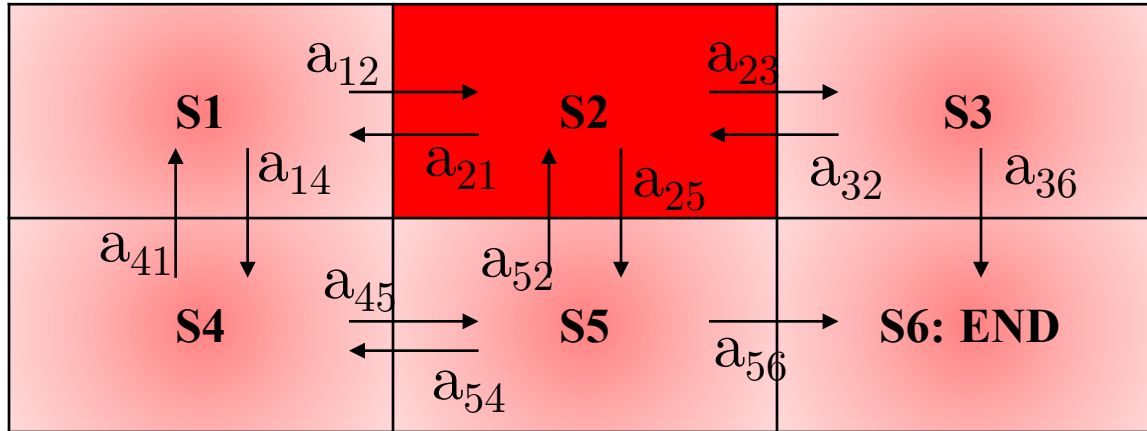


- Available actions:  $a_{12}, a_{14}$
- Choose  $a_{12}$

$Q(S1, a_{12})$	0
$Q(S1, a_{14})$	0
$Q(S2, a_{21})$	0
$Q(S2, a_{25})$	0
$Q(S2, a_{23})$	0
$Q(S3, a_{32})$	0
$Q(S3, a_{36})$	0
$Q(S4, a_{41})$	0
$Q(S4, a_{45})$	0
$Q(S5, a_{52})$	0
$Q(S5, a_{54})$	0
$Q(S5, a_{56})$	0



# NEW STATE, UPDATE



$Q(S1, a_{12})$	0
$Q(S1, a_{14})$	0
$Q(S2, a_{21})$	0
$Q(S2, a_{25})$	0
$Q(S2, a_{23})$	0
$Q(S3, a_{32})$	0
$Q(S3, a_{36})$	0
$Q(S4, a_{41})$	0
$Q(S4, a_{45})$	0
$Q(S5, a_{52})$	0
$Q(S5, a_{54})$	0
$Q(S5, a_{56})$	0

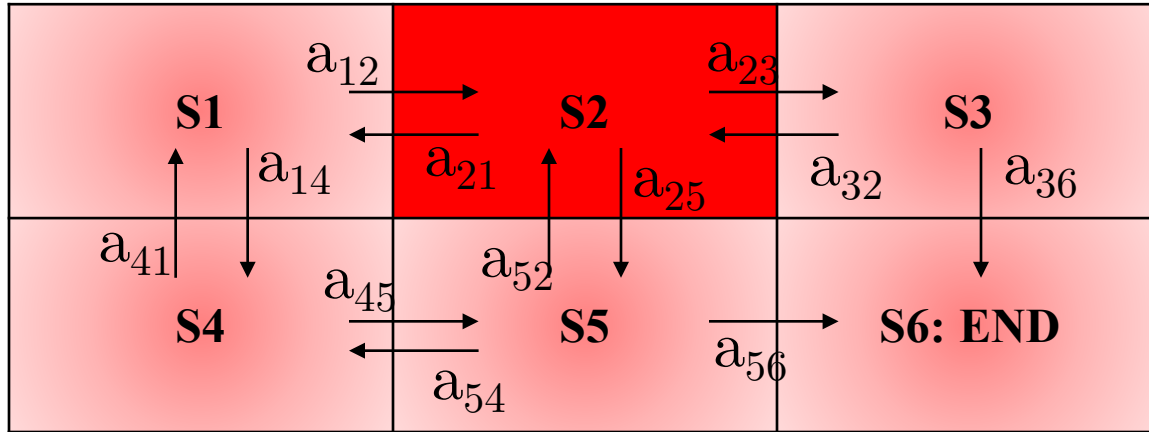
- Reward: 0; New state: S2

- Update state-action:

$$Q(S1, a_{12}) \leftarrow Q(S1, a_{12}) + (0 + 0.5 \cdot \max\{Q(S2, a_{21}), Q(S2, a_{23}), Q(S2, a_{25})\} - Q(S1, a_{12})) = 0$$



# NEW ACTION

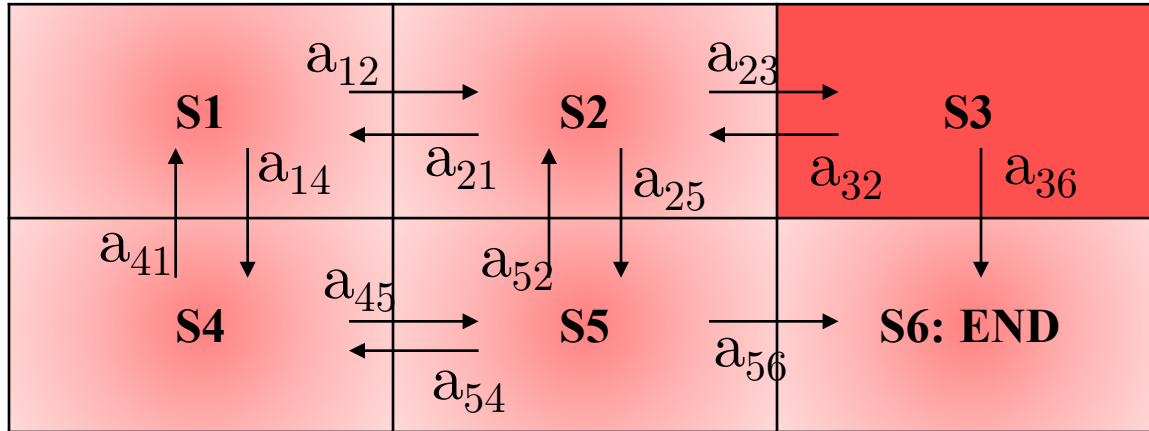


- Available actions:  $a_{21}, a_{23}, a_{25}$
- Choose  $a_{23}$

$Q(S1, a_{12})$	0
$Q(S1, a_{14})$	0
$Q(S2, a_{21})$	0
$Q(S2, a_{25})$	0
$Q(S2, a_{23})$	0
$Q(S3, a_{32})$	0
$Q(S3, a_{36})$	0
$Q(S4, a_{41})$	0
$Q(S4, a_{45})$	0
$Q(S5, a_{52})$	0
$Q(S5, a_{54})$	0
$Q(S5, a_{56})$	0



