

# Reinforcement learning

Episode 2

## Value-based methods, Temporal Difference



Yandex  
Data Factory

LAMBDA



**British Hedgehog  
Preservation Society**

# Reinforcement learning

Episode 2

Oh gosh,  
we have 101 slides!



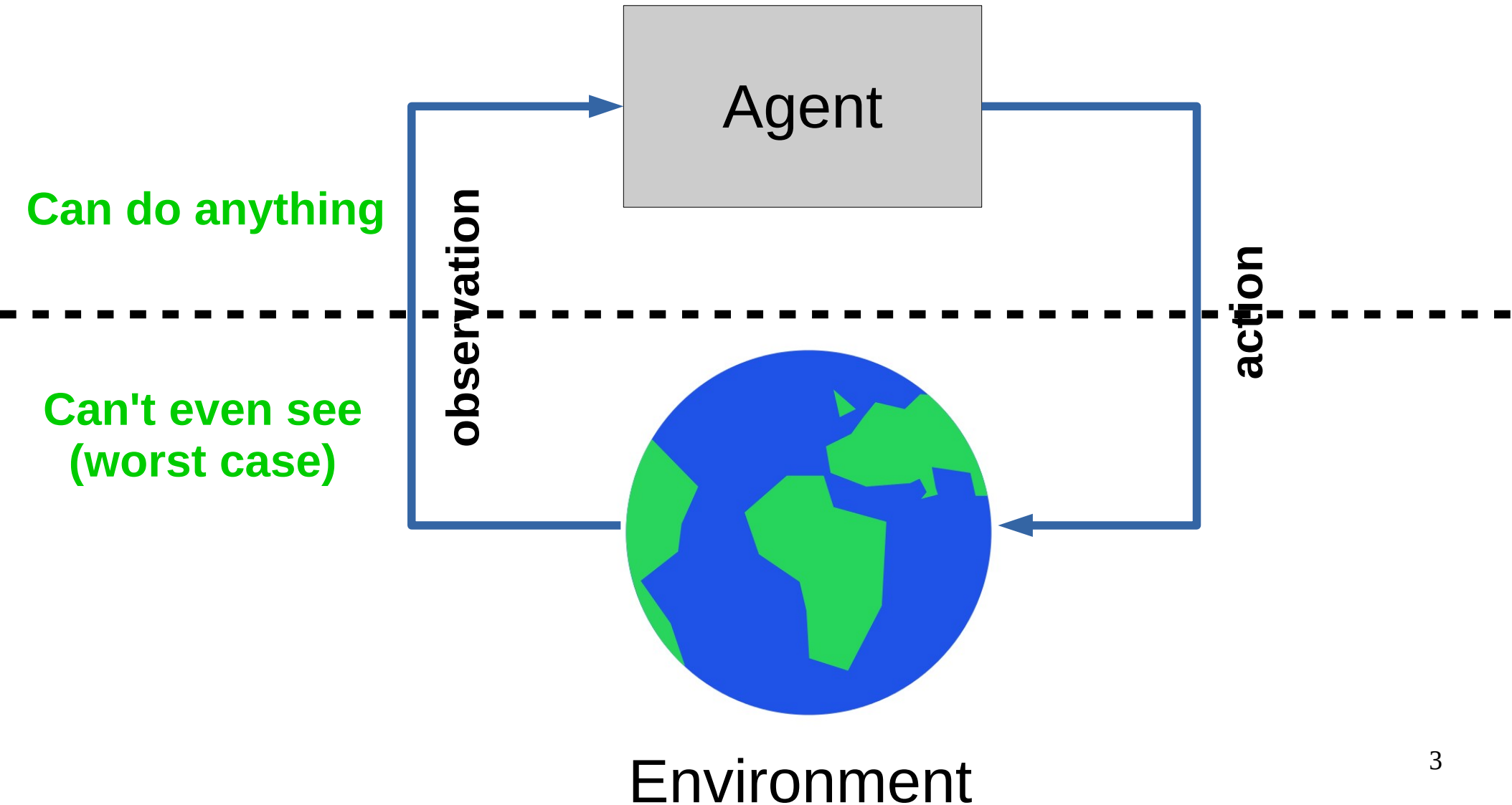
Yandex  
Data Factory

LAMBDA 

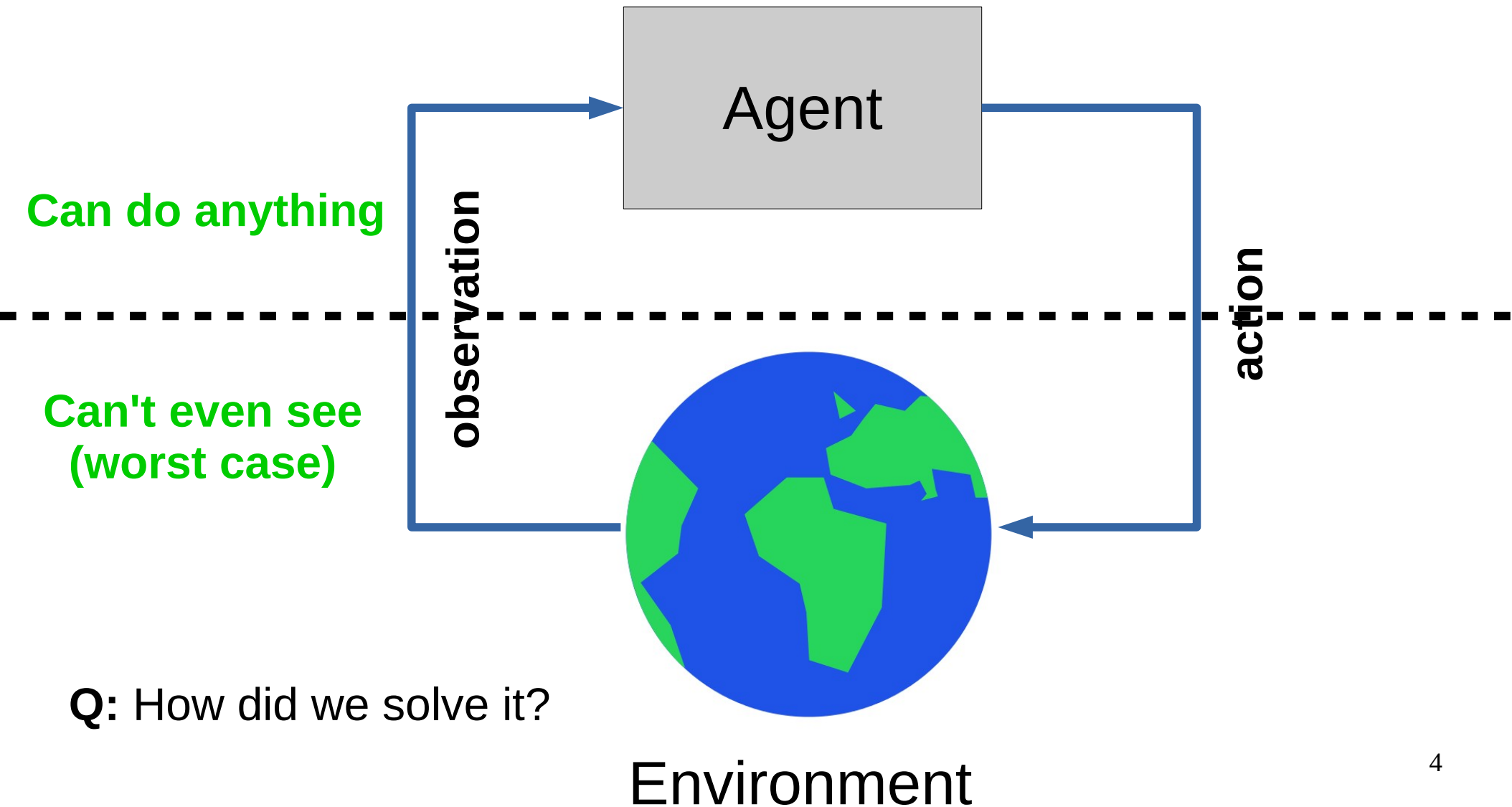


**British Hedgehog  
Preservation Society**

# Recap: reinforcement learning



# Recap: reinforcement learning



# Black box methods

- Genetic algorithms
- Evolution strategies
- Crossentropy method
- ...

# Black box drawbacks

- Both need a full session to start learning
- Requires a lot of interaction
  - A lot of crashed robots / simulations



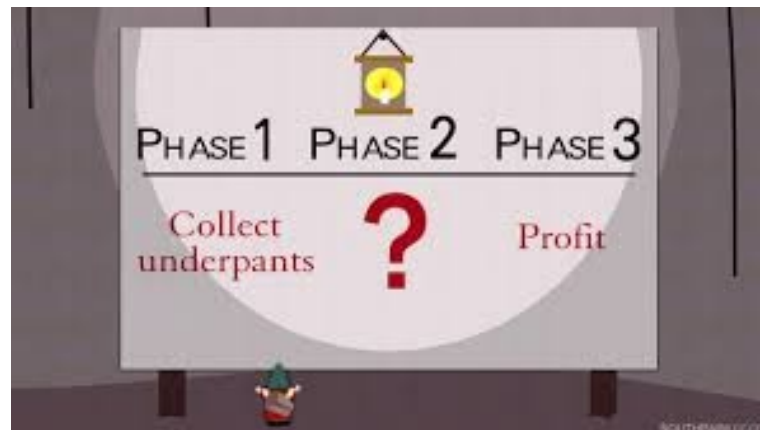
# Today: value-based methods

- Core idea:
  1. You are at state  $s$
  2. Compute expected reward for session if you take each possible action ( $a_1, a_2, a_3, \dots$ )

Then what? (e.g. chess)

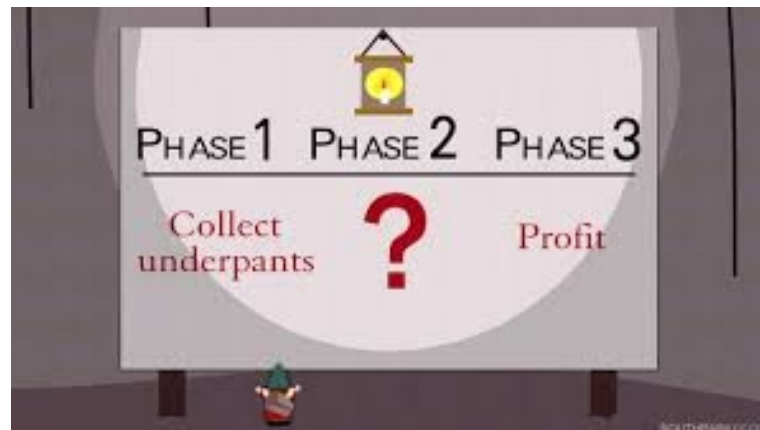
You can compute expected win rate for each move.

What do you do?



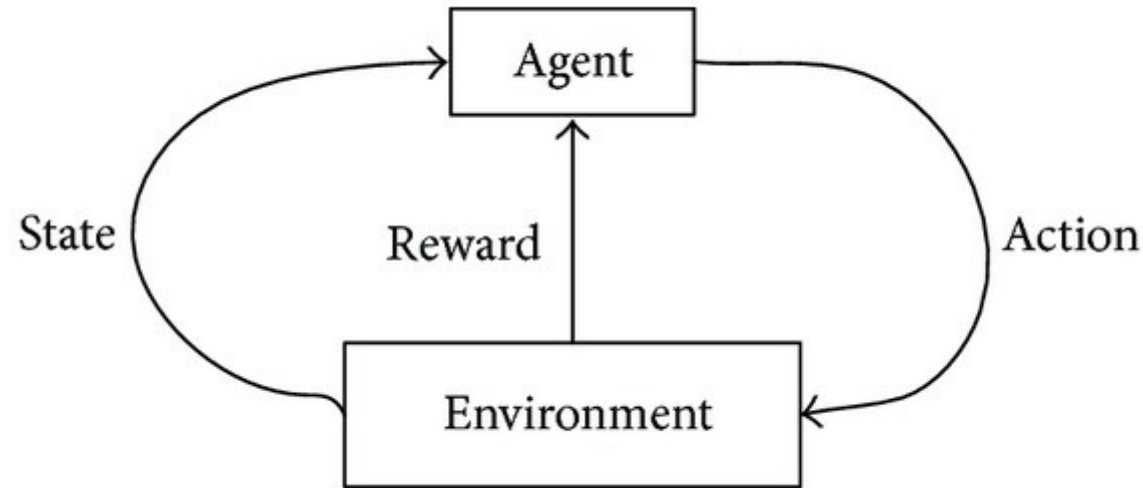
# Today: value-based methods

- Core idea:
  1. You are at state  $s$
  2. Compute expected reward for session if you take each possible action ( $a_1, a_2, a_3, \dots$ )
  3. Take action with highest expected reward!





# MDP formalism: reward on each tick



Classic MDP(Markov Decision Process)

Agent interacts with environment

- Environment states:  $s \in S$
- Agent actions:  $a \in A$
- State transition:  $P(s_{t+1}|s_t, a_t)$
- Reward:  $r_t = r(s_t, a_t)$

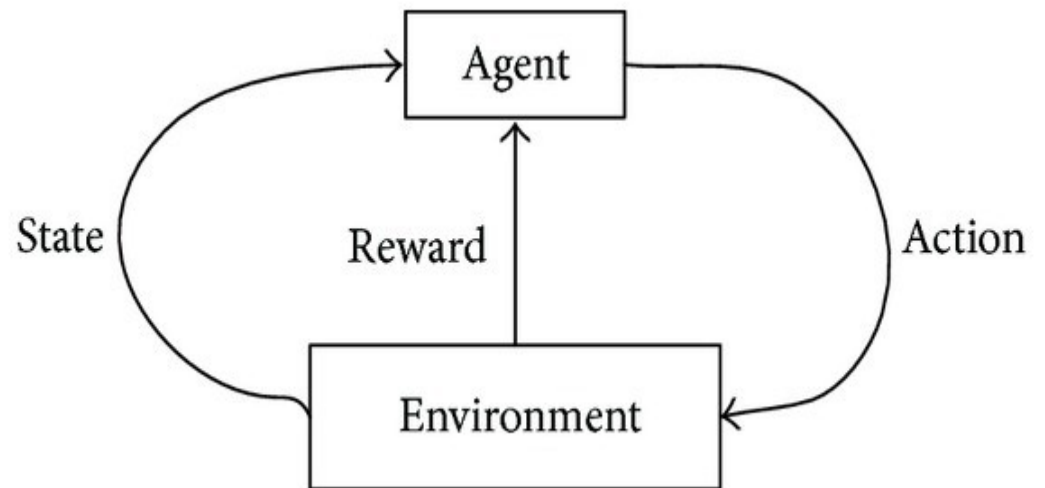
# Model-based setup

## What we know

- State transitions

$$P(s_{next}|s, a) \quad \text{or} \quad s_{next} = T(s, a)$$

- Rewards  $r(s, a)$



# Model-based setup

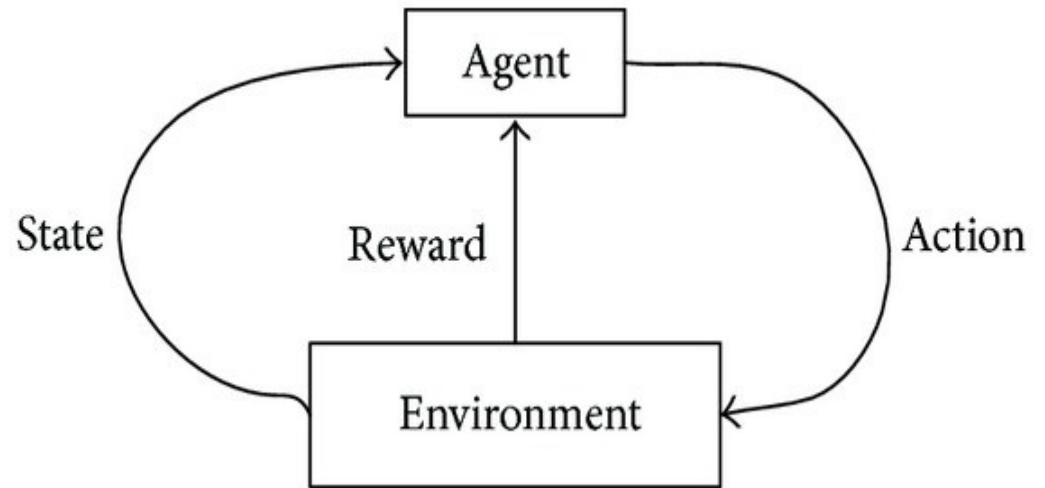
## What we know

- State transitions

$$P(s_{next}|s, a) \quad \text{or} \quad s_{next} = T(s, a)$$

Weaker version: we can only  
sample from  $P(s'|s, a)$

- Rewards  $r(s, a)$



# Discounted reward MDP



Objective:

**Discounted return  $G$**

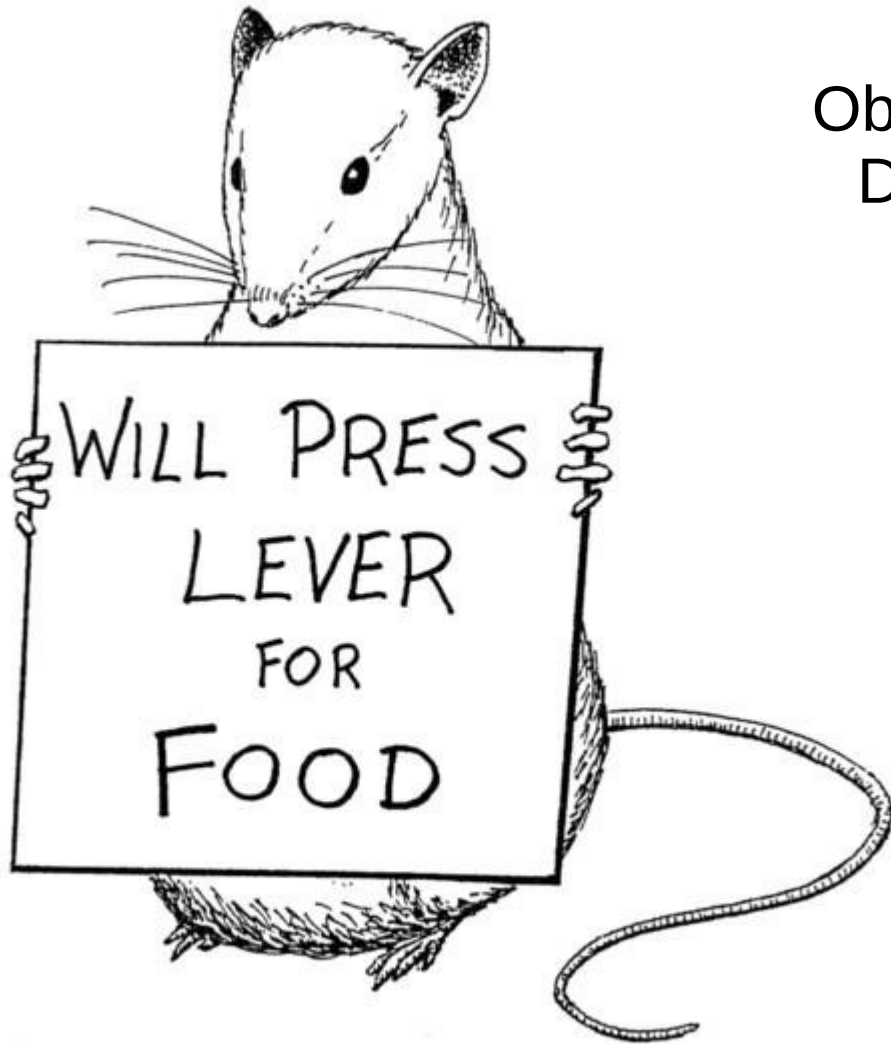
$$G_t = r_t + \gamma \cdot r_{t+1} + \gamma^2 \cdot r_{t+2} + \dots + \gamma^n \cdot r_{t+n}$$

$$G_t = \sum_i \gamma^i \cdot r_{t+i} \quad \gamma \in (0,1) \text{ const}$$

$\gamma \sim$  patience

Cake tomorrow is  $\gamma$  as good as now

# Discounted reward MDP



Objective:

