# Temporal Difference & Q Learning
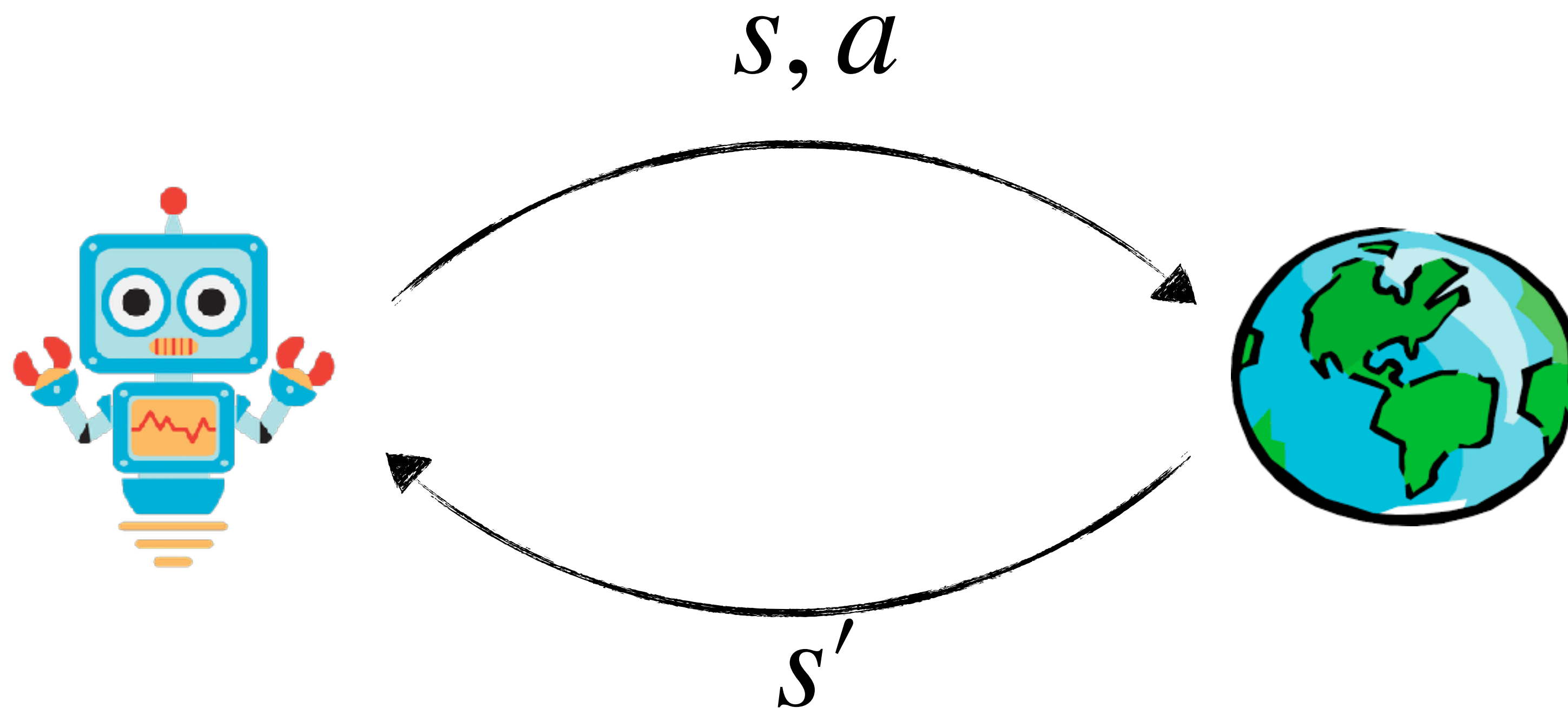
Sanjiban Choudhury