# NEW STATE, UPDATE



$Q(S1, a_{12})$	0
$Q(S1, a_{14})$	0
$Q(S2, a_{21})$	0
$Q(S2, a_{25})$	0
$Q(S2, a_{23})$	0
$Q(S3, a_{32})$	0
$Q(S3, a_{36})$	0
$Q(S4, a_{41})$	0
$Q(S4, a_{45})$	0
$Q(S5, a_{52})$	0
$Q(S5, a_{54})$	0
$Q(S5, a_{56})$	0

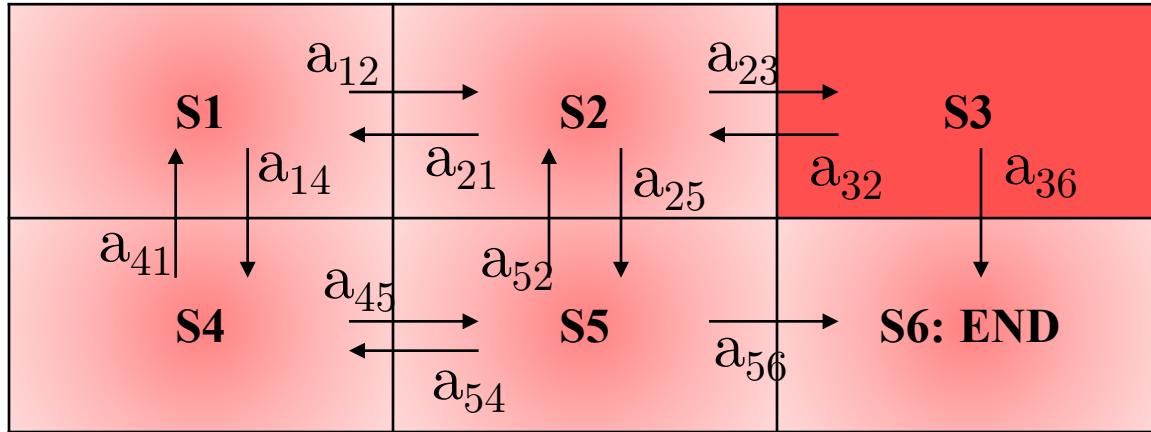
- Reward: 0; New state: S3

- Update state-action:

$$Q(S2, a_{23}) \leftarrow Q(S2, a_{23}) + (0 + 0.5 \cdot \max\{Q(S3, a_{32}), Q(S3, a_{36})\} - Q(S2, a_{23})) = 0$$



# NEW ACTION



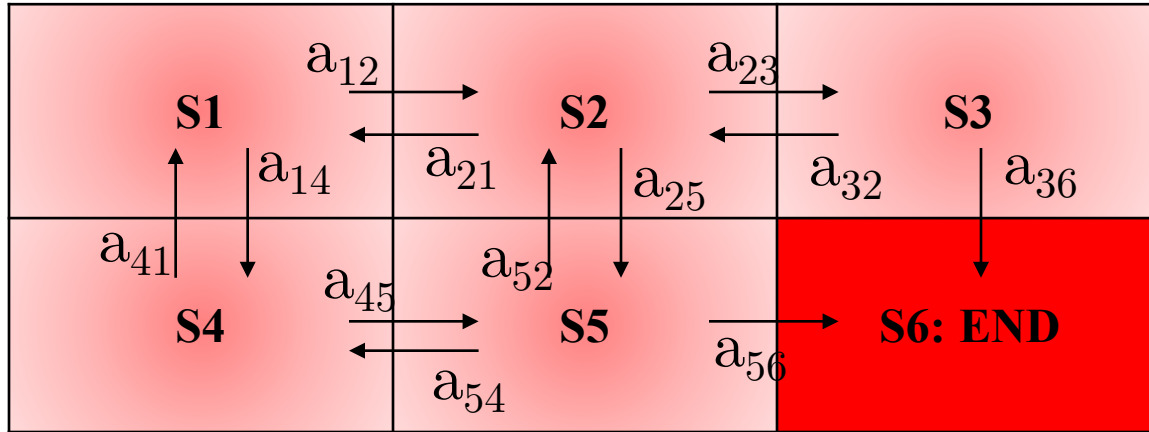
$Q(S1, a_{12})$	0
$Q(S1, a_{14})$	0
$Q(S2, a_{21})$	0
$Q(S2, a_{25})$	0
$Q(S2, a_{23})$	0
$Q(S3, a_{32})$	0
$Q(S3, a_{36})$	0
$Q(S4, a_{41})$	0
$Q(S4, a_{45})$	0
$Q(S5, a_{52})$	0
$Q(S5, a_{54})$	0
$Q(S5, a_{56})$	0

- Available actions:  $a_{36}, a_{32}$
- Choose  $a_{36}$





# NEW STATE, UPDATE



$Q(S1, a_{12})$	0
$Q(S1, a_{14})$	0
$Q(S2, a_{21})$	0
$Q(S2, a_{25})$	0
$Q(S2, a_{23})$	0
$Q(S3, a_{32})$	0
$Q(S3, a_{36})$	100
$Q(S4, a_{41})$	0
$Q(S4, a_{45})$	0
$Q(S5, a_{52})$	0
$Q(S5, a_{54})$	0
$Q(S5, a_{56})$	0

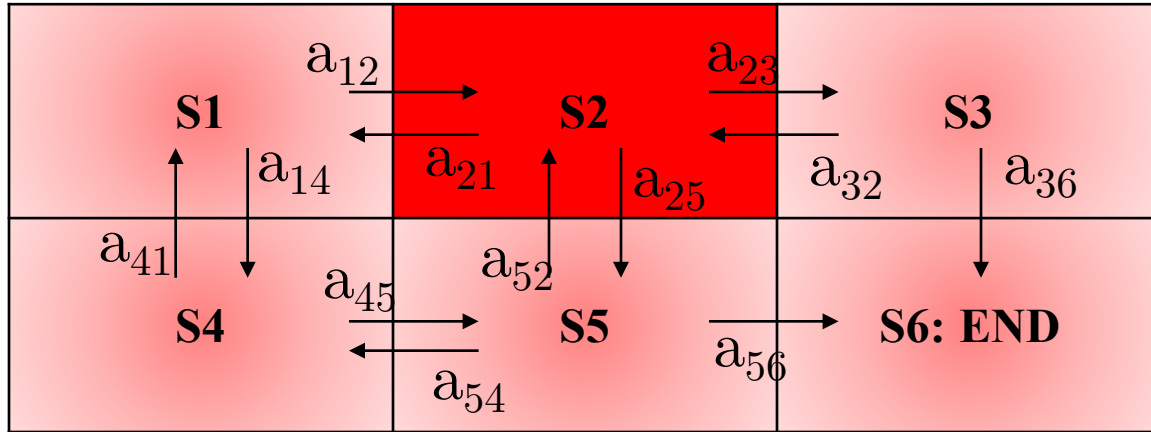
- Reward: 100; New state S6

- Update state-action:

$$Q(S3, a_{36}) \leftarrow Q(S3, a_{36}) + (100 + 0.5 \cdot \max\{\} - Q(S3, a_{36})) = 100$$