Discounted return  $G$

$$G_t = r_t + \gamma \cdot r_{t+1} + \gamma^2 \cdot r_{t+2} + \dots + \gamma^n \cdot r_{t+n}$$

$$G_t = \sum_i \gamma^i \cdot r_{t+i} \quad \gamma \in (0,1) \text{ const}$$

$\gamma \sim$  patience

Cake tomorrow is  $\gamma$  as good as now

**Q1:** which  $\gamma$  corresponds to “only current reward matters”?

**Q2:** with which  $\gamma$   $G$  is just a sum of rewards?

# Discounted reward MDP



Objective:  
Discounted return  $G$

$$G_t = r_t + \gamma \cdot r_{t+1} + \gamma^2 \cdot r_{t+2} + \dots + \gamma^n \cdot r_{t+n}$$

$$G_t = \sum_i \gamma^i \cdot r_{t+i} \quad \gamma \in (0,1) \text{ const}$$

Reinforcement learning:

- Find policy that maximizes the expected reward

$$\pi = P(a|s) : E[G] \rightarrow \max$$

# Discounted reward MDP



Objective:

Discounted return  $G$

$$G_t = r_t + \gamma \cdot r_{t+1} + \gamma^2 \cdot r_{t+2} + \dots + \gamma^n \cdot r_{t+n}$$

$$G_t = \sum_i \gamma^i \cdot r_{t+i} \quad \gamma \in (0,1) \text{ const}$$

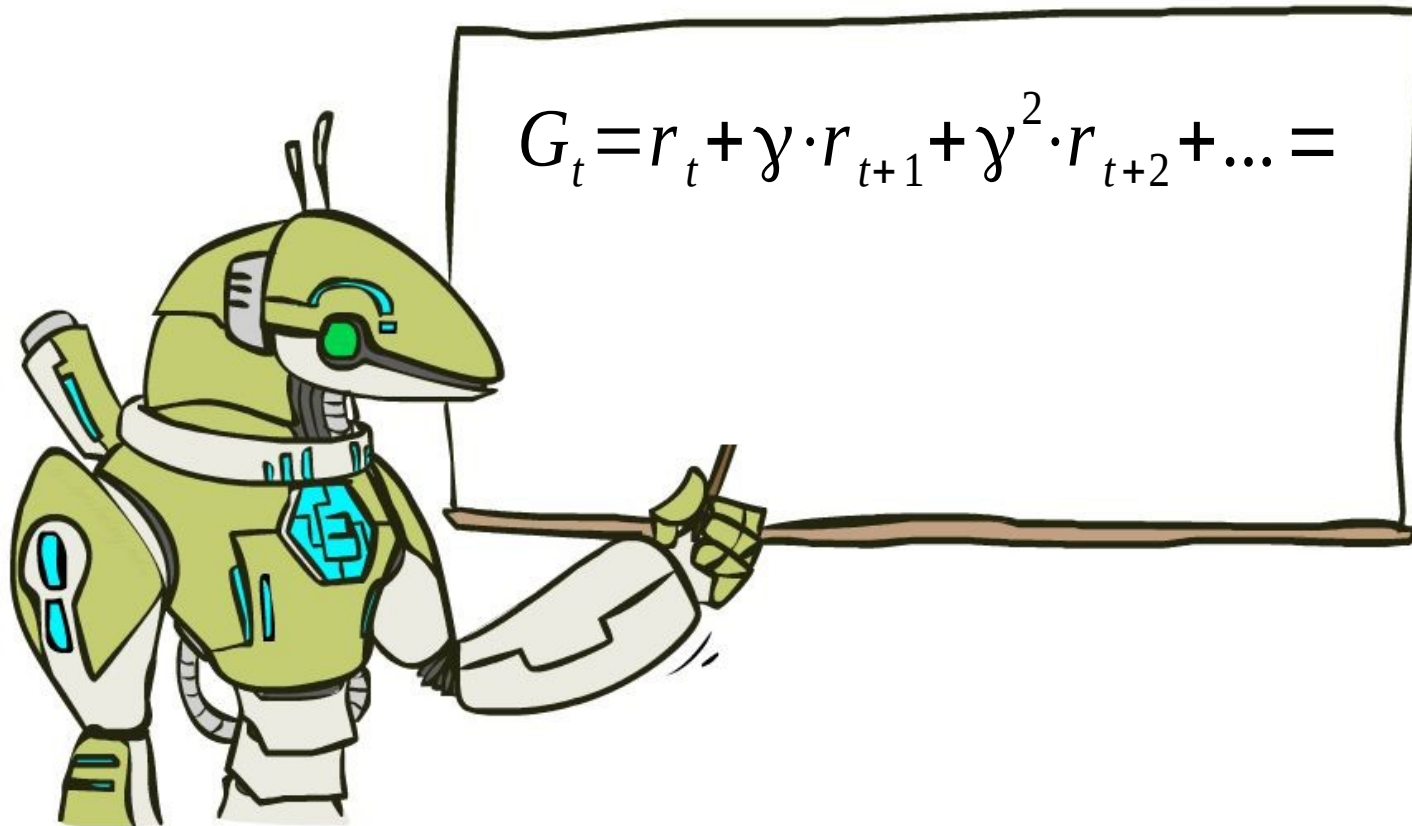
Reinforcement learning:

- Find policy that maximizes the expected reward

$$\pi = P(a|s) : E[G] \rightarrow \max$$

Does not maximize sum of rewards  
Unless  $\gamma = 1$

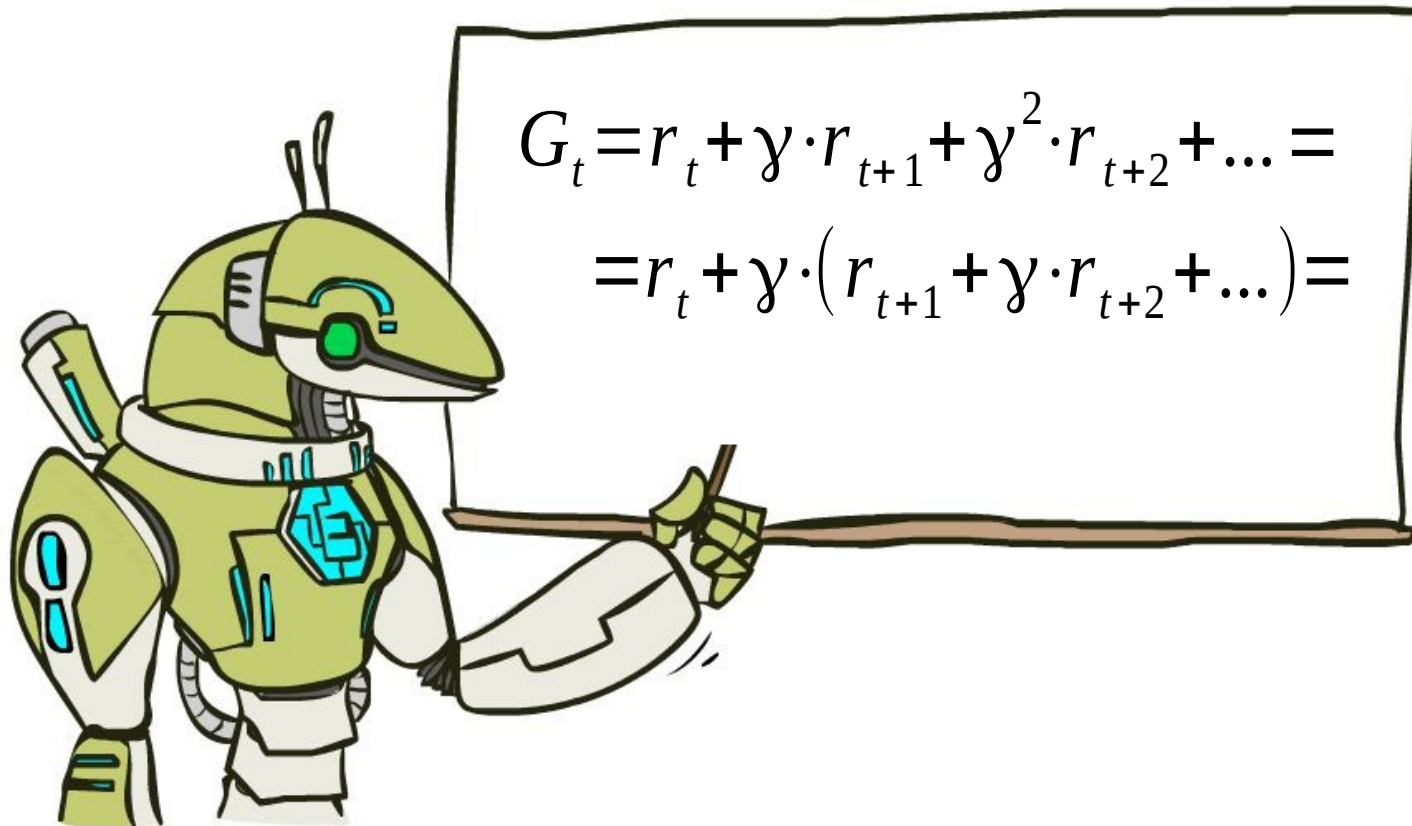
# Dealing with G



We rewrite G with sheer power of math!

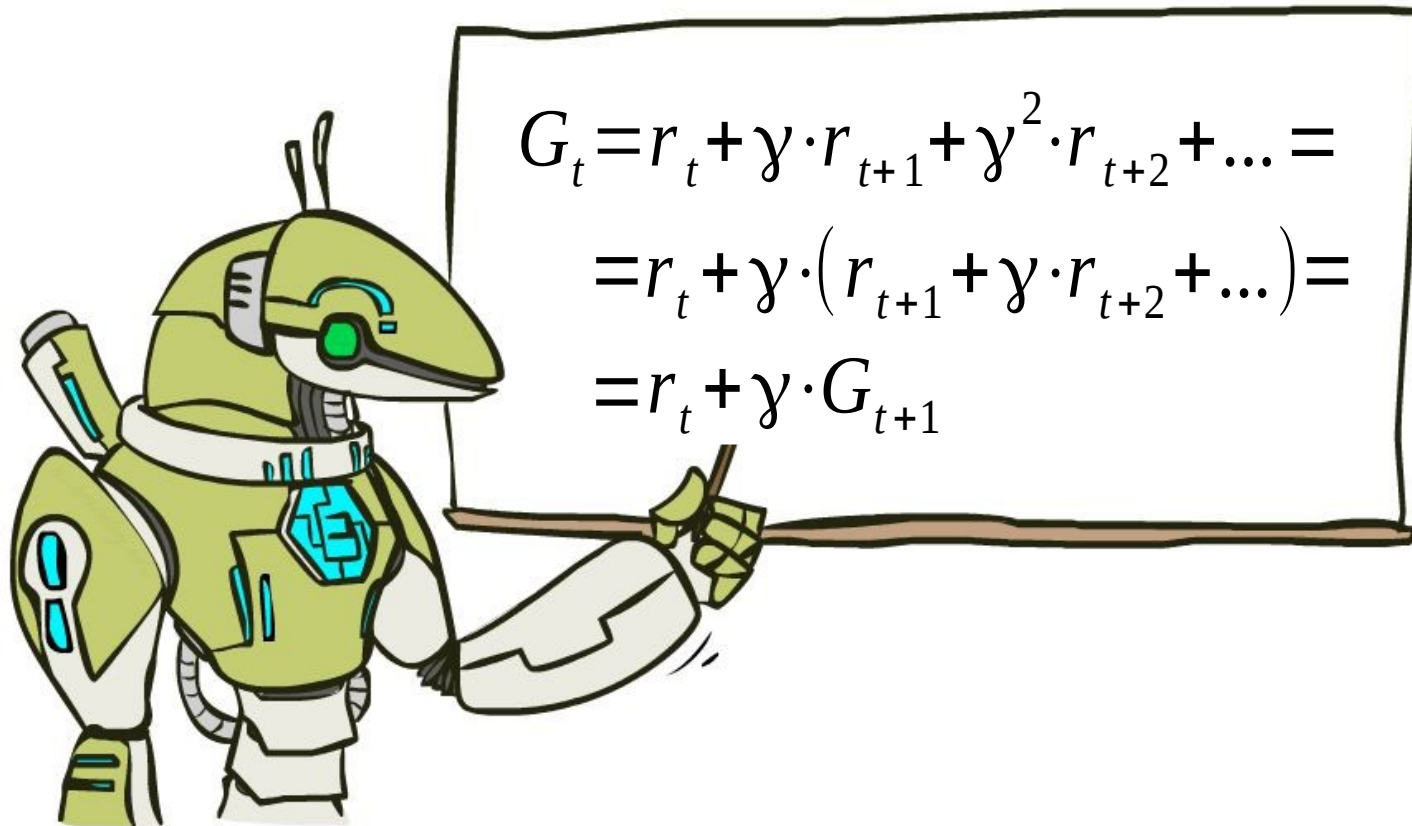


# Dealing with G



We rewrite G with sheer power of math!

# Dealing with G



We rewrite G with sheer power of math!

# More new letters!

State values :  $V_{\pi}(s)$

- **Definition:**  $V_{\pi}(s)$  – expected total reward  $G$  that can be obtained starting from state  $s$  and following policy  $\pi$ .

large potential rewards



vs

small potential rewards



# More new letters!

State values :  $V_{\pi}(s)$

- **Definition:**  $V_{\pi}(s)$  – expected total reward  $G$  that can be obtained starting from state  $s$  and following policy  $\pi$ .

$$V_{\pi}(s) = E_{a \sim \pi(a|s)} \dots$$

# State value

State values :  $V_{\pi}(s)$

- **Definition:**  $V_{\pi}(s)$  – expected total reward  $G$  that can be obtained starting from state  $s$  and following policy  $\pi$ .

$$V_{\pi}(s) = \mathop{E}_{a \sim \pi(a|s)} \mathop{E}_{s', r \sim P(s'|s, a)} \dots$$

# State value

State values :  $V_{\pi}(s)$

- **Definition:**  $V_{\pi}(s)$  – expected total reward  $G$  that can be obtained starting from state  $s$  and following policy  $\pi$ .

$$V_{\pi}(s) = \underset{a \sim \pi(a|s)}{E} \underset{s', r \sim P(s'|s, a)}{E} \underset{a', r', s'', \dots}{E} \dots$$

# State value

State values :  $V_{\pi}(s)$

- **Definition:**  $V_{\pi}(s)$  – expected total reward  $G$  that can be obtained starting from state  $s$  and following policy  $\pi$ .

$$V_{\pi}(s) = E_{a \sim \pi(a|s)} E_{s', r \sim P(s'|s, a)} E_{a', r', s'', \dots} r_t + \gamma \cdot r_{t+1} + \gamma^2 \cdot r_{t+2} + \dots + \gamma^n \cdot r_{t+n}$$

# State value

State values :  $V_{\pi}(s)$

- **Definition:**  $V_{\pi}(s)$  – expected total reward  $G$  that can be obtained starting from state  $s$  and following policy  $\pi$ .

$$V_{\pi}(s) = E_{\tau \sim p(\tau|s)} G(\tau)$$

trajectory  $\tau = (s, a, r, s', a', r', \dots)$



# Bellman equation

State values :  $V_{\pi}(s)$

**Definition:**  $V_{\pi}(s)$  – expected total reward  $R$  that can be obtained starting from state  $s$  and following policy  $\pi$ .

$$V_{\pi}(s) = E_{a \sim \pi(a|s)} E_{s', r \sim P(s', r|s, a)} [r + \gamma \cdot V_{\pi}(s')]$$

# Bellman equation

State values :  $V_{\pi}(s)$

**Definition:**  $V_{\pi}(s)$  – expected total reward  $R$  that can be obtained starting from state  $s$  and following policy  $\pi$ .