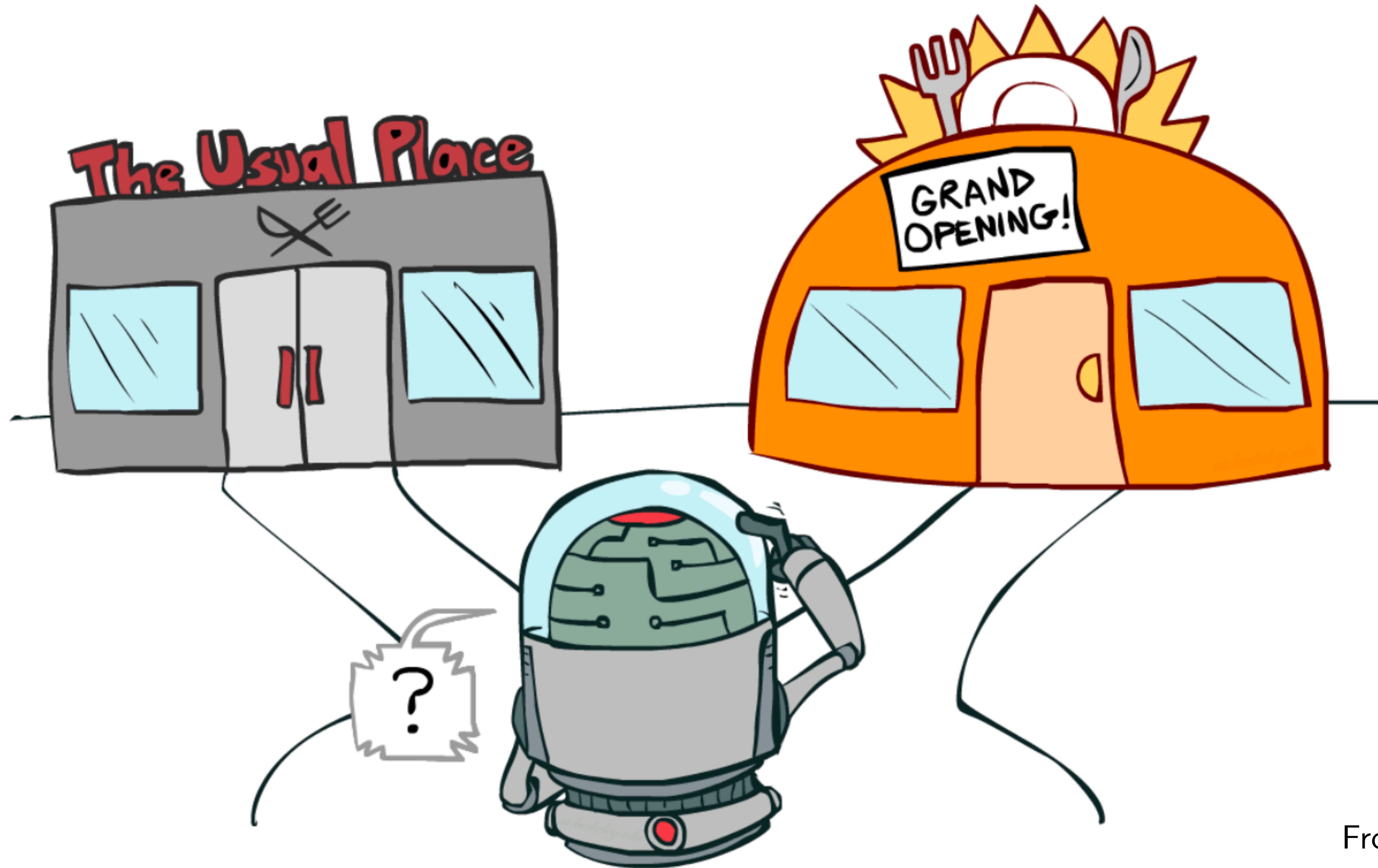
# What if the transitions are unknown?