# NEW EPISODE

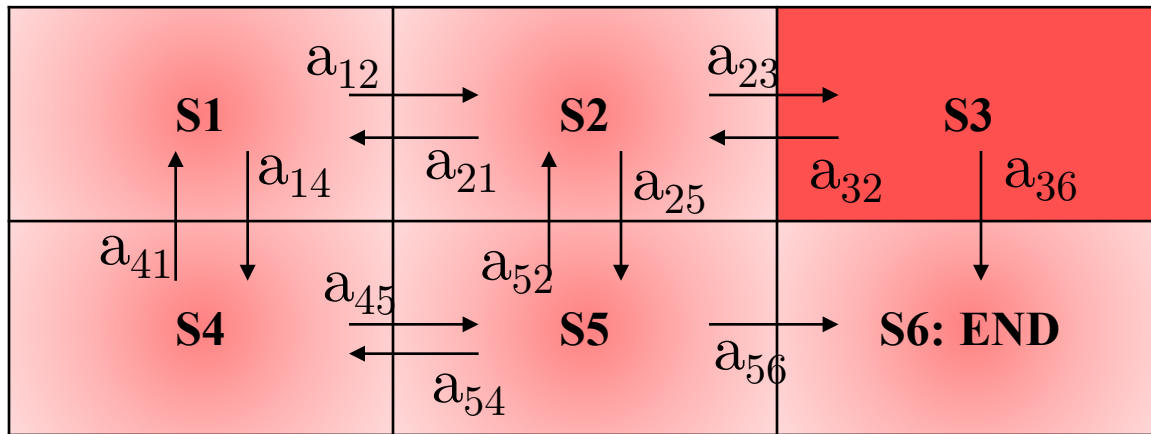


$Q(S1, a_{12})$	0
$Q(S1, a_{14})$	0
$Q(S2, a_{21})$	0
$Q(S2, a_{25})$	0
$Q(S2, a_{23})$	0
$Q(S3, a_{32})$	0
$Q(S3, a_{36})$	100
$Q(S4, a_{41})$	0
$Q(S4, a_{45})$	0
$Q(S5, a_{52})$	0
$Q(S5, a_{54})$	0
$Q(S5, a_{56})$	0

- Available actions:  $a_{21}, a_{23}, a_{25}$
- Choose  $a_{23}$



# NEW STATE, UPDATE



$Q(S1, a_{12})$	0
$Q(S1, a_{14})$	0
$Q(S2, a_{21})$	0
$Q(S2, a_{25})$	0
$Q(S2, a_{23})$	50
$Q(S3, a_{32})$	0
$Q(S3, a_{36})$	100
$Q(S4, a_{41})$	0
$Q(S4, a_{45})$	0
$Q(S5, a_{52})$	0
$Q(S5, a_{54})$	0
$Q(S5, a_{56})$	0

- Reward: 0; New state: S3

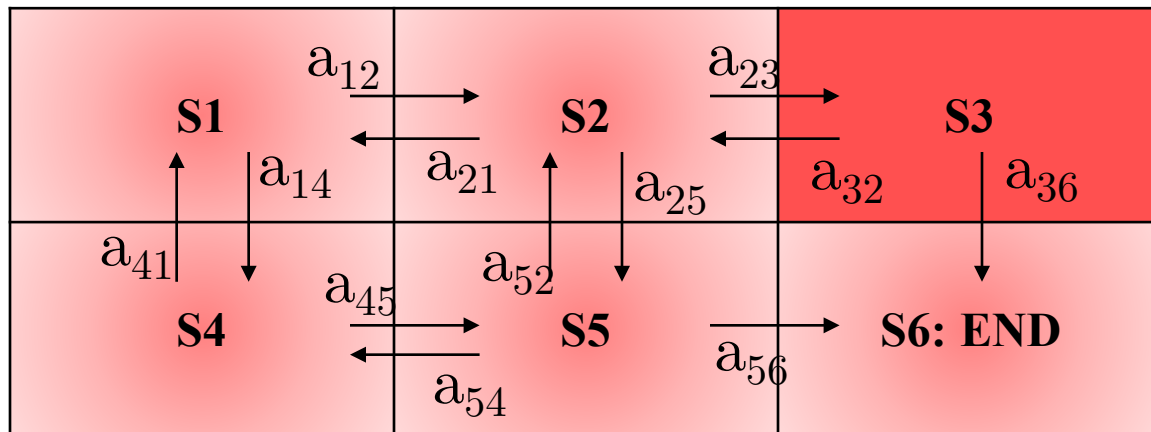
- Update state-action:

$$Q(S2, a_{23}) \leftarrow Q(S2, a_{23}) + (0 + 0.5 \cdot \max\{Q(S3, a_{32}), Q(S3, a_{36})\} - Q(S2, a_{23})) = 50$$



# AFTER MANY EPISODES ...

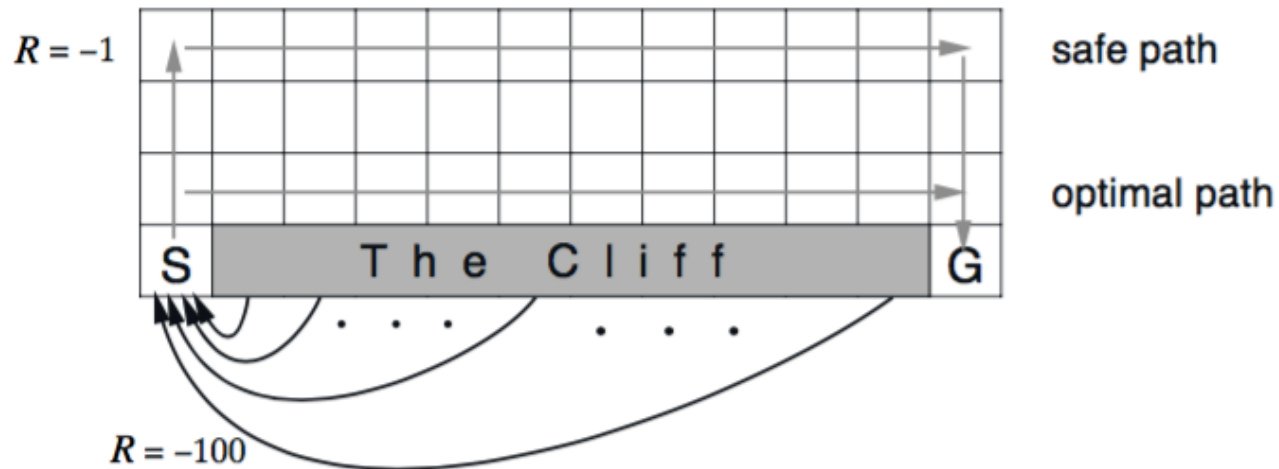
The optimal Q-values for the discount factor  $\gamma=0.5$