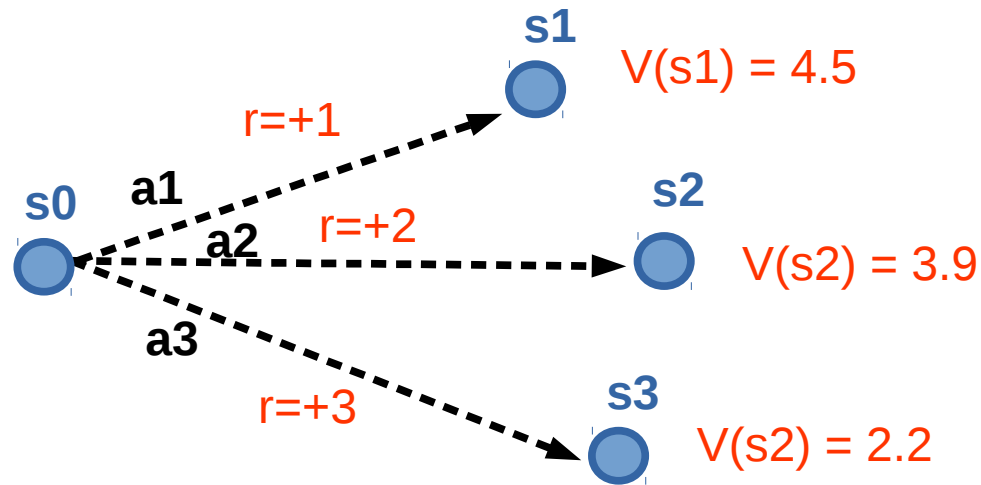
$$V_{\pi}(s) = E_{a \sim \pi(a|s)} E_{s', r \sim P(s', r|s, a)} [r + \gamma \cdot V_{\pi}(s')]$$

how I like  $s$  = what  $r$  I can get right now +  $\gamma$  \* how I like  $s'$

# Optimal policy

Imagine we know exact value of any  $V_{\pi}(s)$ .  
(also  $P(s',r|s,a)$ )

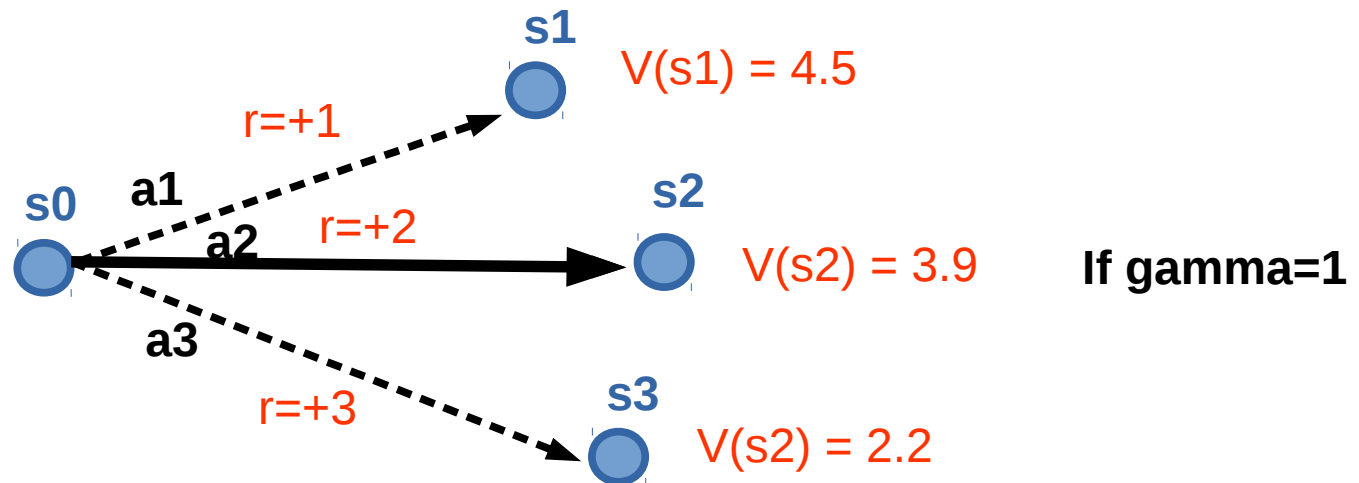
Q: how do we pick optimal action?



# Optimal policy

Imagine we know exact value of any  $V_{\pi}(s)$ .  
(also  $P(s',r|s,a)$ )

Q: how do we pick optimal action?



$$\pi^*(a|s) = \underset{a}{\operatorname{argmax}} E_{s', r \sim P(s', r|s, a)} [r + \gamma \cdot V_{\pi^*}(s')]$$

# $\pi^*(a|s)$ and $V^*(s)$

Optimal policy :  $\pi^*(a|s)$

$$\pi^*(a|s) = \underset{a}{\operatorname{argmax}} \ E_{s', r \sim P(s', r|s, a)} [r + \gamma \cdot V_{\pi^*}(s')]$$

Optimal state value :  $V^*(s)$

$$V^*(s) = V_{\pi^*}(s) = \underset{a}{\operatorname{max}} \ E_{s', r \sim P(s', r|s, a)} [r + \gamma \cdot V(s')]$$

# Value iteration (TD)

Idea:

- Iterative updates

$$\forall s, V_0(s) := 0$$

$$V_{i+1}(s) := \max_a [r(s, a) + \gamma \cdot E_{s' \sim P(s'|s, a)} V_i(s')]$$

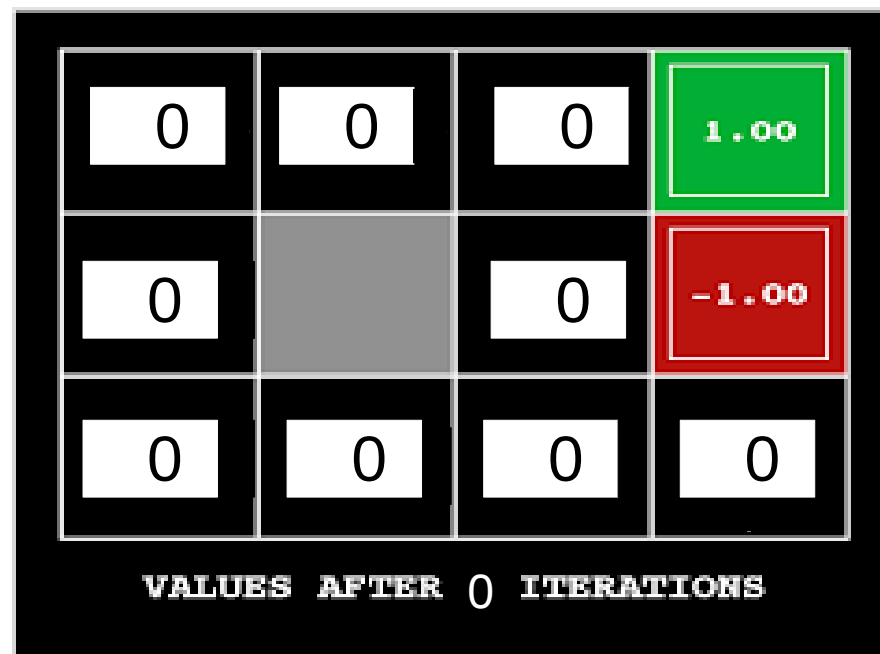
# Value iteration (TD)

Idea:

- Iterative updates

$$\forall s, V_0(s) := 0$$

$$V_{i+1}(s) := \max_a [r(s, a) + \gamma \cdot E_{s' \sim P(s'|s, a)} V_i(s')]$$



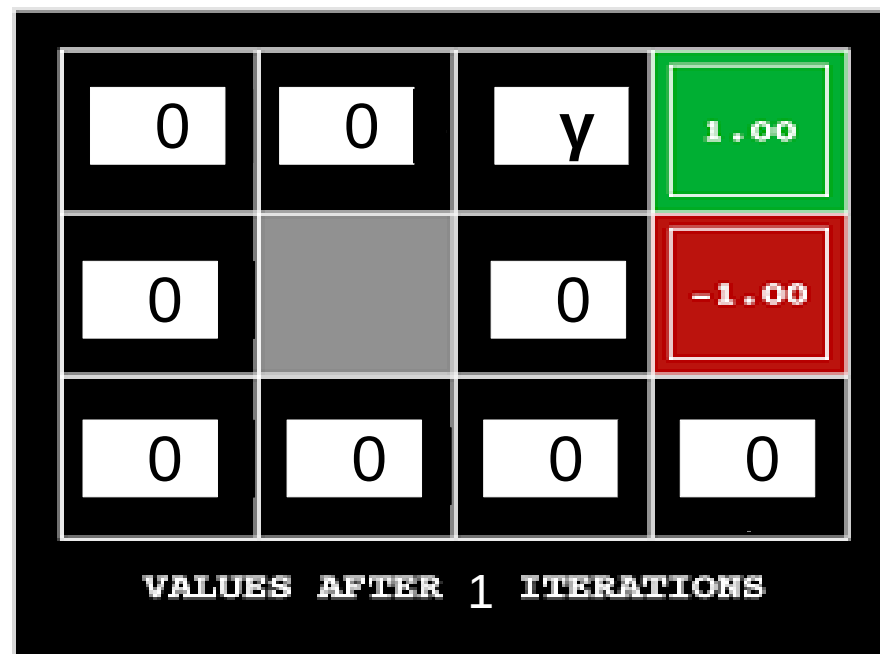
# Value iteration (TD)

Idea:

- Iterative updates

$$\forall s, V_0(s) := 0$$

$$V_{i+1}(s) := \max_a [r(s, a) + \gamma \cdot E_{s' \sim P(s'|s, a)} V_i(s')]$$





# Value iteration (TD)

Idea:

- Iterative updates

$$\forall s, V_0(s) := 0$$

$$V_{i+1}(s) := \max_a [r(s, a) + \gamma \cdot E_{s' \sim P(s'|s, a)} V_i(s')]$$

0	$\gamma^2$	$\gamma$	1.00
0		$\gamma^2$	-1.00
0	0	0	0

VALUES AFTER 2 ITERATIONS

# Value iteration (TD)

Idea:

- Iterative updates

$$\forall s, V_0(s) := 0$$

$$V_{i+1}(s) := \max_a [r(s, a) + \gamma \cdot E_{s' \sim P(s'|s, a)} V_i(s')]$$

$\gamma^3$	$\gamma^2$	$\gamma$	1.00
0		$\gamma^2$	-1.00
0	0	$\gamma^3$	0

VALUES AFTER 3 ITERATIONS

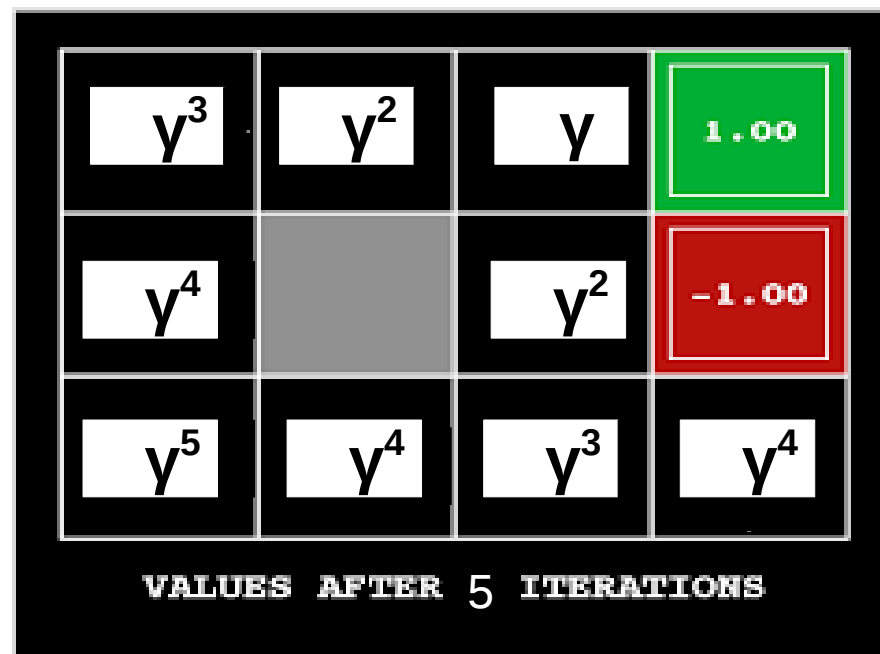
# Value iteration (TD)

Idea:

- Iterative updates

$$\forall s, V_0(s) := 0$$

$$V_{i+1}(s) := \max_a [r(s, a) + \gamma \cdot E_{s' \sim P(s'|s, a)} V_i(s')]$$



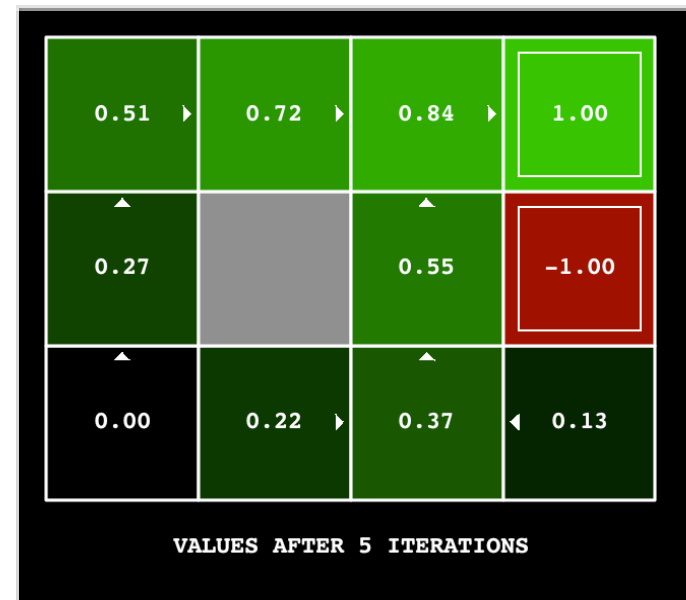
# Value iteration (TD)

Idea:

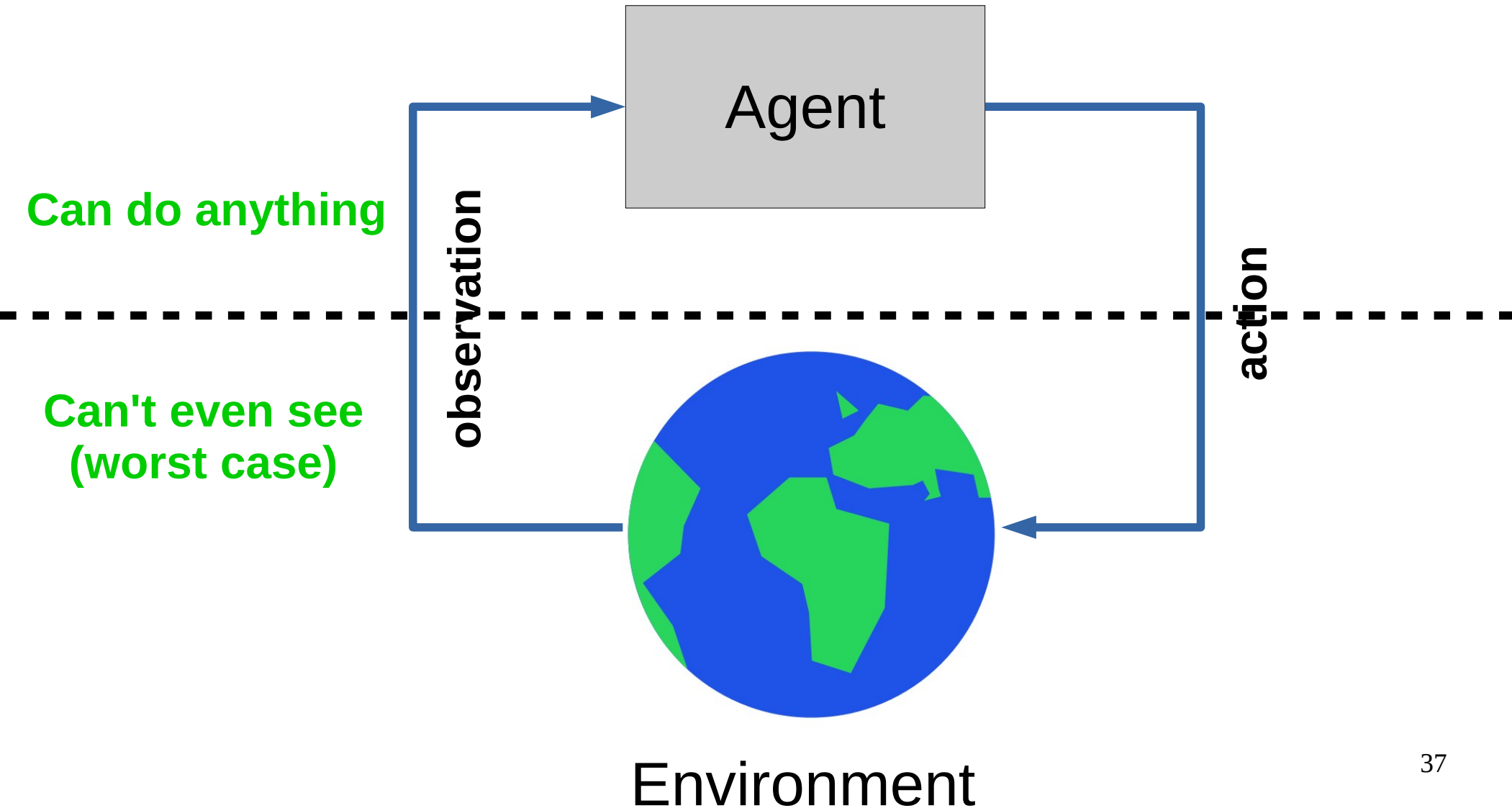
- Iterative updates

$$\forall s, V_0(s) := 0$$

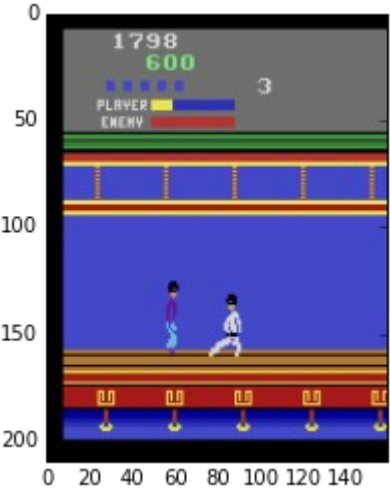
$$V_{i+1}(s) := \max_a [r(s, a) + \gamma \cdot E_{s' \sim P(s'|s, a)} V_i(s')]$$



# Decision process in the wild



# Decision process in the wild



Model-free setting:

We don't know actual

$$P(s',r|s,a)$$

Whachagonnado?

Model-free setting:

We don't know actual

$$P(s',r|s,a)$$

Learn it?

Get rid of it?