$$< S, A, C, \mathcal{T} >$$

$$s, a$$

$$s'$$

# Exploration vs Exploitation



From Dan Klein

3

Doys



Doors

$a^1$     ?

$a^2$     ?

$a^3$     ?

-100

-1

1000

4

Doors

Round 1  Round 2  Round 3

$a^1$

$a^2$

$a^3$

-100

-1

1000
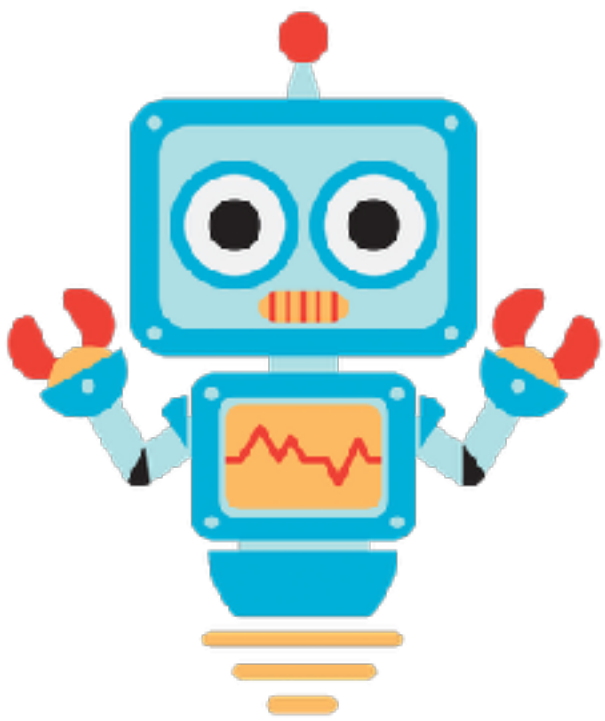
5

Doors

Round 1     Round 2     Round 3

$a^1$

$a^2$

$a^3$

-100

-1
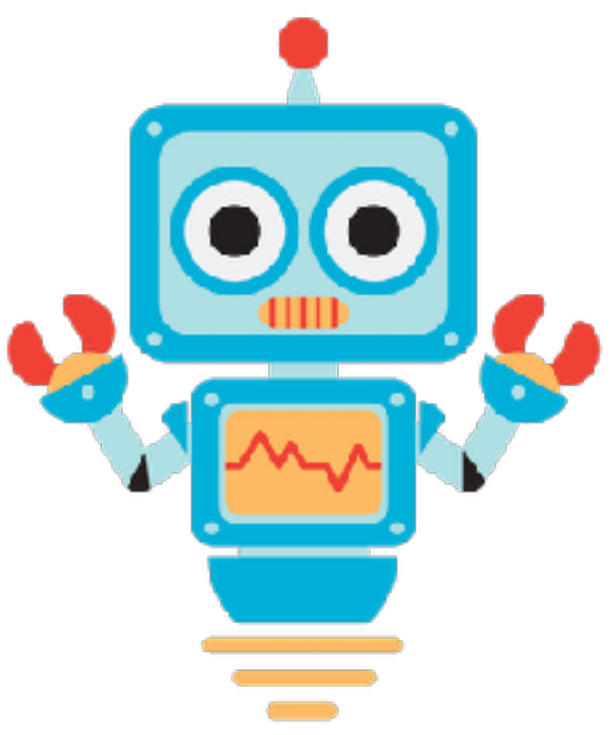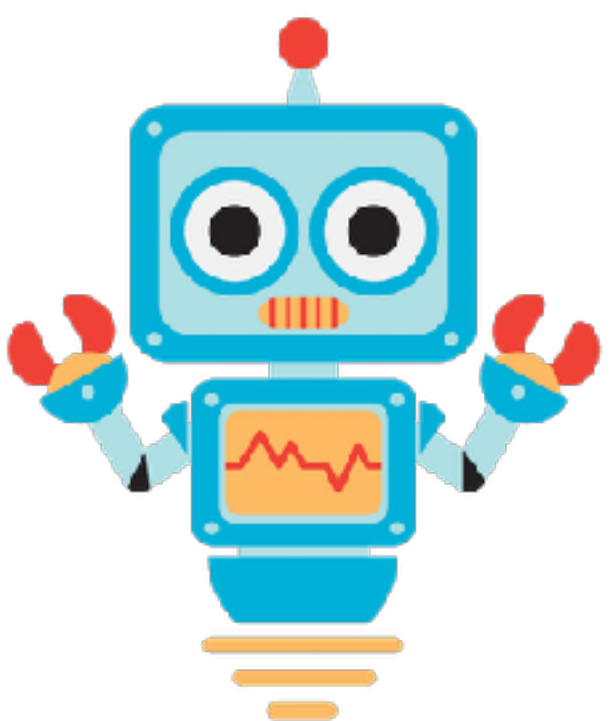
1000
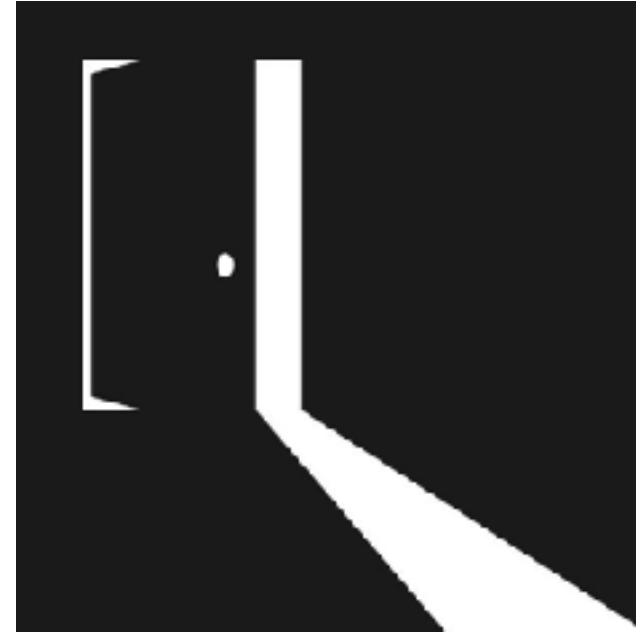
Doors