$Q(S1, a_{12})$	25
$Q(S1, a_{14})$	25
$Q(S2, a_{21})$	12.5
$Q(S2, a_{23})$	25
$Q(S2, a_{25})$	50
$Q(S3, a_{32})$	25
$Q(S3, a_{36})$	100
$Q(S4, a_{41})$	12.5
$Q(S4, a_{45})$	50
$Q(S5, a_{52})$	25
$Q(S5, a_{54})$	25
$Q(S5, a_{56})$	100

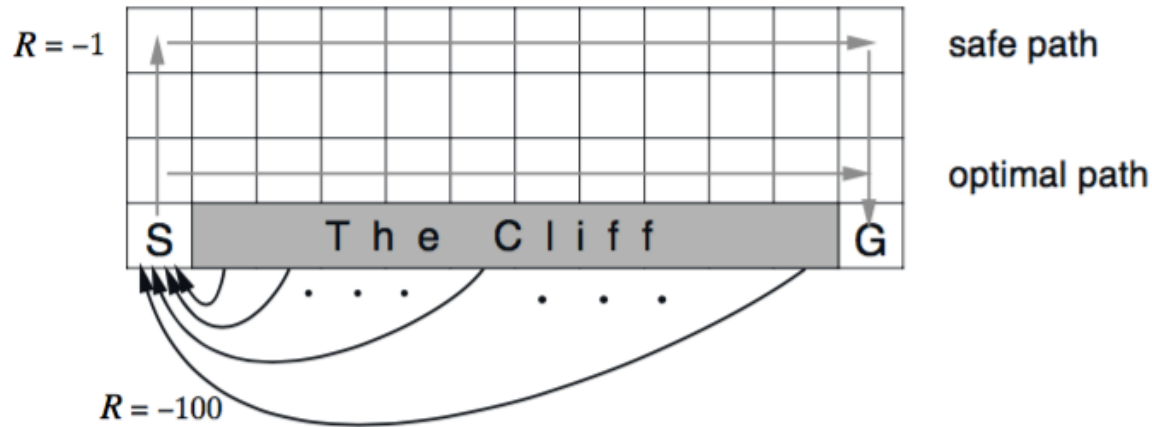


# SARSA vs. Q-LEARNING



- Episodic task in a grid world with a start state  $S$  and a goal state  $G$
- $R = -1$  on all state transitions except when stepping on the *Cliff*, where  $R = -100$  and the agent is sent back to  $S$
- Up, down left, right deterministic actions are available at each state
- Both algorithms employ an  $\epsilon$ -greedy behavior policy,  $\epsilon = 0.1$
- The optimal policy is that moving along the edge of the cliff

# SARSA VS. Q-LEARNING

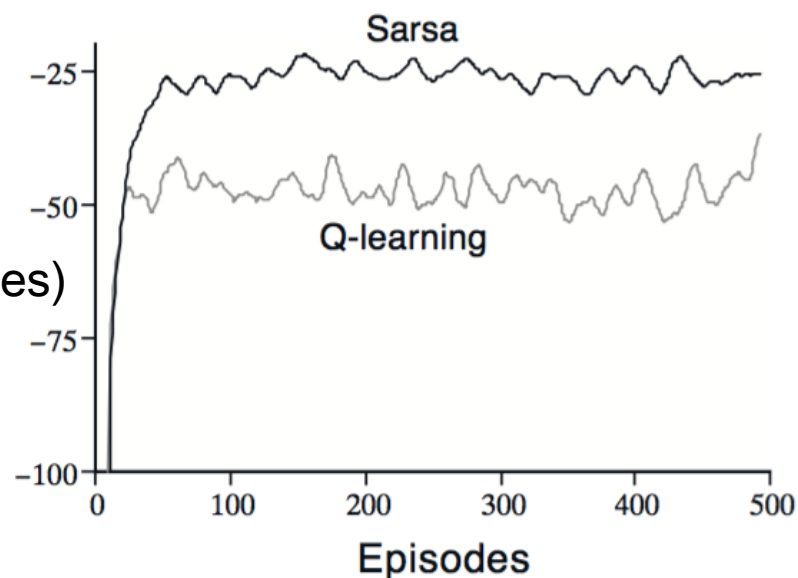


- **Q-learning** learns the optimal policy, however, the *online* (i.e., while learning) performance is not that good since the agent occasionally falls off the cliff because of the  $\epsilon$ -greedy action selection
- **SARSA** does not learn the optimal policy, but rather the action values that accounts for  $\epsilon$ -greedy action selections, and therefore it learns the longer but safer path that minimizes the risk of  $R = -100$



# SARSA VS. Q-LEARNING

Sum of rewards during episode (avg over 10 episodes)



After an initial transient, Q-learning learns the optimal Q-values, corresponding to traveling along the edge of the cliff

Q-learning *online* performance is lower than SARSA because it occasionally falls into the cliff, while rarely happens for SARSA

# Q-LEARNING PROPERTIES

- Acting randomly, Q-learning asymptotically converges to optimal state-action values → Finds optimal policy
- More in general, any behavior policy that guarantees that *all state-action pairs continue to be updated* guarantees asymptotic convergence
- In stochastic environments, optimality also requires that the **learning rate**  $\alpha$  decreases to zero over time. In deterministic environments a constant  $\alpha=1$  suffices
- We are “free” to choose the behavior policy that is best to perform an effective exploration of the environment





# TOWARDS GATHERING HIGH REWARDS

- Fortunately, acting randomly is sufficient, but not necessary, to learn the optimal values and policy