# More new letters

- $V_{\pi}(\mathbf{s})$  – expected G from state  $\mathbf{s}$  if you follow  $\pi$
- $V^*(\mathbf{s})$  – expected G from state  $\mathbf{s}$  if you follow  $\pi^*$

# More new letters

- $V_{\pi}(\mathbf{s})$  – expected G from state  $\mathbf{s}$  if you follow  $\pi$
- $V^*(\mathbf{s})$  – expected G from state  $\mathbf{s}$  if you follow  $\pi^*$
- $Q_{\pi}(\mathbf{s}, \mathbf{a})$  – expected G from state  $\mathbf{s}$ 
  - if you start by taking action  $\mathbf{a}$
  - and follow  $\pi$  from next state on
- $Q^*(\mathbf{s}, \mathbf{a})$  – guess what it is :)

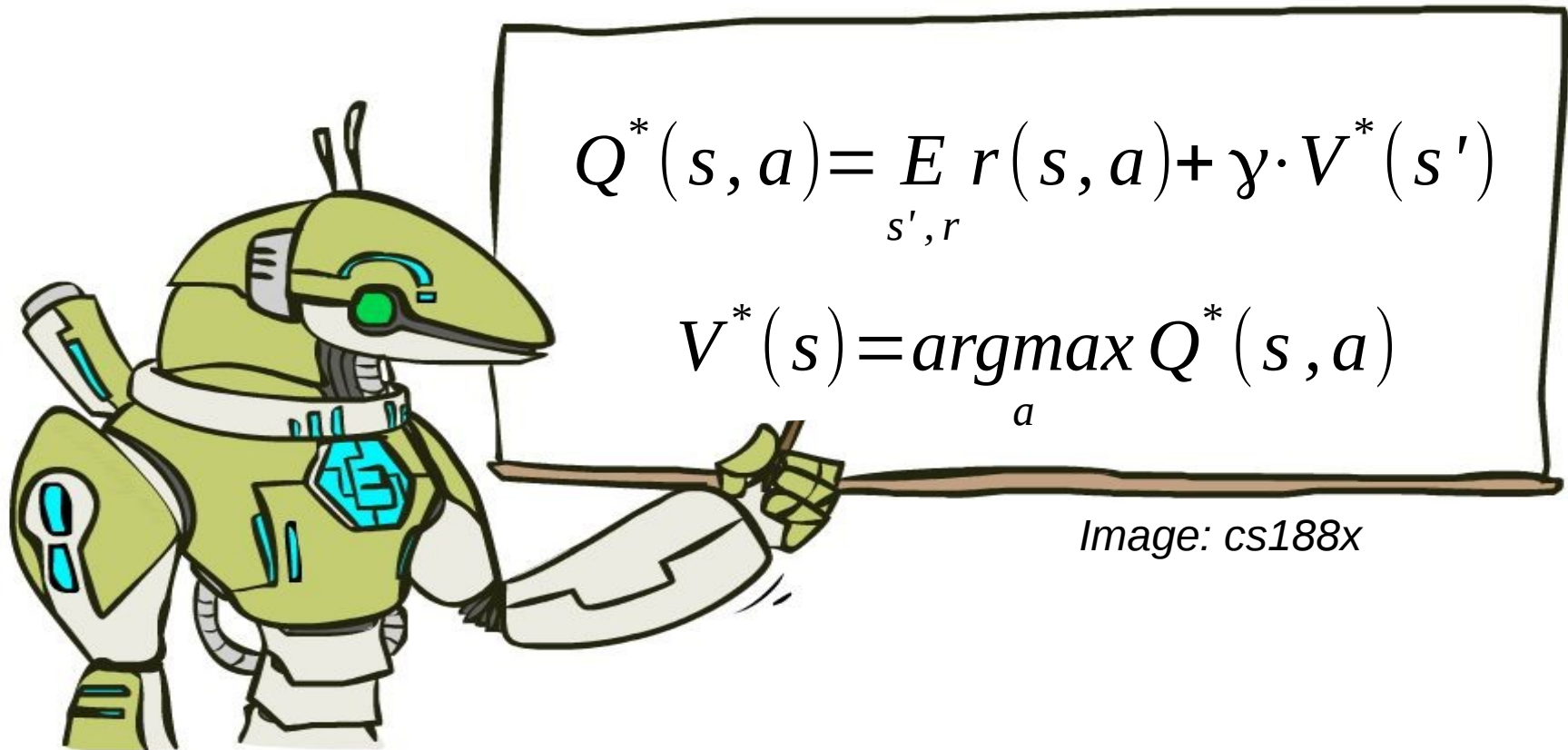
# More new letters

- $V_{\pi}(\mathbf{s})$  – expected G from state  $\mathbf{s}$  if you follow  $\pi$
- $V^*(\mathbf{s})$  – expected G from state  $\mathbf{s}$  if you follow  $\pi^*$
- $Q_{\pi}(\mathbf{s}, \mathbf{a})$  – expected G from state  $\mathbf{s}$ 
  - if you start by taking action  $\mathbf{a}$
  - and follow  $\pi$  from next state on
- $Q^*(\mathbf{s}, \mathbf{a})$  – same as  $Q_{\pi}(\mathbf{s}, \mathbf{a})$  where  $\pi = \pi^*$

# Any ideas?

- Assuming you know  $Q^*(s,a)$ ,
  - how do you compute  $\pi^*$
  - how do you compute  $V^*(s)$ ?
- Assuming you know  $V(s)$ 
  - how do you compute  $Q(s,a)$ ?

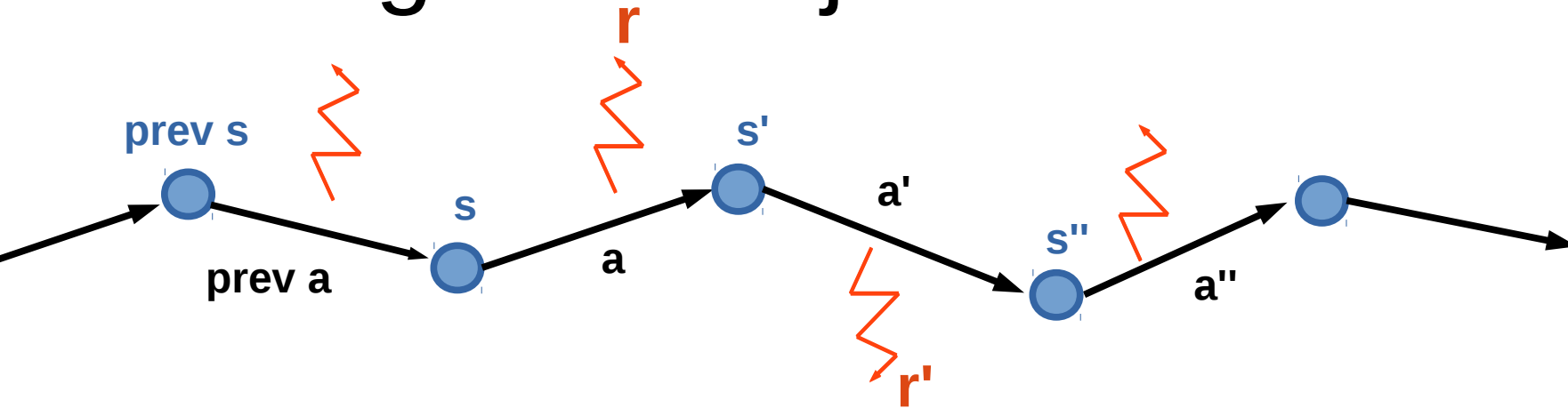
# To sum up



**Action value  $Q_{\pi}(s, a)$**  is the expected total reward **G** agent gets from state **s** by taking action **a** and following policy  **$\pi$**  from next state.

$$\pi(s) : \operatorname{argmax}_a Q(s, a)$$

# Learning from trajectories



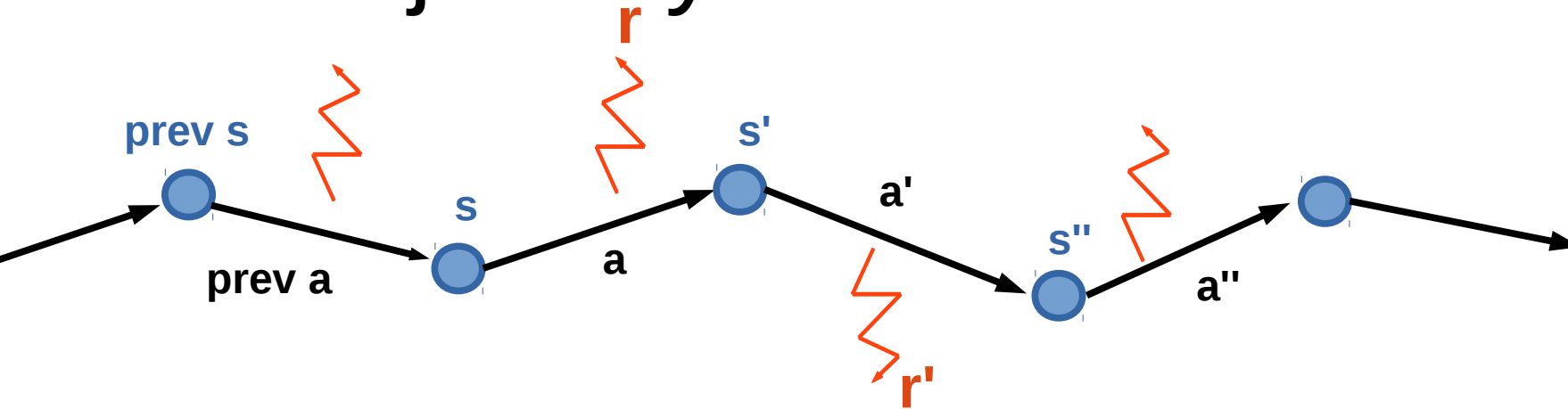
**Model-based:** you know  $P(s'|s,a)$

- can apply dynamic programming
- can plan ahead

**Model-free:** you can sample trajectories

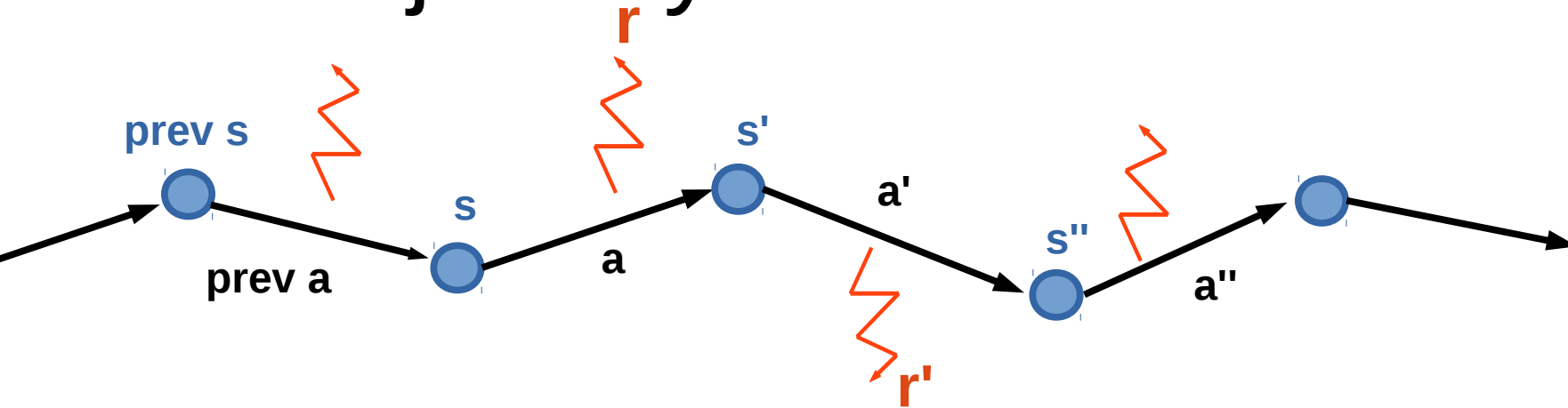
- can try stuff out
- insurance not included

# MDP trajectory



- Trajectory is a sequence of
  - states ( $s$ )
  - actions ( $a$ )
  - rewards ( $r$ )
- We can only sample trajectories

# MDP trajectory



- Trajectory is a sequence of

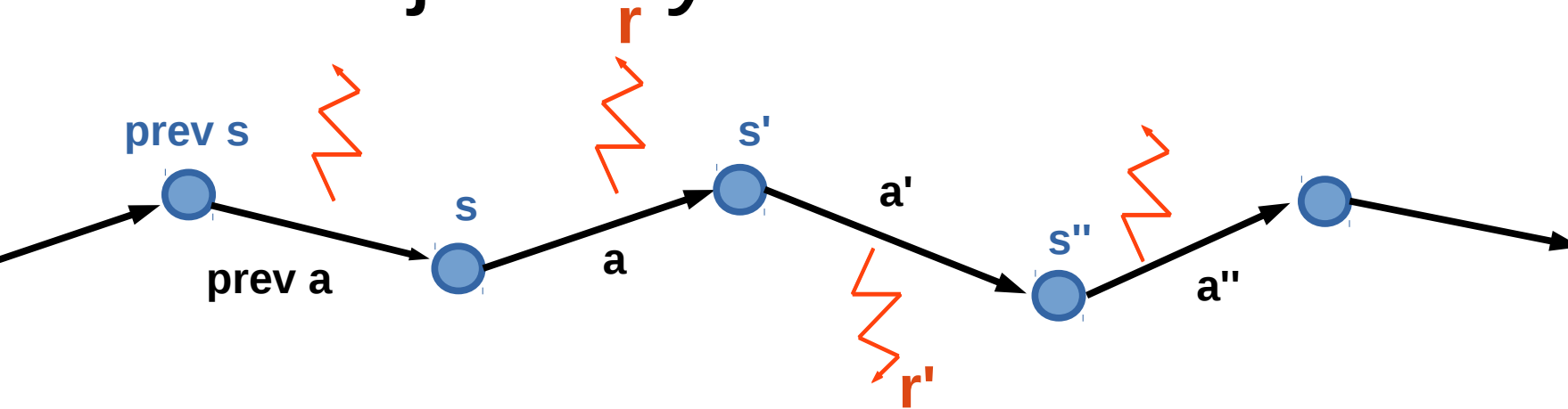
- states ( $s$ )
- actions ( $a$ )
- rewards ( $r$ )

**Q:** What to learn?  
 $V(s)$  or  $Q(s,a)$

- We can only sample trajectories



# MDP trajectory



- Trajectory is a sequence of

- states ( $s$ )
- actions ( $a$ )
- rewards ( $r$ )

**Q:** What to learn?  
 $V(s)$  or  $Q(s,a)$

$V(s)$  is useless  
without  $P(s'|s,a)$

- We can only sample trajectories

# More temporal difference

- Remember we can improve  $Q(s,a)$  iteratively!

$$Q(s_t, a_t) \leftarrow E_{r_t, s_{t+1}} r_t + \gamma \cdot \max_{a'} Q(s_{t+1}, a')$$

# More temporal difference

- Remember we can improve  $Q(s,a)$  iteratively!

$$Q(s_t, a_t) \leftarrow E_{r_t, s_{t+1}} r_t + \gamma \cdot \max_{a'} Q(s_{t+1}, a')$$

↑  
That's  $Q^*(s,a)$

↑  
That's value for  $\pi^*$   
aka optimal policy

# More temporal difference

- Remember we can improve  $Q(s,a)$  iteratively!

$$Q(s_t, a_t) \leftarrow E_{r_t, s_{t+1}} r_t + \gamma \cdot \max_{a'} Q(s_{t+1}, a')$$

↑  
That's  $Q^*(s,a)$

↑  
That's value for  $\pi^*$   
aka optimal policy

↑  
That's something  
we don't have

What do we do?

# More temporal difference



# More temporal difference

- Replace expectation with sampling

$$E_{r_t, s_{t+1}} r_t + \gamma \cdot \max_{a'} Q(s_{t+1}, a') \approx \frac{1}{N} \sum_i r_i + \gamma \cdot \max_{a'} Q(s_i^{\text{next}}, a')$$

# More temporal difference

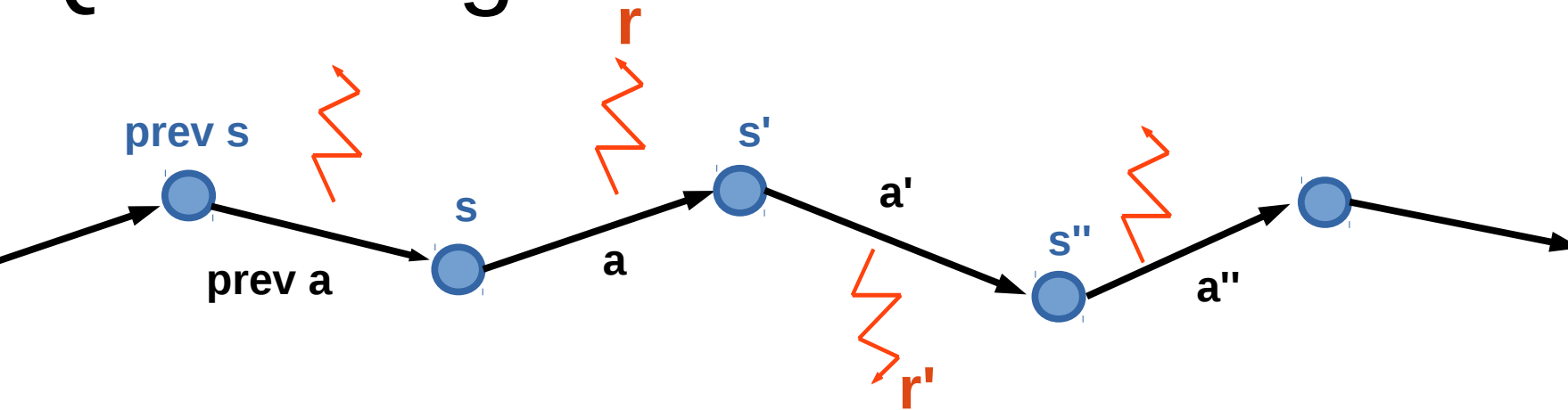
- Replace expectation with sampling

$$E_{r_t, s_{t+1}} r_t + \gamma \cdot \max_{a'} Q(s_{t+1}, a') \approx \frac{1}{N} \sum_i r_i + \gamma \cdot \max_{a'} Q(s_i^{\text{next}}, a')$$

- Use moving average with just one sample!

$$Q(s_t, a_t) \leftarrow \alpha \cdot (r_t + \gamma \cdot \max_{a'} Q(s_{t+1}, a')) + (1 - \alpha) Q(s_t, a_t)$$

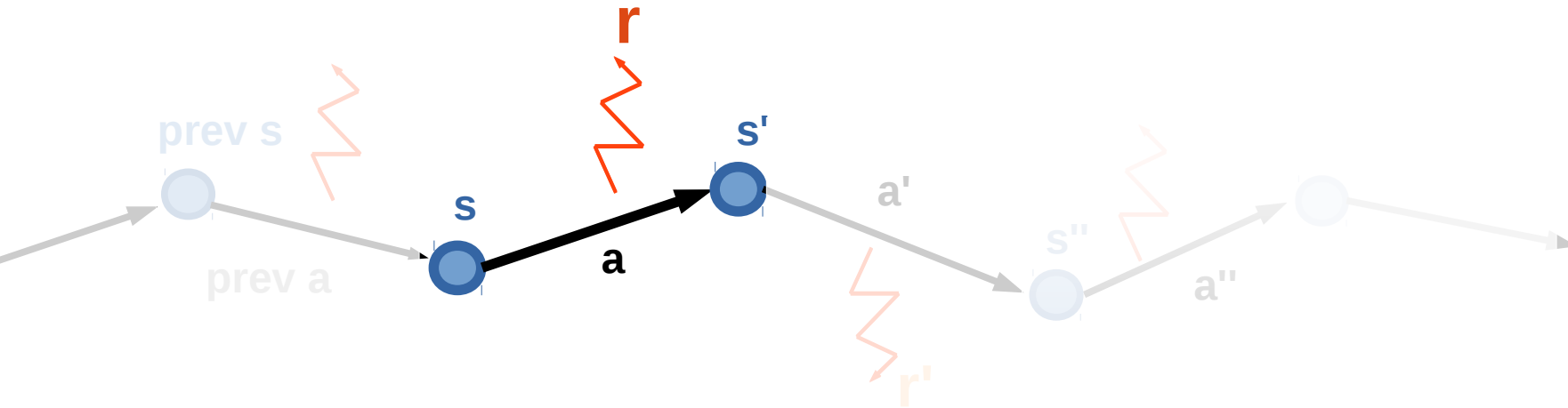
# Q-learning



- Works on a sequence of
  - states ( $s$ )
  - actions ( $a$ )
  - rewards ( $r$ )