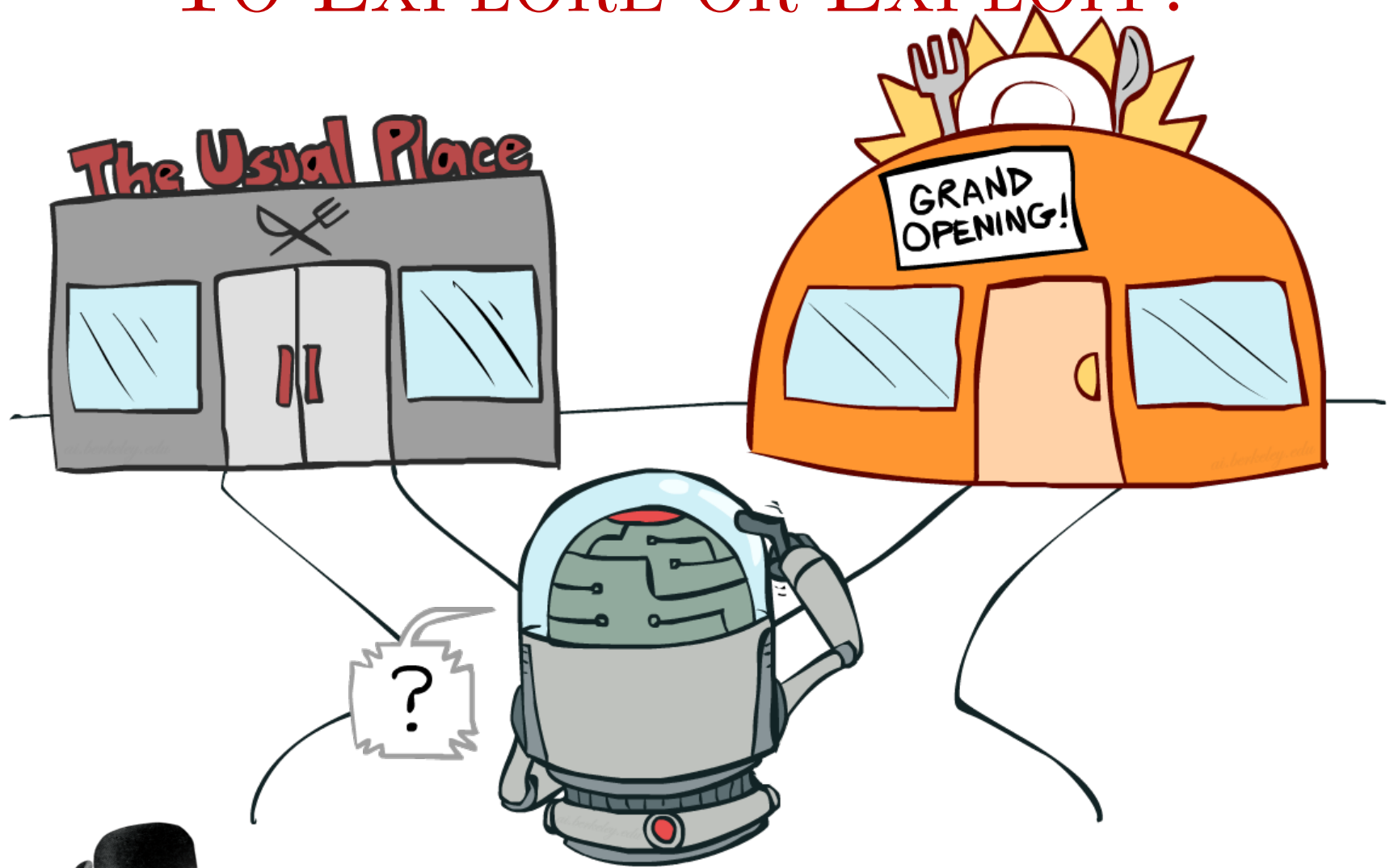


# HOW TO ACT?

- Initialize  $s$  to a starting state
- Initialize  $Q(s,a)$  values
- For  $t=1,2,\dots$ 
  - Choose  $a = \operatorname{argmax} Q(s,a)$
  - Observe  $s', r(s,a,s')$
  - Update/Compute  $Q$  values (using model-based or Q-learning approach)
- Is this guaranteed to find the optimal policy?
  1. Yes
  2. No
  3. Not sure



# TO EXPLORE OR EXPLOIT?



# SIMPLE APPROACH: $\epsilon$ -GREEDY

- With probability  $1-\epsilon$ 
  - Choose  $\operatorname{argmax}_a Q(s,a)$
- With probability  $\epsilon$ 
  - Select random action
  
- Guaranteed to compute optimal policy
- But even after millions of steps still won't always be following policy computed (the  $\operatorname{argmax} Q(s,a)$ )



# GREEDY IN LIMIT OF INFINITE EXPLORATION (GLIE)

- $\epsilon$ -greedy approach
- But decay  $\epsilon$  over time
- Eventually will be following optimal policy almost all the time



# PERFORMANCE EVALUATION

How should we evaluate the performance of an algorithm?

- Computational efficiency
- How much reward gathered under algorithm?



# OBJECTIVES FOR RL ALGORITHMS

- **Asymptotic guarantees**
  - In *limit* converge to a policy identical to the optimal policy if knew unknown model parameters
  - Q-learning! (under what conditions?)
- **Probably Approximately Correct (PAC)**
  - On all but *finite* number of samples, choose action whose expected reward is close to expected reward of action take if knew model parameters
  - E<sup>3</sup> (Kearns & Singh), R-MAX (Brafman & Tennenholtz)



# SAMPLE COMPLEXITY AND PAC LEARNING

- **Sample complexity** of exploration of a reinforcement-learning algorithm: the number of steps in which the algorithm does *not* choose near-optimal actions
- **A PAC algorithm** is such that the sample complexity is bounded by some function **polynomial** in the states, action, and error quantities, with a certain (high) **probability**





# GENERAL IDEA OF PAC LEARNING

- **Give up (a bit) being exact**
  - We may spare the learner from learning the exact function/model/concept, an approximation would be fine.
- **Allow a small probability that the learning algorithm fails**
  - This would be in some worst cases (e.g., samples are badly chosen). Try to reduce the probability of getting such unlucky examples by considering sufficient #examples
  - A PAC learning strategy places an upper bound on the probability of committing an error by placing a min bound on the number of examples/samples it takes to learn a target