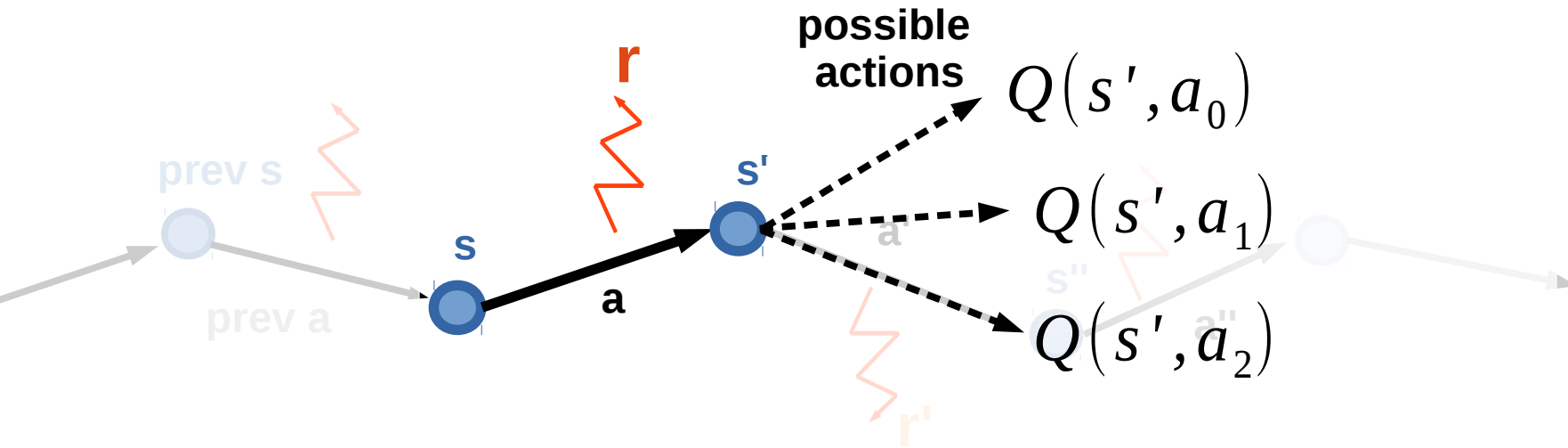
# Q-learning



Initialize  $Q(s,a)$  with zeros

- Loop:
  - Sample  $\langle s, a, r, s' \rangle$  from env

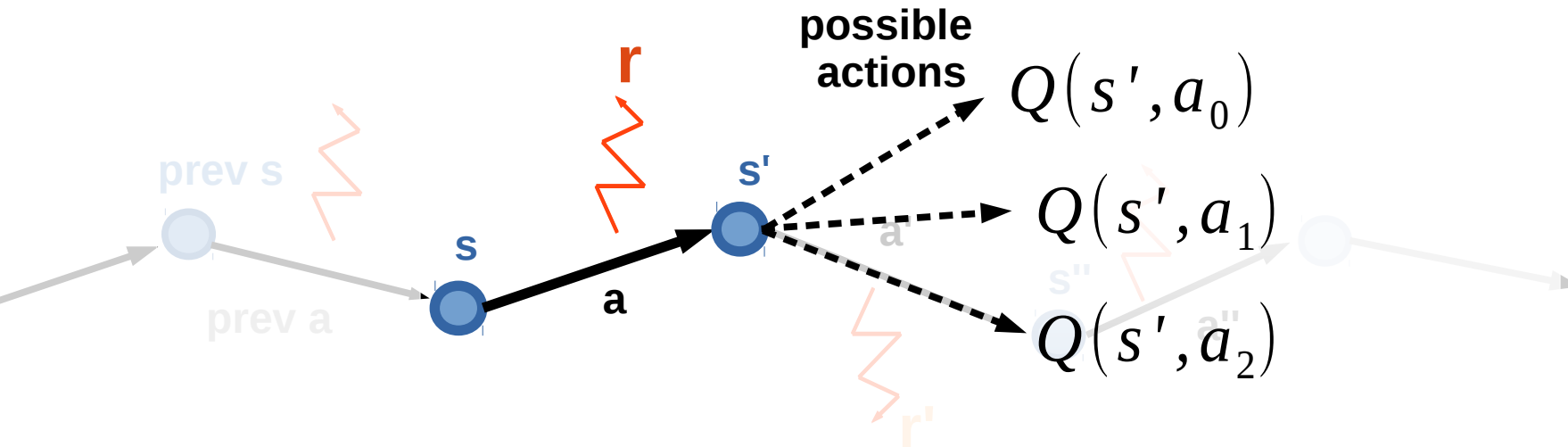
# Q-learning



Initialize  $Q(s,a)$  with zeros

- Loop:
  - Sample  $\langle \mathbf{s}, \mathbf{a}, \mathbf{r}, \mathbf{s}' \rangle$  from env
  - Compute 
$$\hat{Q}(s,a) = r(s,a) + \gamma \max_{a_i} Q(s',a_i)$$

# Q-learning

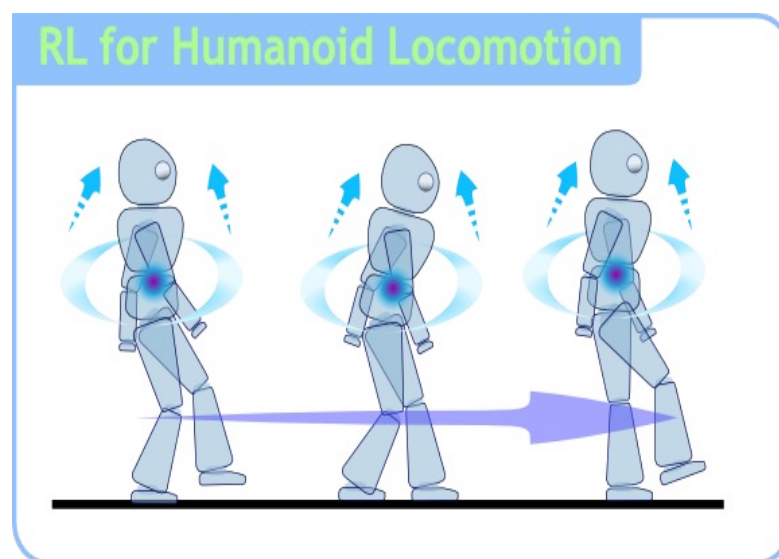


Initialize  $Q(s,a)$  with zeros

- Loop:
  - Sample  $\langle \mathbf{s}, \mathbf{a}, \mathbf{r}, \mathbf{s}' \rangle$  from env
  - Compute  $\hat{Q}(s,a) = r(s,a) + \gamma \max_{a_i} Q(s',a_i)$
  - Update  $Q(s,a) \leftarrow \alpha \cdot \hat{Q}(s,a) + (1-\alpha) Q(s,a)$

# What could possibly go wrong?

Our mobile robot learns to walk.

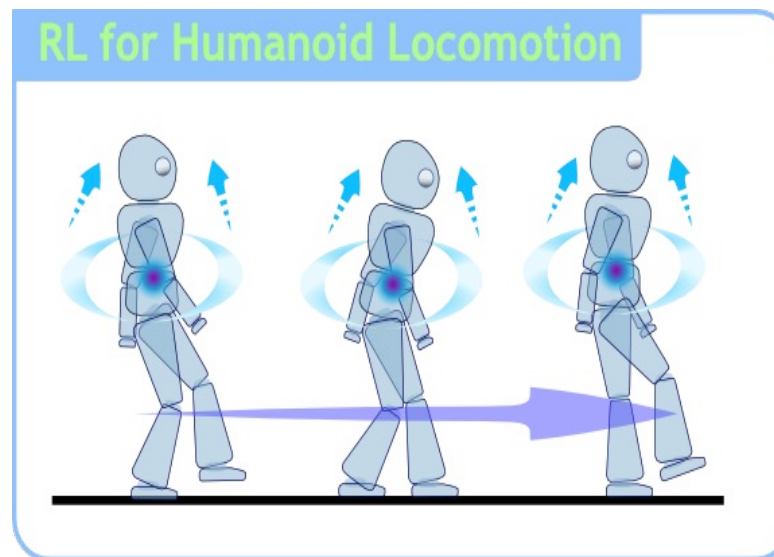


Initial  $Q(s,a)$  are zeros  
robot uses  $\operatorname{argmax} Q(s,a)$

He has just learned to crawl with positive reward! <sup>60</sup>

# What could possibly go wrong?

Our mobile robot learns to walk.



Initial  $Q(s,a)$  are zeros  
robot uses  $\operatorname{argmax} Q(s,a)$

*Too bad, now he will never learn to walk upright = (*<sup>81</sup>

# What could possibly go wrong?

New problem:

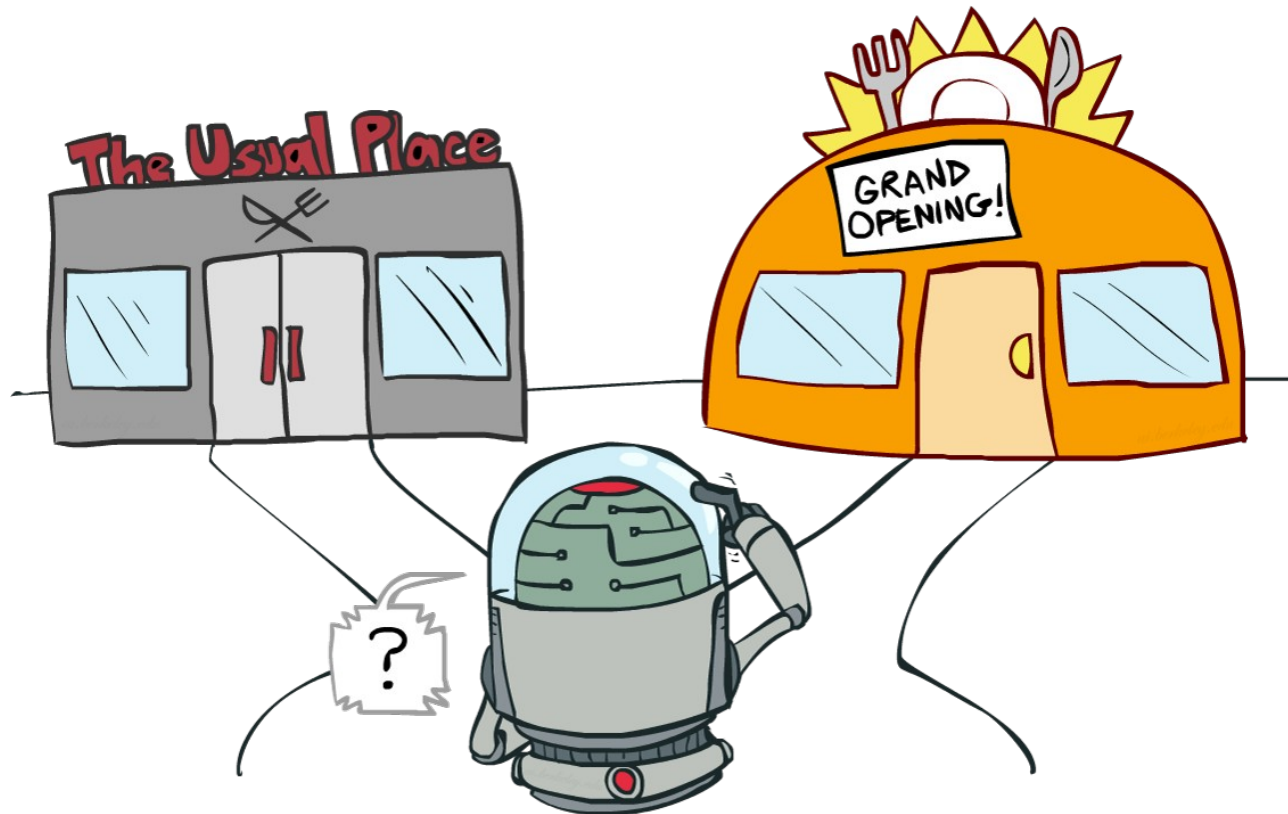
If our agent always takes “best” actions  
from his current point of view,

How will he ever learn that other actions  
may be better than his current best one?

Ideas?

# Exploration Vs Exploitation

Balance between using what you learned and trying to find something even better



# Exploration Vs Exploitation

Strategies:

- $\epsilon$ -greedy
  - With probability  $\epsilon$  take random action; otherwise take optimal action.



# Exploration Vs Exploitation

Strategies:

- $\epsilon$ -greedy
  - With probability  $\epsilon$  take random action; otherwise take optimal action.
- Softmax
  - Pick action proportional to softmax of shifted normalized Q-values.

$$\pi(a|s) = \text{softmax}\left(\frac{Q(s, a)}{\tau}\right)$$

- More cool stuff coming later

# Exploration over time

## Idea:

If you want to converge to optimal policy, you need to gradually reduce exploration

## Example:

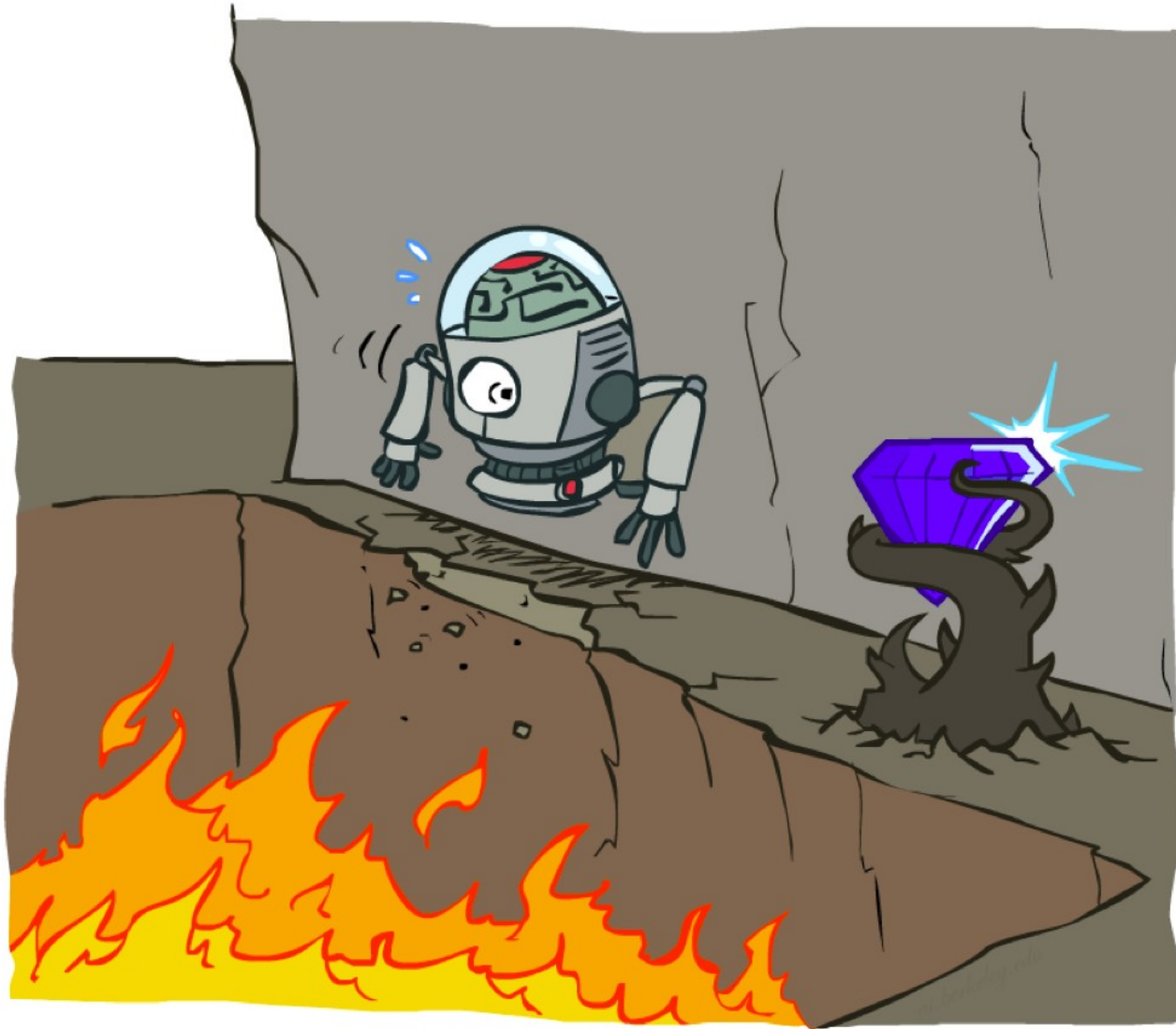
Initialize  $\epsilon$ -greedy  $\epsilon = 0.5$ , then gradually reduce it

- If  $\epsilon \rightarrow 0$ , it's **greedy in the limit**
- Be careful with non-stationary environments

# A popular vote

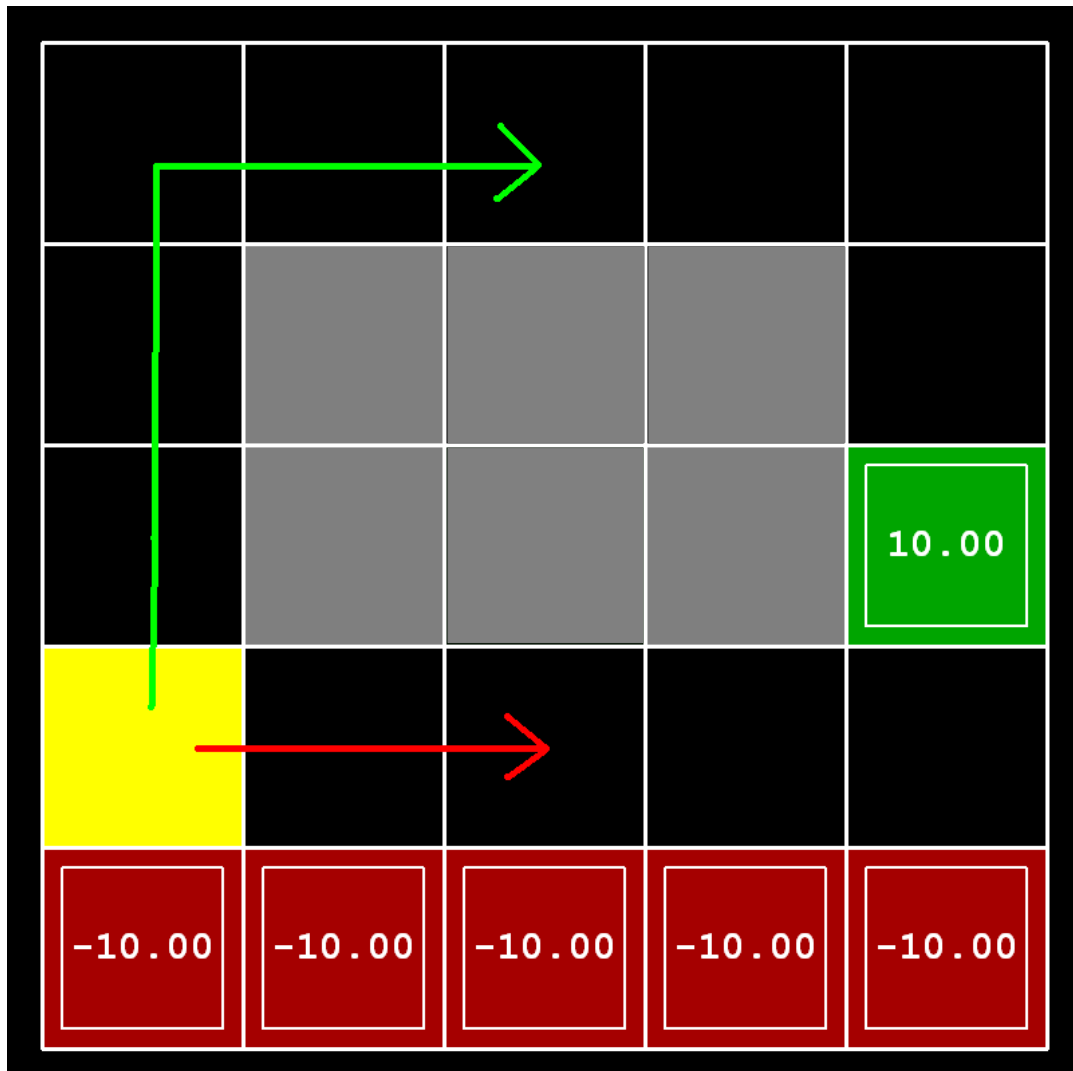
More lecture Vs some seminar?