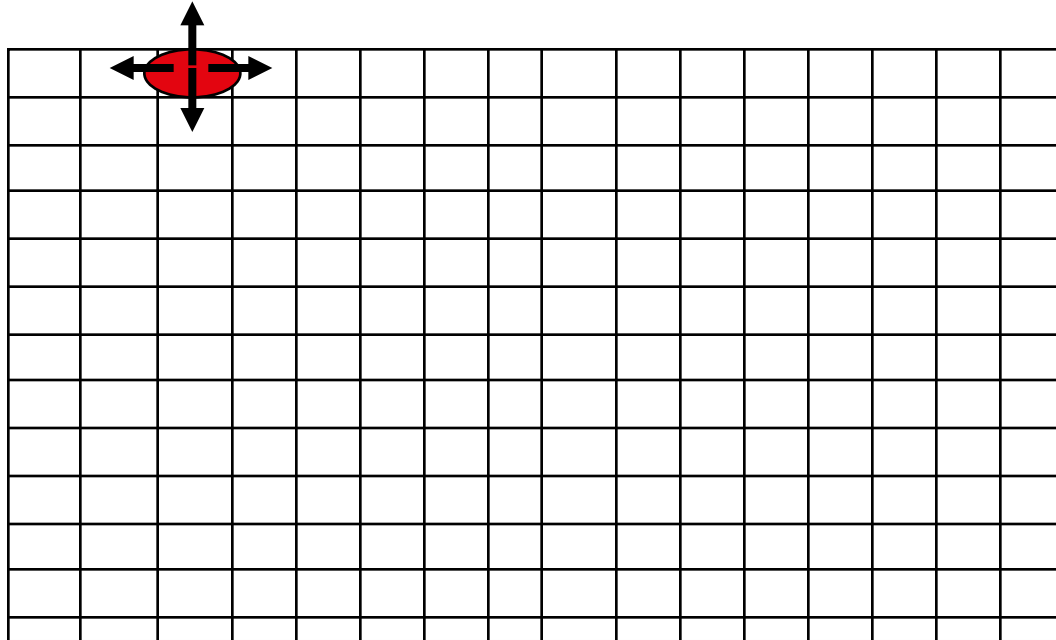
# RECALL MODEL-BASED RL

- Given data seen so far
- Build an explicit model of the MDP
- Compute policy for it
- Select action for current state given policy, observe next state and reward
- Repeat



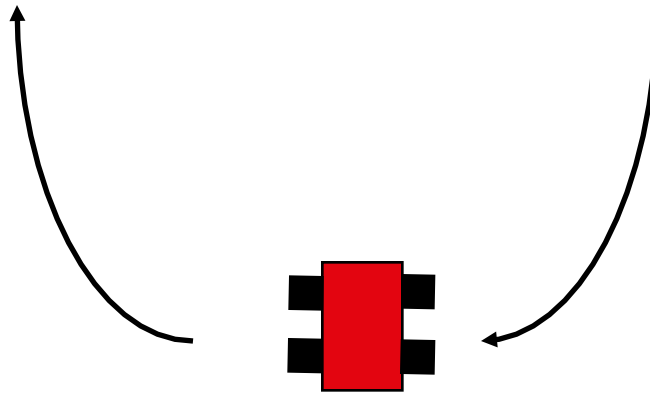
# R-MAX (BRAFMAN & TENNENHOLTZ)

S1 S2 ...



# R-MAX IS MODEL-BASED RL

Think hard: estimate models & compute policies



Act in world

**R-max leverages optimism under uncertainty!**

If I don't know about  $s$ , I assume I can get a good reward there  $\rightarrow$  I will go there  $\rightarrow$  I will explore it!



# R-MAX ALGORITHM

- The key idea in R-max is to track if we have enough experience to **trust model estimates** for rewards,  $R$ , and dynamics,  $T$
- Set up  $T$  and  $R$  models in order to favor exploration as long as no enough evidence has been gathered
- Initially mark all state-action pairs as *Unknown*, since no experience has been gathered yet
- All unknown dynamics are *self loops*
- Set reward to *R-max* to encourage **exploration**
- **Transition counts are 0**  $\rightarrow$  **Transition dynamics is flat**



# R-MAX ALGORITHM:

## INITIALIZE: DEFINE “KNOWN” MDP

(S,A)

Known/  
Unknown

	S1	S2	S3	S4	...
↑	U	U	U	U	
→	U	U	U	U	
↓	U	U	U	U	
←	U	U	U	U	

T

Transition  
Counts

	S1	S2	S3	S4	...
↑	0	0	0	0	
→	0	0	0	0	
↓	0	0	0	0	
←	0	0	0	0	

R

Reward

	S1	S2	S3	S4	...
↑	$R_{max}$	$R_{max}$	$R_{max}$	$R_{max}$	
→	$R_{max}$	$R_{max}$	$R_{max}$	$R_{max}$	
↓	$R_{max}$	$R_{max}$	$R_{max}$	$R_{max}$	
←	$R_{max}$	$R_{max}$	$R_{max}$	$R_{max}$	

**In the “known” MDP,  
any unknown (s,a) pair  
has its dynamics set as  
a self loop &  
reward =  $R_{max}$**

# R-MAX ALGORITHM

Plan in known MDP



# R-MAX: PLANNING

- Compute optimal policy  $\pi_{\text{known}}$  for “known” MDP





# EXERCISE: WHAT WILL INITIAL VALUE OF $Q(s,a)$ BE FOR EACH $(s,a)$ PAIR IN THE KNOWN MDP?

WHAT IS THE POLICY? Reward

Known/  
Unknown

	S1	S2	S3	S4	...
↑	U	U	U	U	
→	U	U	U	U	
↓	U	U	U	U	
←	U	U	U	U	

	S1	S2	S3	S4	...
↑	$R_{max}$	$R_{max}$	$R_{max}$	$R_{max}$	
→	$R_{max}$	$R_{max}$	$R_{max}$	$R_{max}$	
↓	$R_{max}$	$R_{max}$	$R_{max}$	$R_{max}$	
←	$R_{max}$	$R_{max}$	$R_{max}$	$R_{max}$	

Transition  
Counts

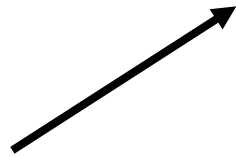
	S1	S2	S3	S4	...
↑	0	0	0	0	
→	0	0	0	0	
↓	0	0	0	0	
←	0	0	0	0	