# Cliff world



Picture from Berkeley CS188x

# Cliff world



## Conditions

- Q-learning

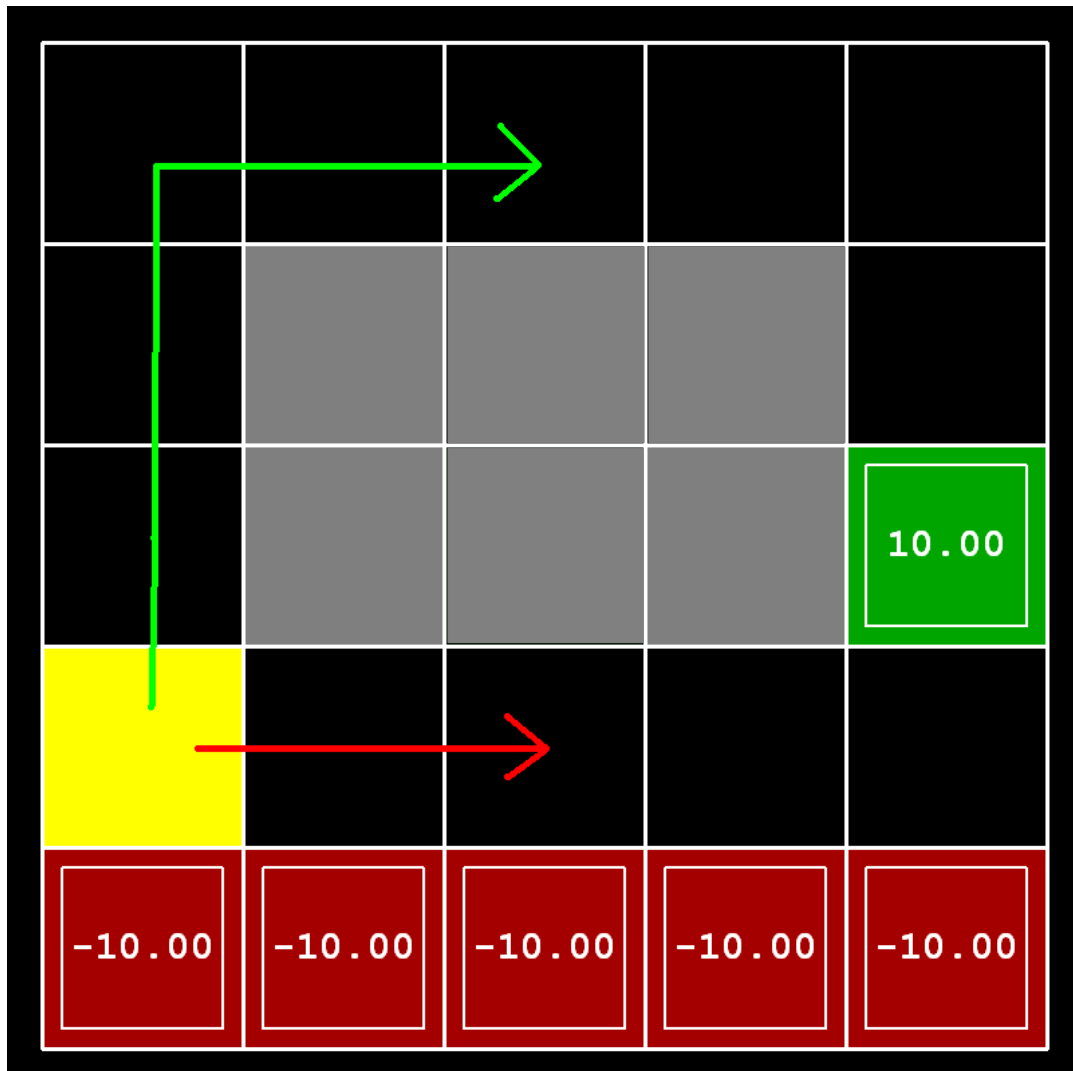
$$\gamma = 0.99 \quad \epsilon = 0.1$$

- no slipping

## Trivia:

What will q-learning learn?

# Cliff world



## Conditions

- Q-learning

$$\gamma = 0.99 \quad \epsilon = 0.1$$

- no slipping

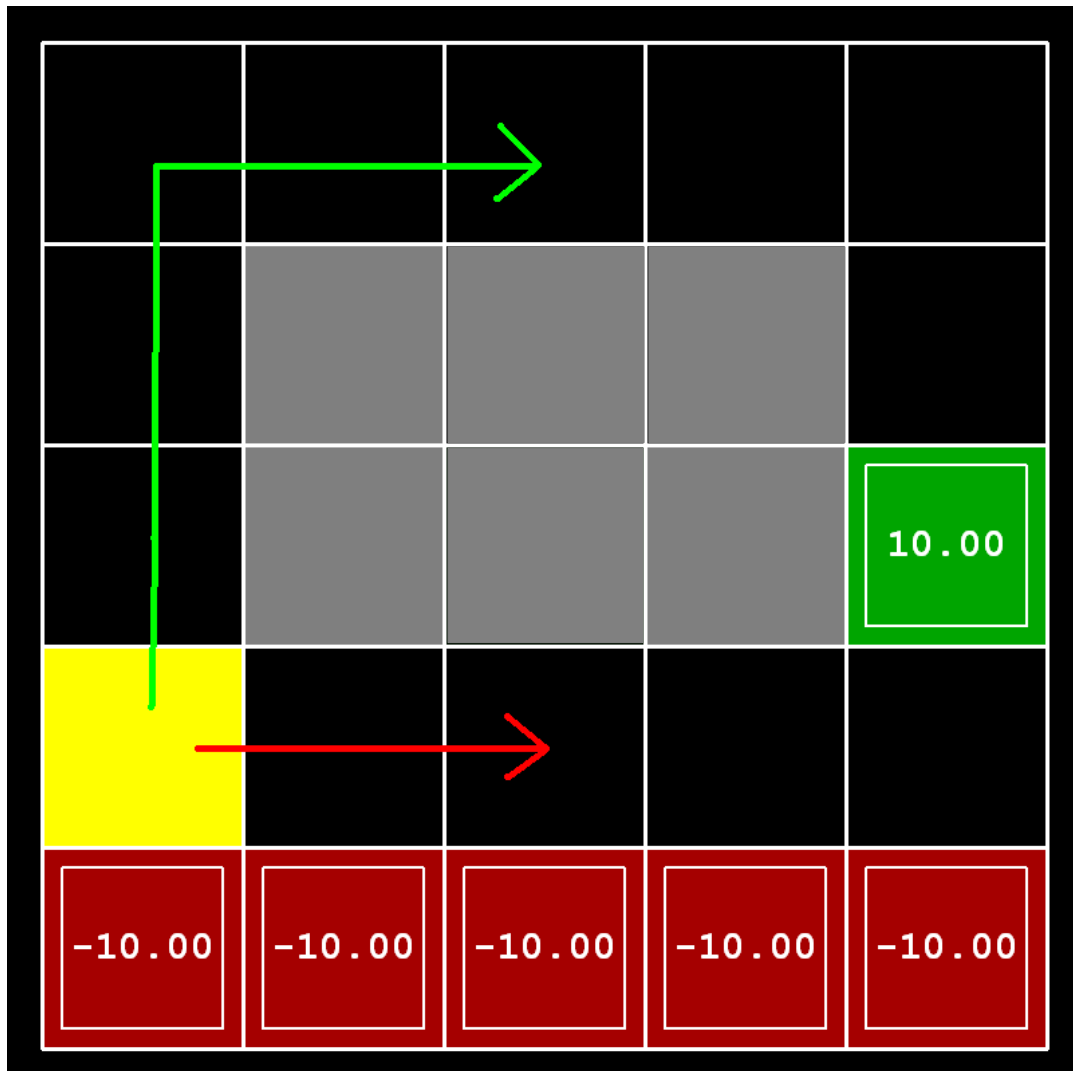
## Trivia:

What will q-learning learn?

**follow the short path**

Will it maximize reward?

# Cliff world



## Conditions

- Q-learning

$$\gamma = 0.99 \quad \epsilon = 0.1$$

- no slipping

## Trivia:

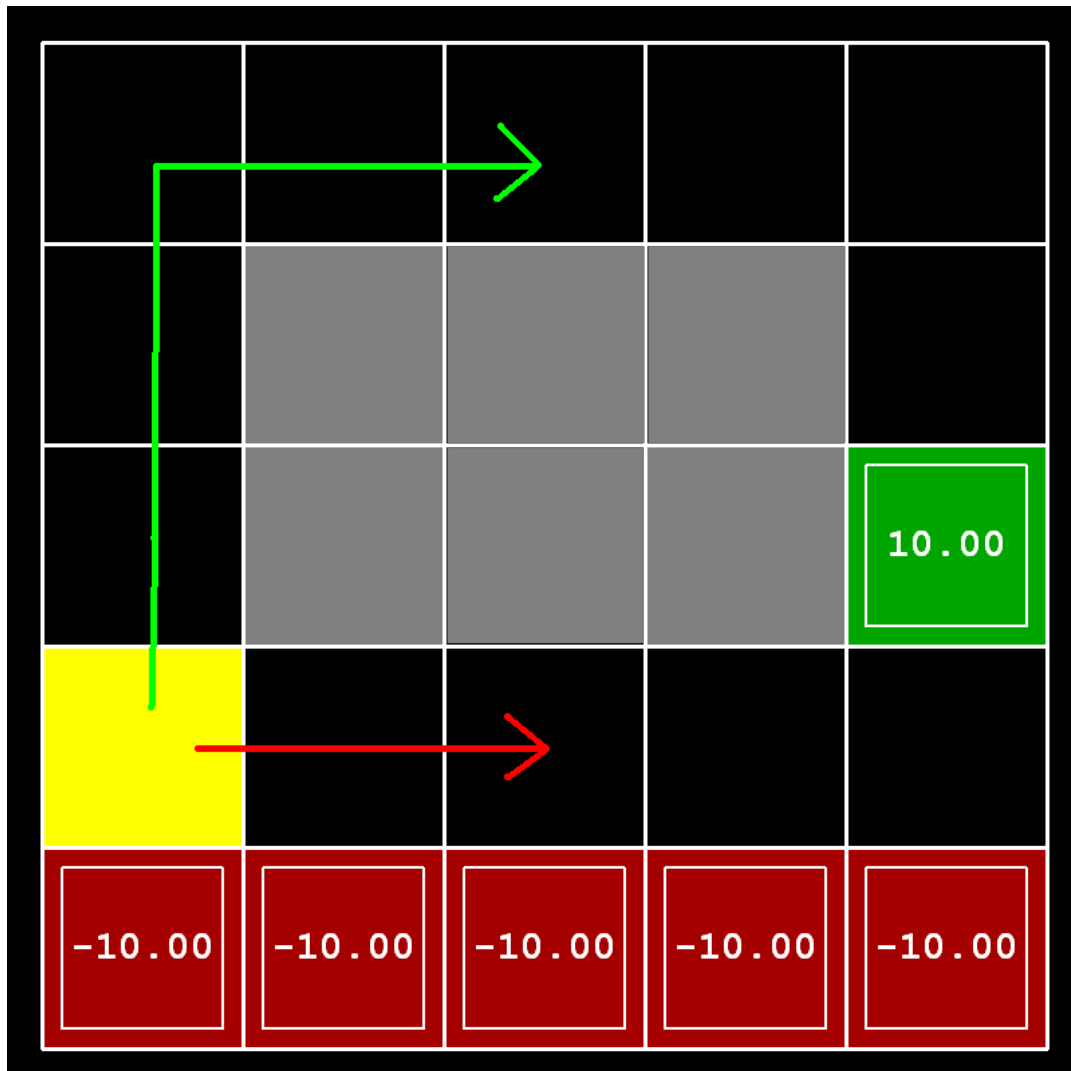
What will q-learning learn?

**follow the short path**

Will it maximize reward?

**no, robot will fall due to epsilon-greedy “exploration”**

# Cliff world



## Conditions

- Q-learning

$$\gamma = 0.99 \quad \epsilon = 0.1$$

- no slipping

**Decisions must account  
for actual policy!**

e.g.  $\epsilon$ -greedy policy



# Generalized update rule

Update rule (from Bellman eq.)

$$Q(s_t, a_t) \leftarrow \alpha \cdot \hat{Q}(s_t, a_t) + (1 - \alpha) Q(s_t, a_t)$$

 “better  $Q(s,a)$ ”

# Q-learning VS SARSA

Update rule (from Bellman eq.)

$$Q(s_t, a_t) \leftarrow \alpha \cdot \hat{Q}(s_t, a_t) + (1 - \alpha) Q(s_t, a_t)$$

Q-learning

“better  $Q(s, a)$ ”



$$\hat{Q}(s, a) = r(s, a) + \gamma \cdot \max_{a'} Q(s', a')$$

# Q-learning VS SARSA

Update rule (from Bellman eq.)

$$Q(s_t, a_t) \leftarrow \alpha \cdot \hat{Q}(s_t, a_t) + (1 - \alpha) Q(s_t, a_t)$$

Q-learning

“better  $Q(s, a)$ ”

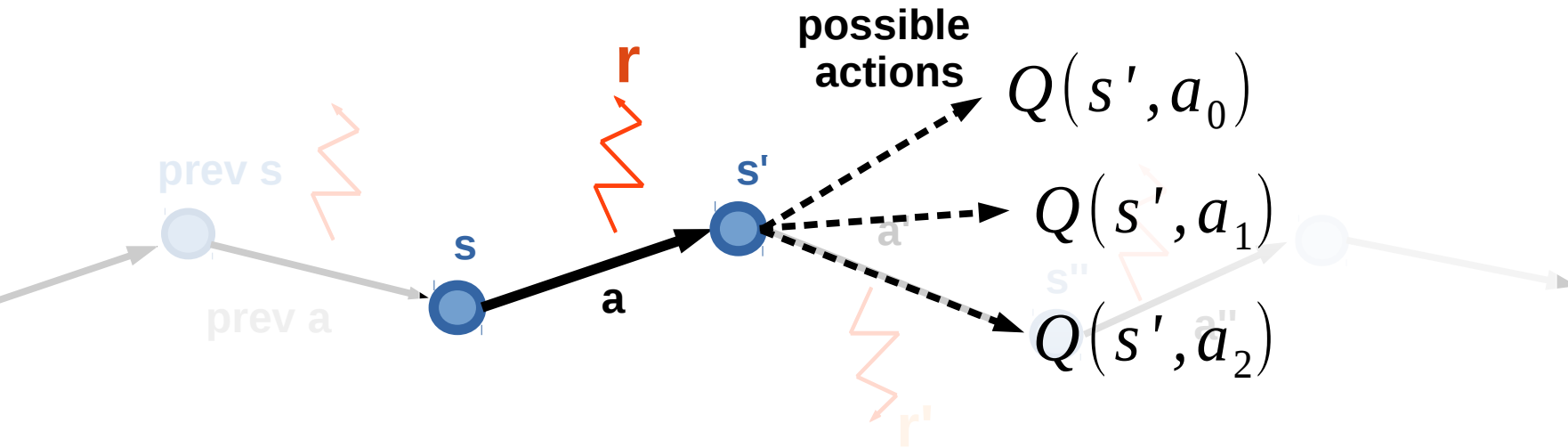


$$\hat{Q}(s, a) = r(s, a) + \gamma \cdot \max_{a'} Q(s', a')$$

SARSA

$$\hat{Q}(s, a) = r(s, a) + \gamma \cdot E_{a' \sim \pi(a'|s')} Q(s', a')$$

# Recap: Q-learning



$$\forall s \in S, \forall a \in A, Q(s, a) \leftarrow 0$$

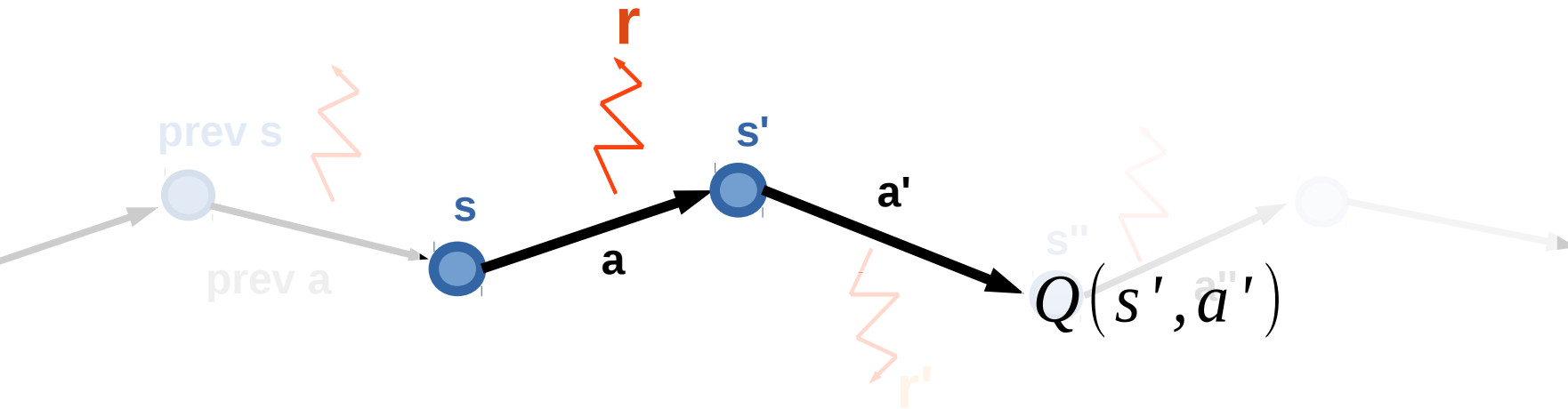
Loop:

- Sample  $\langle \mathbf{s}, \mathbf{a}, \mathbf{r}, \mathbf{s}' \rangle$  from env

- Compute  $\hat{Q}(s, a) = r(s, a) + \gamma \max_{a_i} Q(s', a_i)$

- Update  $Q(s, a) \leftarrow \alpha \cdot \hat{Q}(s, a) + (1 - \alpha) Q(s, a)$

# SARSA



$$\forall s \in S, \forall a \in A, Q(s, a) \leftarrow 0$$

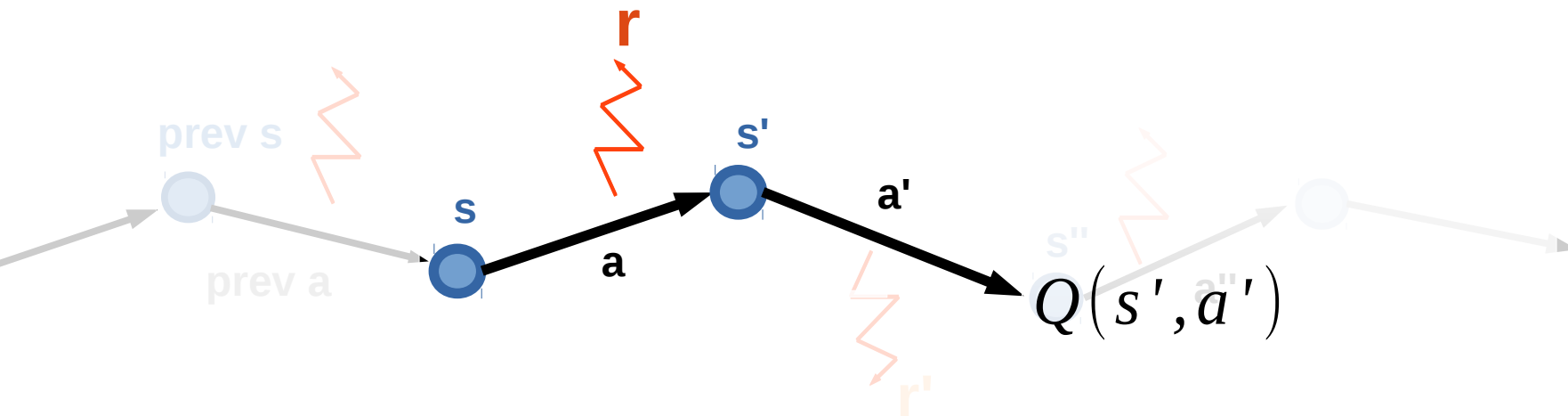
Loop:

- Sample  $\langle \mathbf{s}, \mathbf{a}, \mathbf{r}, \mathbf{s}', \mathbf{a}' \rangle$  from env

- Compute  $\hat{Q}(s, a) = r(s, a) + \gamma Q(s', a')$

- Update  $Q(s, a) \leftarrow \alpha \cdot \hat{Q}(s, a) + (1 - \alpha) Q(s, a)$

# SARSA



$$\forall s \in S, \forall a \in A, Q(s, a) \leftarrow 0$$

Loop:

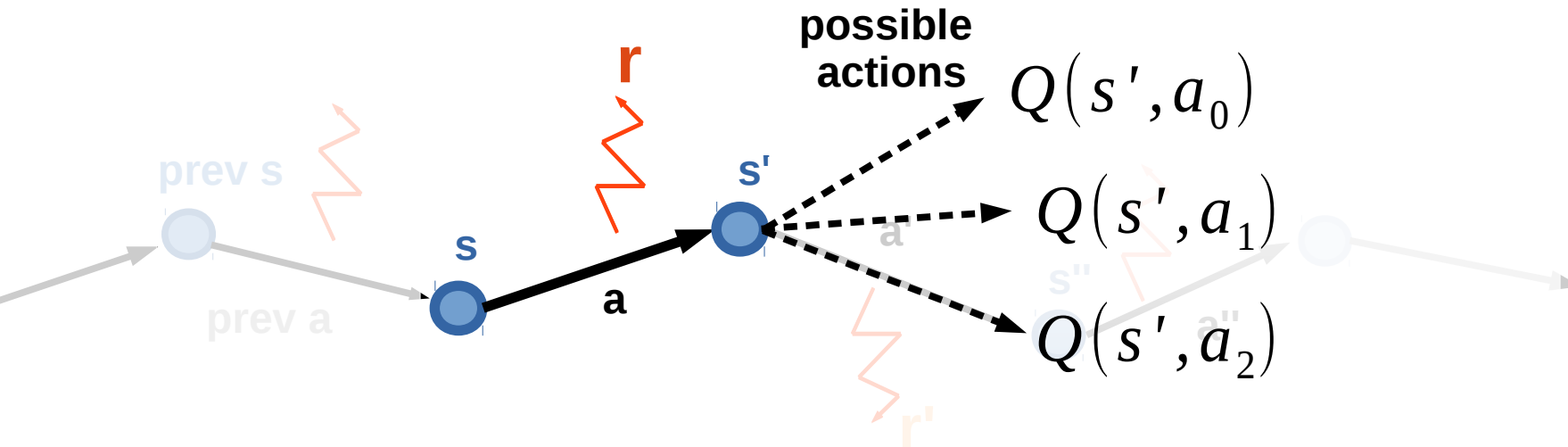
– Sample  $\langle \mathbf{s}, \mathbf{a}, \mathbf{r}, \mathbf{s}', \mathbf{a}' \rangle$  from env

hence “SARSA”

– Compute  $\hat{Q}(s, a) = r(s, a) + \gamma \underline{Q(s', a')}$  **next action (not max)**

– Update  $Q(s, a) \leftarrow \alpha \cdot \hat{Q}(s, a) + (1 - \alpha) Q(s, a)$

# Expected value SARSA



$$\forall s \in S, \forall a \in A, Q(s, a) \leftarrow 0$$

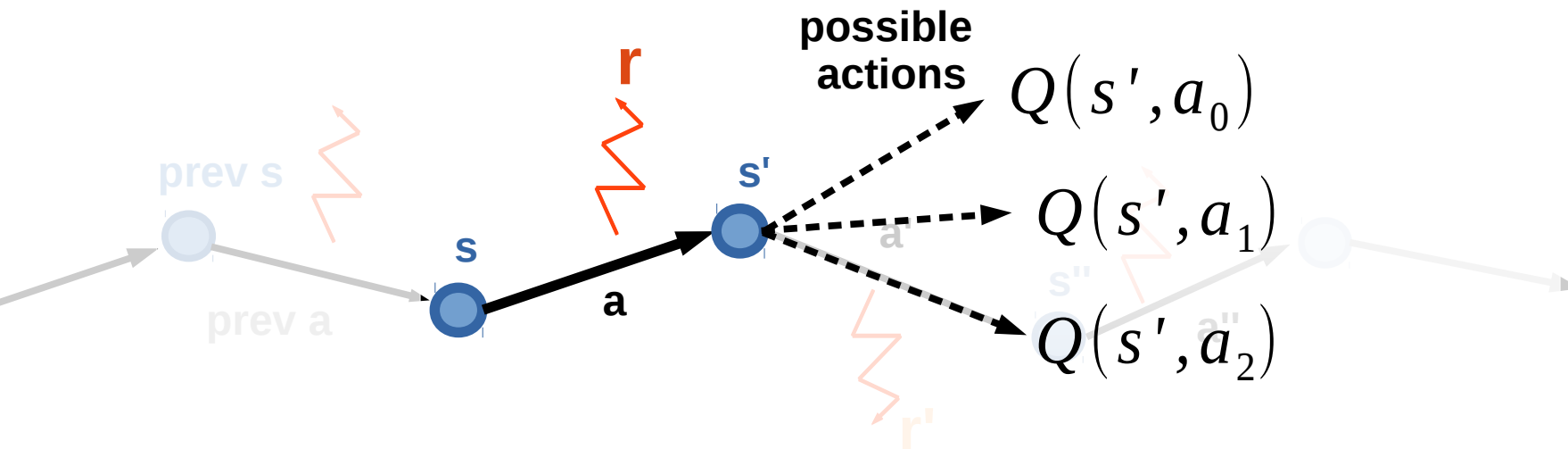
Loop:

- Sample  $\langle \mathbf{s}, \mathbf{a}, \mathbf{r}, \mathbf{s}' \rangle$  from env

- Compute  $\hat{Q}(s, a) = r(s, a) + \gamma \mathop{E}_{a_i \sim \pi(a|s')} Q(s', a_i)$

- Update  $Q(s, a) \leftarrow \alpha \cdot \hat{Q}(s, a) + (1 - \alpha) Q(s, a)$

# Expected value SARSA



$$\forall s \in S, \forall a \in A, Q(s, a) \leftarrow 0$$

Loop:

- Sample  $\langle \mathbf{s}, \mathbf{a}, \mathbf{r}, \mathbf{s}' \rangle$  from env

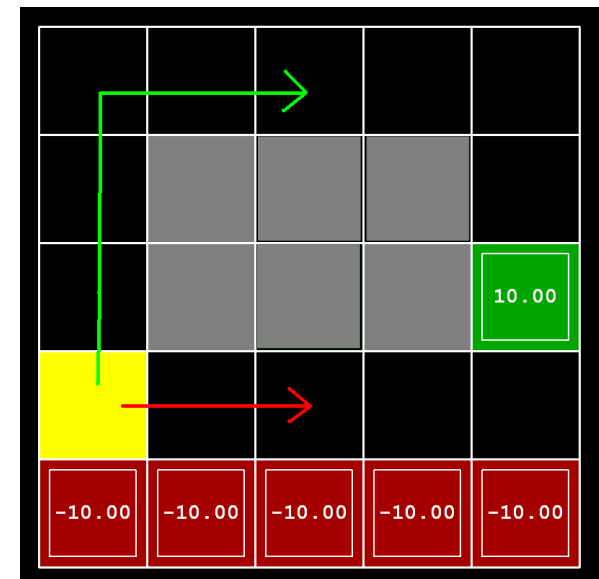
- Compute  $\hat{Q}(s, a) = r(s, a) + \gamma \underset{a_i \sim \pi(a|s')}{E} Q(s', a_i)$

- Update  $Q(s, a) \leftarrow \alpha \cdot \hat{Q}(s, a) + (1 - \alpha) Q(s, a)$



# Difference

- SARSA gets optimal rewards under current policy
- Q-learning policy **would be** optimal under



→ Q-learning  
→ SARSA

# On-policy vs Off-policy

## Two problem setups

### on-policy

Agent **can** pick actions

- Most obvious setup :)
- Agent always follows his **own** policy

### off-policy

Agent **can't** pick actions

- Learning with exploration,  
playing without exploration
- Learning from expert  
(expert is imperfect)
- Learning from sessions  
(recorded data)

# On-policy vs Off-policy

Two problem setups

**on-policy**

Agent **can** pick actions

- On-policy algorithms **can't** learn off-policy

**off-policy**

Agent **can't** pick actions

- Off-policy algorithms **can** learn on-policy

learn optimal policy even if agent takes random actions

**Q:** which of Q-learning, SARSA and exp. val. SARSA will **only** work on-policy?

# On-policy vs Off-policy

## Two problem setups

### on-policy

Agent **can** pick actions

- On-policy algorithms **can't** learn off-policy
- SARSA
- more later

### off-policy

Agent **can't** pick actions

- Off-policy algorithms **can** learn on-policy
- Q-learning
- Expected Value SARSA

# On-policy vs Off-policy

## Two problem setups

### on-policy

Agent **can** pick actions

- On-policy algorithms **can't** learn off-policy
- SARSA
- more coming soon

### off-policy

Agent **can't** pick actions

- Off-policy algorithms **can** learn on-policy
- Q-learning
- Expected Value SARSA

# On-policy vs Off-policy

## Two problem setups

### on-policy

Agent **can** pick actions

- On-policy algorithms **can't** learn off-policy
- SARSA
- more coming soon

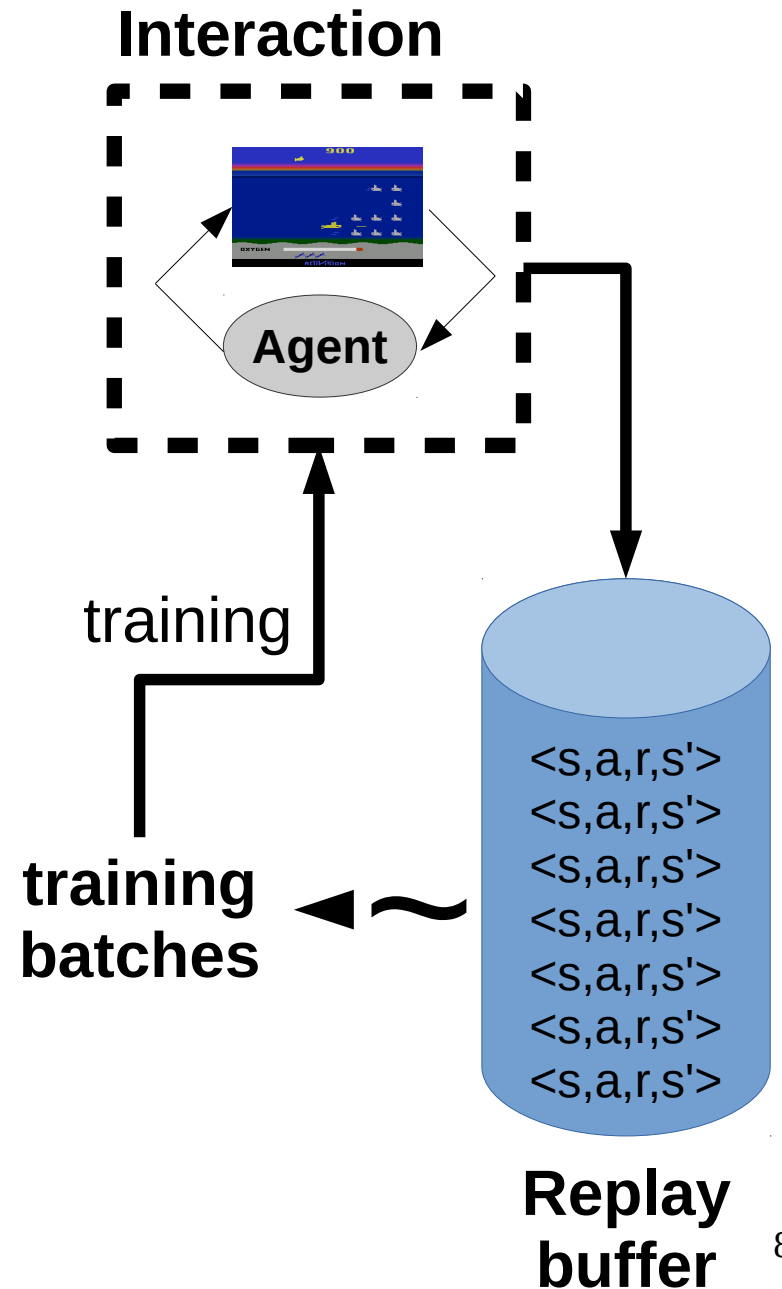
### off-policy

Agent **can't** pick actions

- Off-policy algorithms **can** learn on-policy
- Q-learning
- Expected Value SARSA

# Experience replay

**Idea:** store several past interactions  
 $\langle s, a, r, s' \rangle$   
Train on random subsamples



# Experience replay

**Idea:** store several past interactions  
 $\langle s, a, r, s' \rangle$   
Train on random subsamples

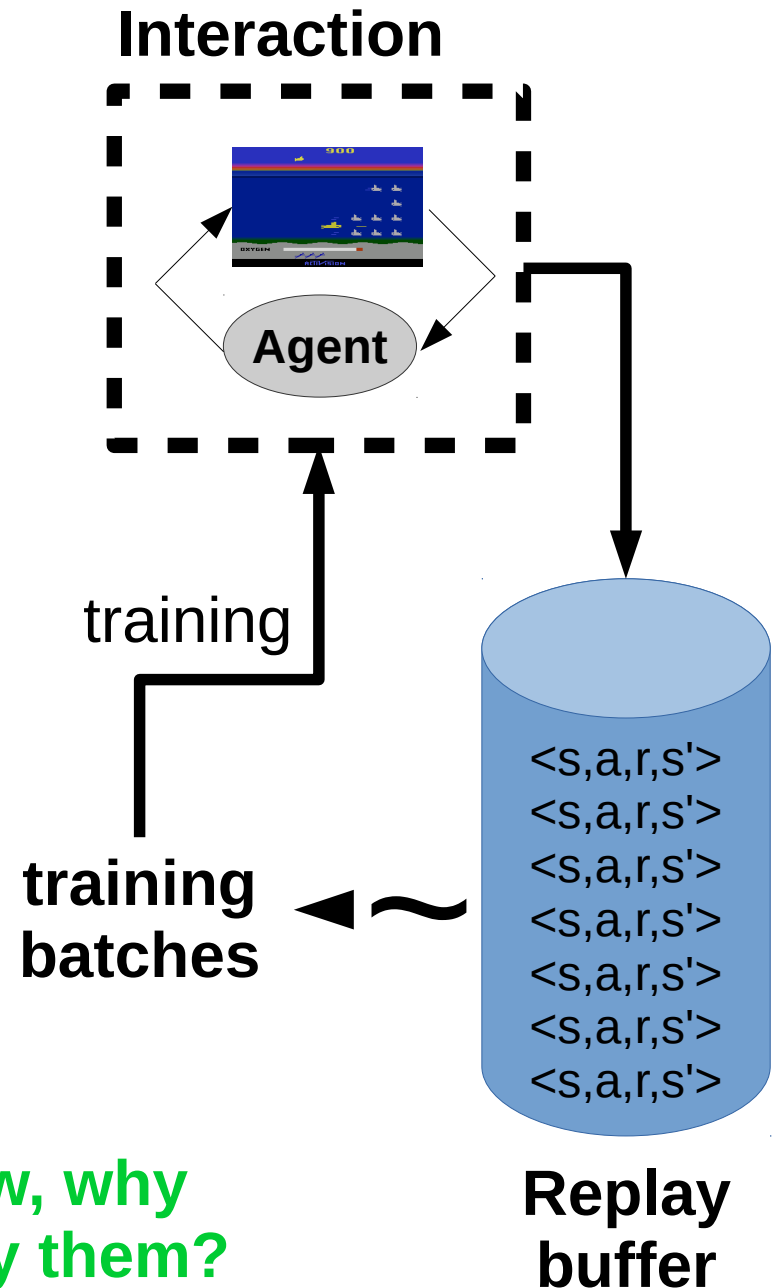
## Training curriculum:

- play 1 step and record it
- pick N random transitions to train

**Profit:** you don't need to re-visit same  
(s,a) many times to learn it.

**Only works with  
off-policy algorithms!**

**Btw, why  
only them?**





# Experience replay

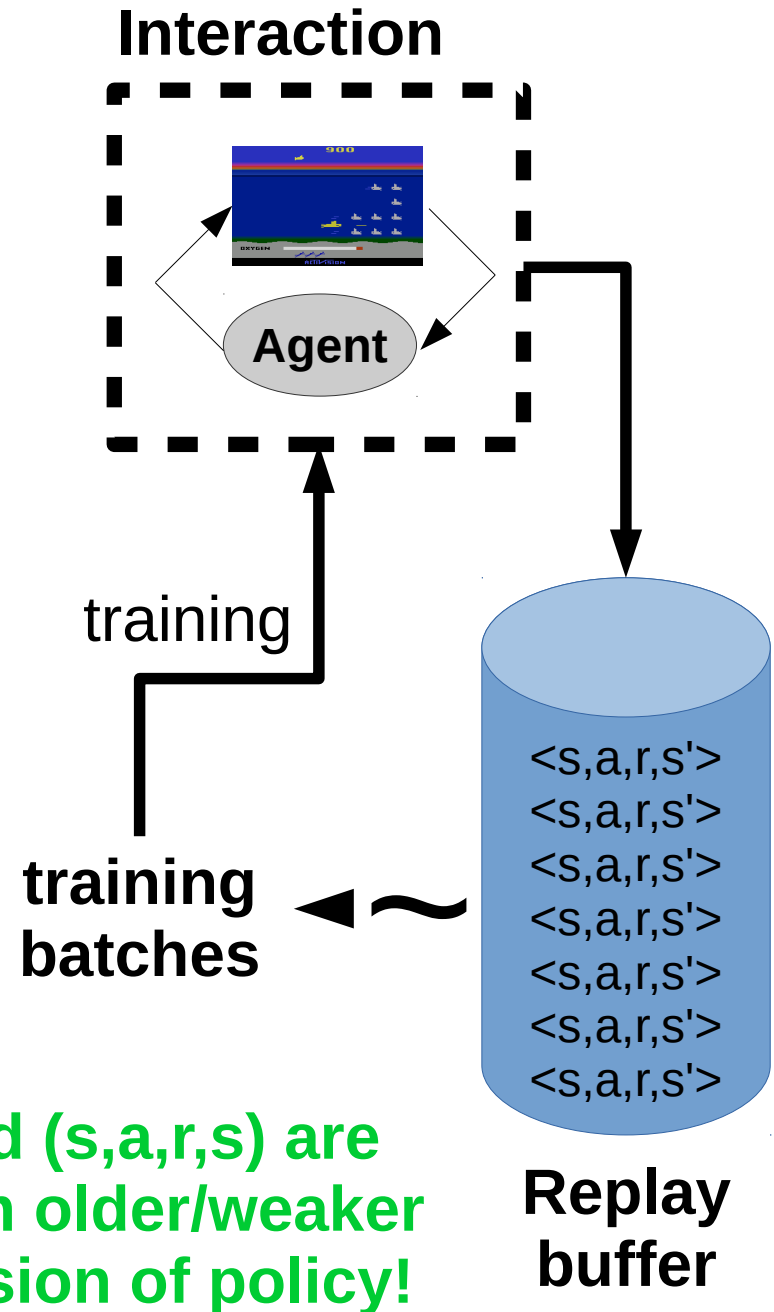
**Idea:** store several past interactions  
 $\langle s, a, r, s' \rangle$   
Train on random subsamples

## Training curriculum:

- play 1 step and record it
- pick N random transitions to train

**Profit:** you don't need to re-visit same  
(s,a) many times to learn it.

**Only works with  
off-policy algorithms!**



# New stuff we learned

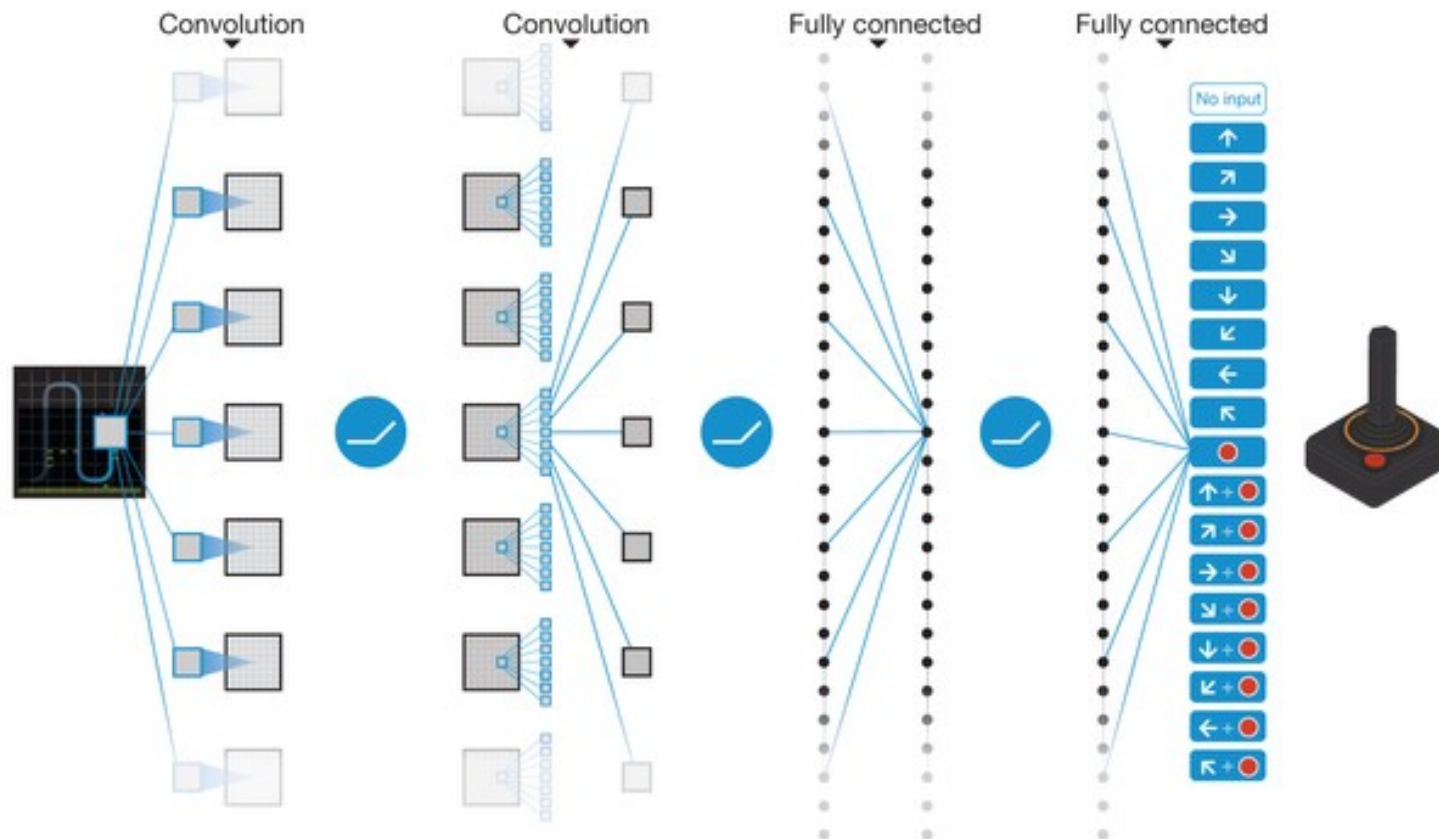
- Anything?

# New stuff we learned

- $V(s)$ ,  $V^*(s)$ ,  $Q(s,a)$ ,  $Q^*(s,a)$
- Value iteration (model-based)
- Q-learning, SARSA (model-free)
- Exploration vs exploitation (basics)
- Learning On-policy vs Off-policy
  - Using experience replay

# Coming next...

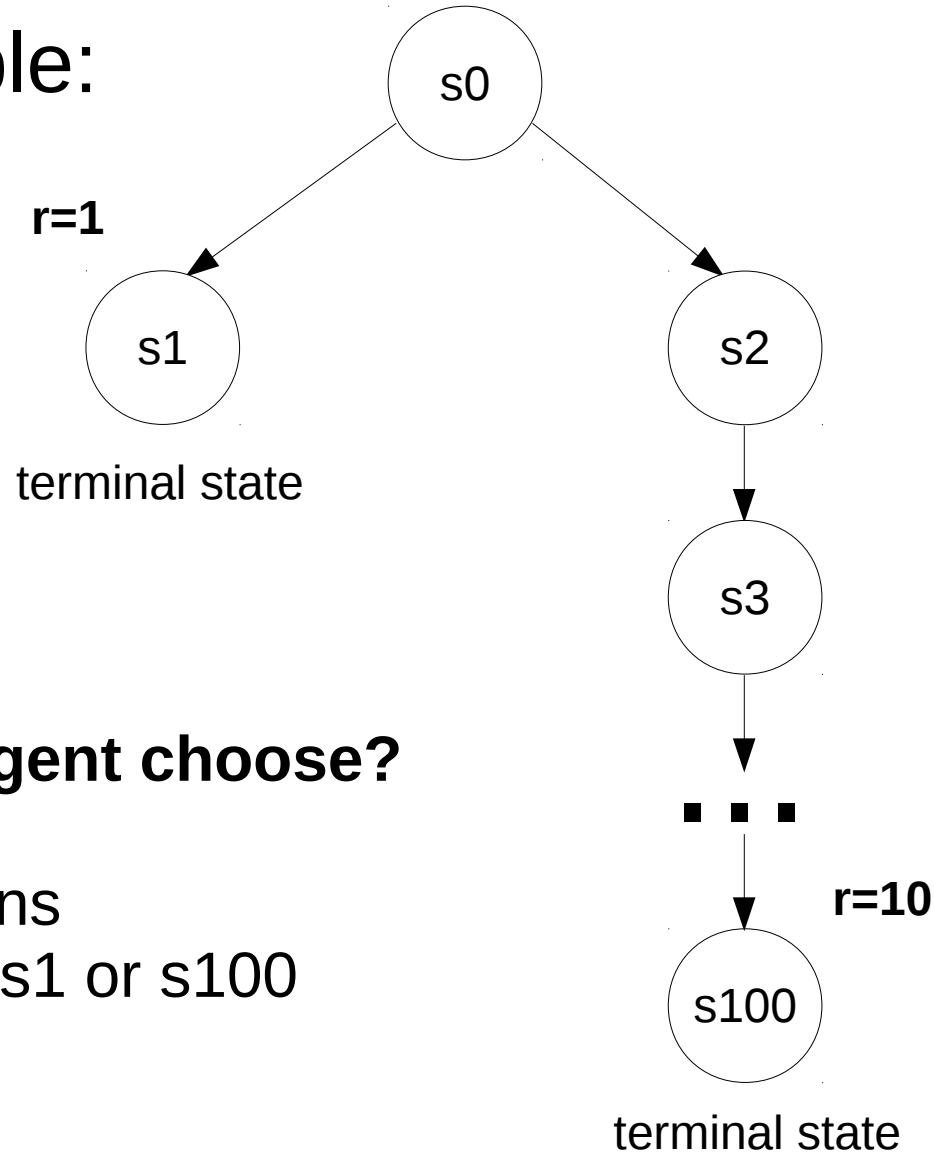
- What if state space is large/continuous
  - Deep reinforcement learning



Remember discounted rewards?

# Discounted reward **fails** #1

Trivial example:

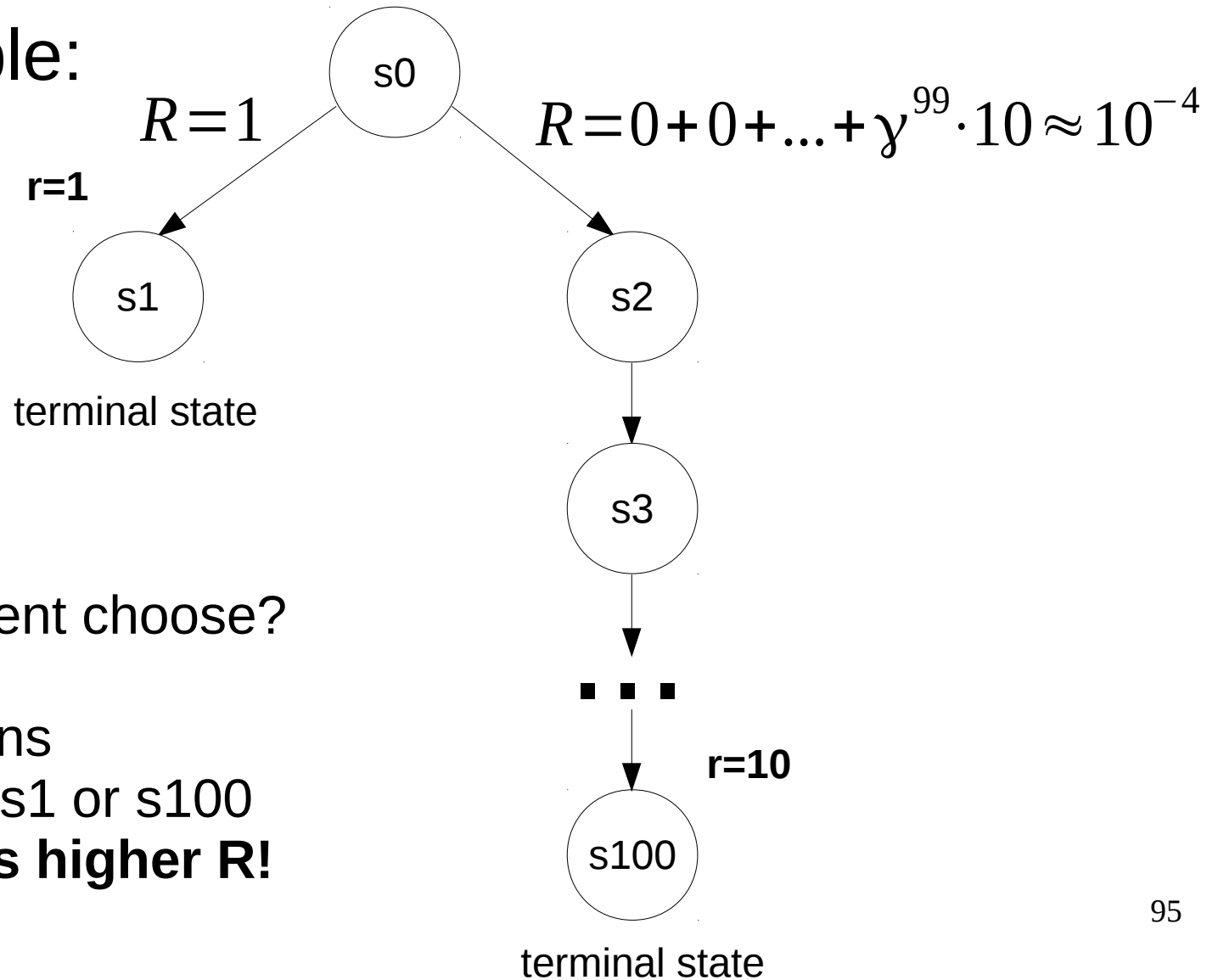


**What path will agent choose?**

- $\gamma=0.9$
- arrows = actions
- game ends at s1 or s100

# Discounted reward **fails** #1

Trivial example:



What path will agent choose?

- $\gamma=0.9$
- arrows = actions
- game ends at  $s_1$  or  $s_{100}$
- **left action has higher  $R$ !**

# Discounted reward **fails** #2

## Deephack'17 qualification round, Atari Skiing



- You steer the red guy
- Session lasts ~5k steps
- You get -3~-7 reward each tick (faster game = better score)
- At the end of session, you get up to  $r=-30k$  (based on passing gates, etc.)
- Q-learning with  $\gamma=0.99$  fails it doesn't learn to pass gates

**What's the problem?**



# Discounted reward **fails** #2

## Deephack'17 qualification round, Atari Skiing



- You steer the red guy
- Session lasts ~5k steps
- You get -3~-7 reward each tick (faster game = better score)
- At the end of session, you get up to  $r=-30k$  (based on passing gates, etc.)
- Q-learning with  $\gamma=0.99$  fails

# Discounted reward **fails** #3

## CoastRunner7 experiment (openAI)



- You control the boat
- Rewards for getting to checkpoints
- Rewards for collecting bonuses
- What could possibly go wrong?
- “Optimal” policy video:  
<https://www.youtube.com/watch?v=tlOIHko8ySg>

# Nuts and bolts: MC vs TD

## Monte-carlo

- Ignores intermediate rewards  
doesn't need  $\gamma$  (discount)
- Needs full episode to learn  
Infinite MDP are a problem
- Doesn't use Markov property  
Works with non-markov envs

## Temporal Difference

- Uses intermediate rewards  
trains faster under right  $\gamma$
- Learns from incomplete episode  
Works with infinite MDP
- Requires markov property  
Non-markov env is a problem



# Nuts and bolts: discount

- Effective horizon  $1 + \gamma + \gamma^2 + \dots = \frac{1}{(1 - \gamma)}$

Heuristic: your agent stops giving a damn in *this many* turns.

Typical values:

- $\gamma=0.9$ , 10 turns
- $\gamma=0.95$ , 20 turns
- $\gamma=0.99$ , 100 turns
- $\gamma=1$ , infinitely long

Higher  $\gamma$  = less stable algorithm.

$\gamma=1$  only works for episodic MDP (finite amount of turns).

# Nuts and bolts: discount

- Effective horizon  $1 + \gamma + \gamma^2 + \dots = \frac{1}{(1 - \gamma)}$

Heuristic: your agent stops giving a damn in *this many* turns.

- Atari Skiing, reward was delayed by in 5k steps
- $\gamma=0.99$  is not enough
- $\gamma=1$  and a few hacks works better
- Or use a better reward function



# #ShamelessSelfPromotion

- More stuff
  - [https://github.com/yandexdataschool/practical\\_rl/tree/fall17](https://github.com/yandexdataschool/practical_rl/tree/fall17)
  - Videos lectures, more theory, more coding
    - Model-based – week2
    - Mode-free – week3

Let's write some code!