**In the “known” MDP,  
any unknown  $(s,a)$  pair  
has its dynamics set as  
a self loop &  
reward =  $R_{max}$**



# R-MAX ALGORITHM

Act using  
policy

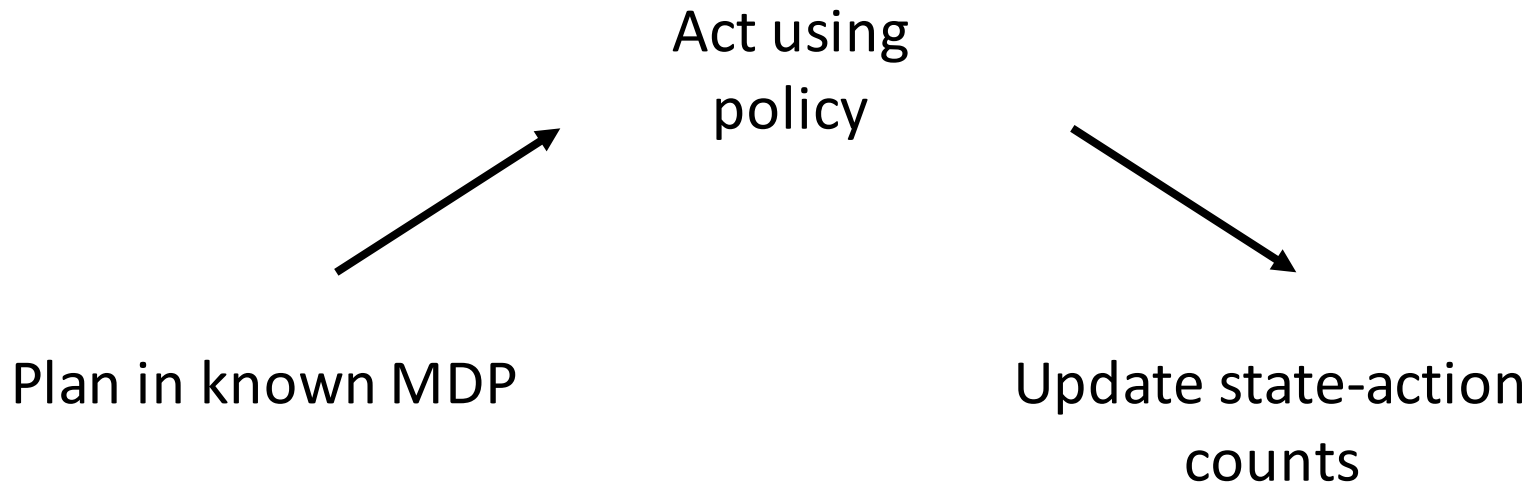


Plan in known MDP

- Given optimal policy  $\pi_{\text{known}}$  for “known” MDP
- Take best action for current state  $\pi_{\text{known}}(s)$ , transition to new state  $s'$  and get reward  $r$



# R-MAX ALGORITHM



# UPDATE KNOWN MDP

Reward

Known/  
Unknown

	S2	S2	S3	S4	...
↑	U	U	U	U	
→	U	U	U	U	
↓	U	U	U	U	
←	U	U	U	U	

	S2	S2	S3	S4	...
↑	$R_{\max}$	$R_{\max}$	$R_{\max}$	$R_{\max}$	
→	$R_{\max}$	$R_{\max}$	$R_{\max}$	$R_{\max}$	
↓	$R_{\max}$	$R_{\max}$	$R_{\max}$	$R_{\max}$	
←	$R_{\max}$	$R_{\max}$	$R_{\max}$	$R_{\max}$	

Transition  
Counts

	S2	S2	S3	S4	...
↑	0	0	0	0	
→	0	0	1	0	
↓	0	0	0	0	
←	0	0	0	0	

Increment counts for  
state-action tuple



# UPDATE KNOWN MDP

Known/  
Unknown

	S2	S2	S3	S4	...
↑	U	U	U	U	
→	U	U	K	U	
↓	U	U	U	U	
←	U	U	U	U	

Reward

	S2	S2	S3	S4	...
↑	$R_{\max}$	$R_{\max}$	$R_{\max}$	$R_{\max}$	
→	$R_{\max}$	$R_{\max}$	$R$	$R_{\max}$	
↓	$R_{\max}$	$R_{\max}$	$R_{\max}$	$R_{\max}$	
←	$R_{\max}$	$R_{\max}$	$R_{\max}$	$R_{\max}$	

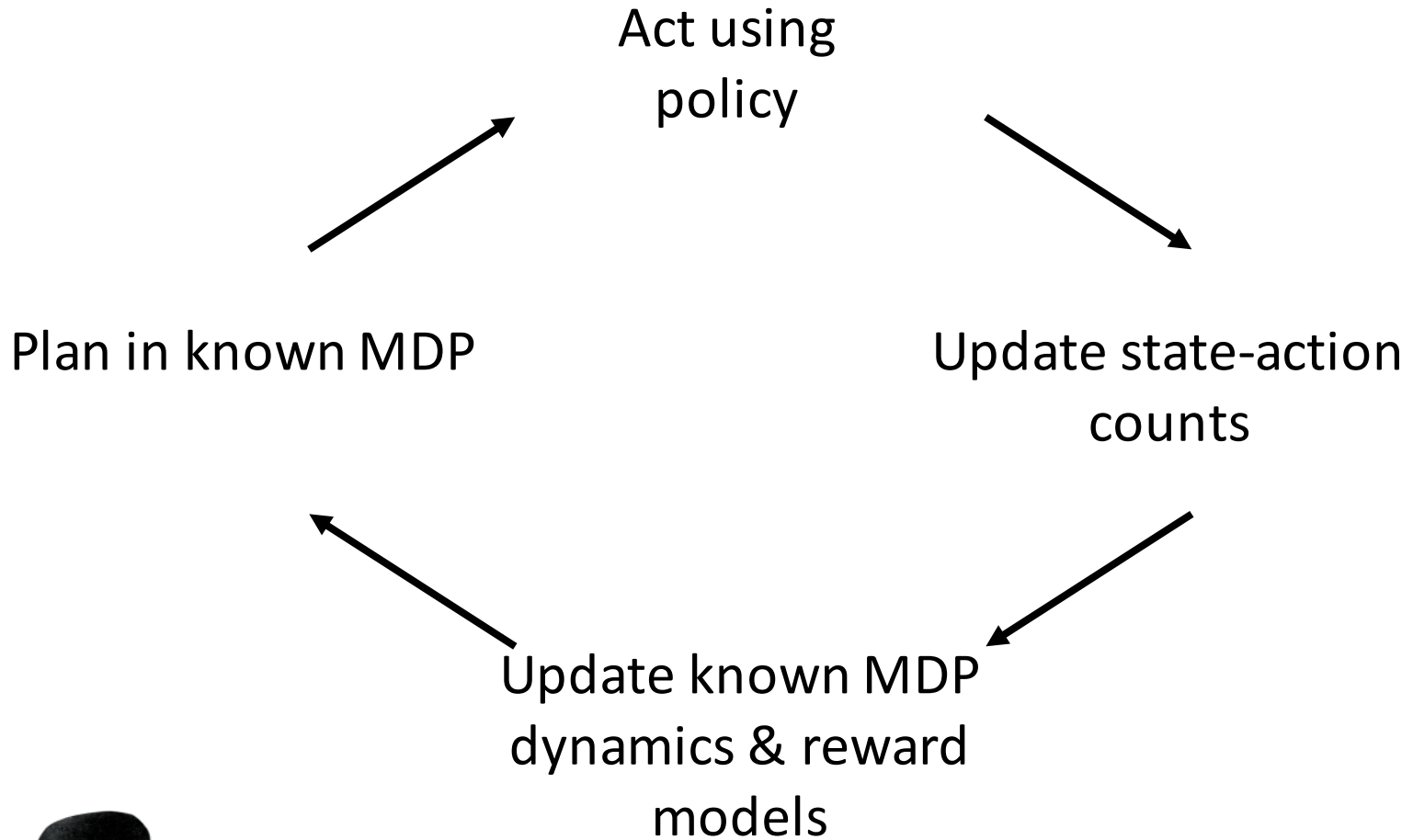
Transition  
Counts

	S2	S2	S3	S4	...
↑	3	3	4	3	
→	2	4	5	0	
↓	4	0	4	4	
←	2	2	4	1	

If counts for  $(s,a) > N$ ,  
 $(s,a)$  becomes *known*:  
**use observed data to  
estimate transition &  
reward model** for  $(s,a)$   
when planning



# R-MAX ALGORITHM



# ESTIMATING MDP MODEL FOR A (S,A) PAIR GIVEN DATA

- Transition model estimation
- Reward model estimation



# SAMPLE COMPLEXITY OF R-MAX

$$\tilde{O}\left(\frac{SA}{\varepsilon(1-\gamma)^2} \underbrace{\frac{S}{\varepsilon^2(1-\gamma)^4}}_{\substack{\text{\# samples need} \\ \text{per (s,a) pair}}}\right)$$

**PAC guarantees:** On all but the above number of steps, chooses action whose expected reward is close to expected reward of action take if knew model parameters, with high probability





# Common Idea of Model-Based PAC RL Algorithms:

Quantify how **many samples** do we need to build a **good model** that we can use to **act well** in the world?



# “GOOD” RL MODELS

- Estimate model parameters from experience
- More experience means our estimated model parameters will closer be to the true unknown parameters, with high probability



# ACTING WELL IN THE WORLD

$p(s' | s, a)$  known  $\rightarrow$  Compute  $\epsilon$ -optimal policy

Bound  $(\tilde{p}(s' | s, a) - p(s' | s, a)) \rightarrow$  Bound error in policy calculated using  $\tilde{p}(s' | s, a)$



# HOW MANY SAMPLES DO WE NEED TO BUILD A GOOD MODEL THAT WE CAN USE TO ACT WELL IN THE WORLD?

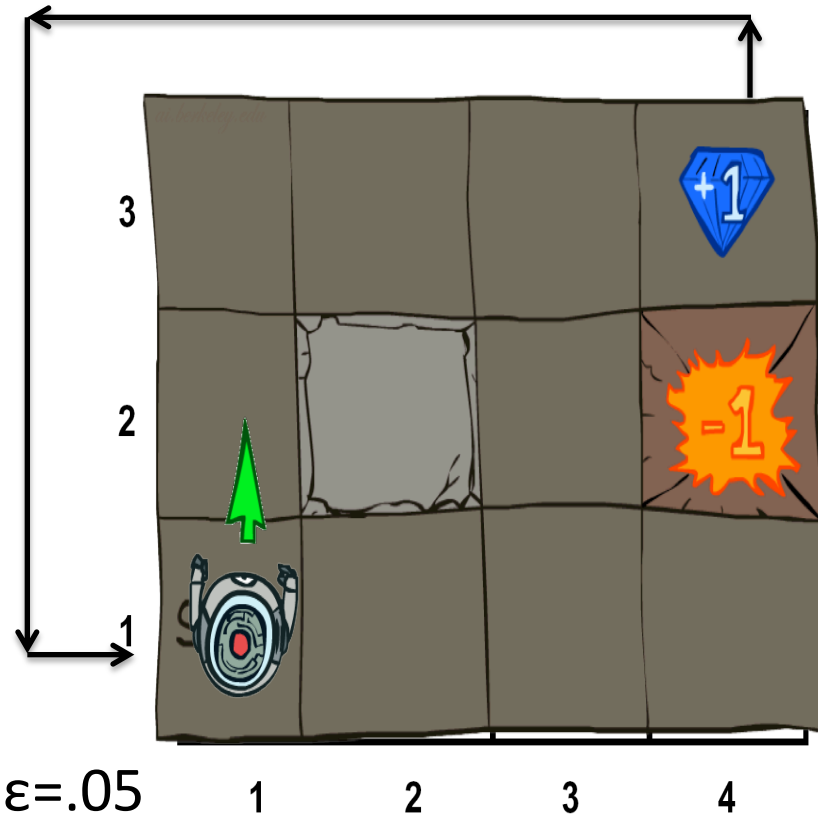
Sample complexity = # steps on which may not act well (could be far from optimal)

(R-MAX and  $E^3$ ) = Poly( # of states)



# IS THIS A SMALL NUMBER OF STEPS?

$$\tilde{O}\left(\frac{SA}{\varepsilon(1-\gamma)^2} \frac{S}{\varepsilon^2(1-\gamma)^4}\right)$$



$\gamma=.99, \varepsilon=.05$



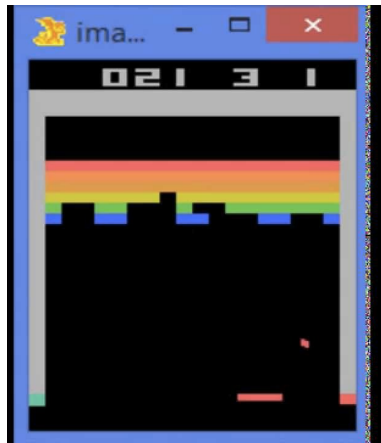
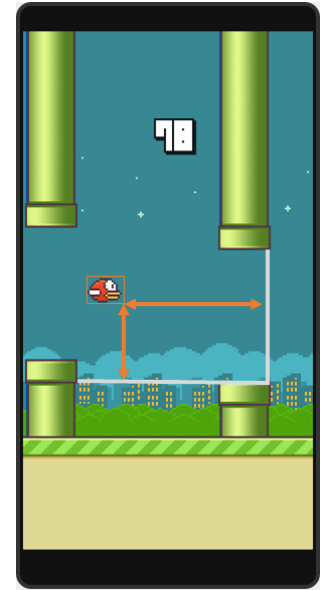
# WHAT YOU SHOULD KNOW ABOUT RL

- Define
  - Exploration vs exploitation tradeoff
  - Model free and model based RL
- Implement Q-learning and R-max
- Describe evaluation criteria for RL algorithms
  - empirical performance, computational complexity sample/data efficiency
  - Contrast strengths & weaknesses of Q-learning vs R-max in terms of these criteria. Give examples of where each might be preferred



# SOME (INTERESTING) RL FUN

- RL in the “real” world: Flappy Bird game
- <http://sarvagyaavaish.github.io/FlappyBirdRL/>



- DeepMind’s Deep Q-learning playing Atari games at human expert level
- <https://www.youtube.com/watch?v=V1eYniJ0Rnk>