Round 1    Round 2    Round 3

$a^1$

$a^2$

$a^3$

How do we explore/ exploit when picking doors?

What if we played the game over multiple time steps?

$t = 1$

$t = 2$

-100

-1

1000

How do we estimate values of each door?

$t = 1$                    $t = 2$

# Two Ingredients of RL



Exploration Exploitation

$s$ $a^1$

$a^2$

$a^3$

Estimate Values $Q(s, a)$

# Recap: The Swamp MDP



Swamp

$$< S , A , C , \mathcal{T} >$$

- Two absorbing states: Goal and Swamp
- Cost of each state is 1 till you reach the goal
- Let's set T = 30

# When the MDP is known!

## Run Value / Policy Iteration

# When MDP is known: Policy Iteration



Iter: 0

$$V^\pi(s) = c(s, \pi(s)) + \gamma \mathbb{E}_{s' \sim \mathcal{T}(s,a)} V^\pi(s')]$$

Estimate value

$$\pi^+(s) = \arg\min_a c(s, a) + \gamma \mathbb{E}_{s' \sim \mathcal{T}(s,a)} V^\pi(s')]$$

Improve policy

What happens when the
MDP is *unknown?*

# Need to *estimate the value* of policy



Value $V^{\pi}(s)$

Policy $\pi$

# *Estimate the value* of policy from sample rollouts



Iter: 0

Roll outs

Policy $\pi$

# *Estimate the value* of policy from sample rollouts



| | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 74 | 75 | 76 | 77 | 77 | 77 | 77 | 2 | 1 | 0 |
| 1 | 74 | 75 | 76 | 77 | 77 | 77 | 77 | 3 | 2 | 1 |
| 2 | 74 | 75 | 76 | 77 | 77 | 77 | 77 | 3.9 | 3 | 2 |
| 3 | 55 | 56 | 56 | 57 | 50 | 40 | 26 | 4.9 | 3.9 | 3 |
| 4 | 74 | 75 | 76 | 77 | 77 | 77 | 77 | 5.9 | 4.9 | 3.9 |
| 5 | 74 | 75 | 76 | 77 | 77 | 77 | 77 | 6.8 | 5.9 | 4.9 |
| 6 | 74 | 75 | 76 | 77 | 77 | 77 | 77 | 7.7 | 6.8 | 5.9 |
| 7 | 15 | 14 | 13 | 12 | 11 | 10 | 9.6 | 8.6 | 7.7 | 6.8 |
| 8 | 16 | 15 | 14 | 13 | 12 | 11 | 10 | 9.6 | 8.6 | 7.7 |
| 9 | 17 | 16 | 15 | 14 | 13 | 12 | 11 | 10 | 9.6 | 8.6 |

Roll outs                    Value $V^\pi(s)$

18

# Activity!

# Think-Pair-Share

Think (30 sec): Given a bunch of roll-outs, how can you estimate value of a state? (Hint: More than one way!)

Pair: Find a partner

Share (45 sec): Partners exchange
ideas



Roll outs

Value $V^\pi(s)$

# Option 1: Just execute the damn policy!



and look at the returns ..

# Monte Carlo Evaluation

Goal: Learn $V^\pi(s)$ from complete rollout $\quad s_1, a_1, c_1, s_2, a_2, c_2, \ldots \sim \pi$

Define: *Return* is the total discounted cost

$$G_t = c_{t+1} + \gamma c_{t+2} + \gamma^2 c_{t+3} + \ldots$$

Value function is the expected return

$$V^\pi(s) = \mathbb{E}_\pi[G_t \mid s_t = s]$$

# Monte Carlo



For episode in rollouts:

 Increment counter $N(s) \leftarrow N(s) + 1$

 Increment total return
 $S(s) \leftarrow S(s) + G_t$

 Update $V(s) = S(s)/N(s)$

Law of large numbers: $V(s) \rightarrow V^{\pi}(s)$ as $N(s) \rightarrow \infty$

# Monte Carlo



For episode in rollouts:

Increment counter $N(s) \leftarrow N(s) + 1$

Increment total return
$S(s) \leftarrow S(s) + G_t$

Update $V(s) = S(s)/N(s)$

Law of large numbers: $V(s) \to V^{\pi}(s)$ as $N(s) \to \infty$

# Monte Carlo

$V(s)$

| | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| | | | | | | | | | 0 |
| | | | | | | | | | 1 |
| | | | | | | | | | 2 |
| 20.5 | 19.5 | 18.5 | 17.5 | 16.5 | 6 | 5 | 4 | 3 | |
| | | | | 25 | | | | | |

For episode in rollouts:

    Increment counter $N(s) \leftarrow N(s) + 1$

    Increment total return
$S(s) \leftarrow S(s) + G_t$

    Update $V(s) = S(s)/N(s)$

Law of large numbers: $V(s) \rightarrow V^{\pi}(s)$ as $N(s) \rightarrow \infty$

# Exponential Moving Average MC

$V(s)$ | 20.5 | 19.5 | 18.5 | 17.5 | 16.5 | 6 | 5 | 4 | 3

0
1
2

25

For episode in rollouts:

Update $V(s) \leftarrow V(s) + \alpha(G_t - V(s))$

Law of large numbers: $V(s) \rightarrow V^{\pi}(s)$ as $N(s) \rightarrow \infty$

# Can we do better than Monte Carlo?

What if we want quick updates?
(No patience to wait till end)


What if we don't have complete episodes?

# Option 2: Trust your value estimate

# Temporal Difference (TD) learning

Goal: Learn $V^\pi(s)$ from traces

$$\left(s_t, a_t, c_t, s_{t+1}\right) \qquad \left(s_t, a_t, c_t, s_{t+1}\right) \qquad \left(s_t, a_t, c_t, s_{t+1}\right) \qquad \left(s_t, a_t, c_t, s_{t+1}\right)$$

Recall value function $V^\pi(s)$ satisfies

$$V^\pi(s) = c(s, \pi(s)) + \gamma \mathbb{E}_{s'} V^\pi(s')$$

TD Idea: Update value using estimate of next state value

$$V(s_t) \leftarrow V(s_t) + \alpha \left( c_t + \gamma V(s_{t+1}) - V(s_t) \right)$$

Temporal Difference Error

# TD Learning

For every $\left(s_t, a_t, c_t, s_{t+1}\right)$

$$V(s_t) \leftarrow V(s_t) + \alpha(c_t + \gamma V(s_{t+1}) - V(s_t))$$

# Did you spot the trick?

$$V^\pi(s) = c(s, \pi(s)) + \gamma \mathbb{E}_{s'} V^\pi(s')$$

$$\Updownarrow$$

$$V(s_t) \leftarrow V(s_t) + \alpha(c_t + \gamma V(s_{t+1}) - V(s_t))$$

# Monte-Carlo

$$V(s) \leftarrow V(s) + \alpha(G_t - V(s))$$

Zero Bias

High Variance

Always convergence
(Just have to wait till heat death of the universe)

# Temporal Difference

$$V(s) \leftarrow V(s) + \alpha(c + \gamma V(s') - V(s))$$

Can have bias

Low Variance

May *not* converge if
using function approximation

We have been talking about trying to learn the value of a given policy $\pi$

$$V^\pi(s) \, / \, Q^\pi(s, a)$$

What if we wanted to learn the optimal value function

$$V^*(s) \, / \, Q^*(s, a)$$

# Q-learning: Learning off-policy



$(s, a, s', c)$

dataset of transitions
("replay buffer")

$\pi(\mathbf{a}|\mathbf{s})$ (e.g., $\epsilon$-greedy)

off-policy
Q-learning

For every $\left(s_t, a_t, c_t, s_{t+1}\right)$

Can learn from any data!

$$Q^*(s_t, a_t) = Q^*(s_t, a_t) + \alpha(c(s_t, a_t) + \gamma \min_{a'} Q^*(s_{t+1}, a') - Q^*(s_t, a_t))$$

# Is this ... magic?

We just learned in IL how distribution shift is a big deal ...

It's not magic. Q-learning relies on a set of assumptions:

1. Each state-action is visited *infinite* times

2. Learning rate $\alpha$ must be annealed over time

# QT-Opt: Scalable Deep Reinforcement Learning for Vision-Based Robotic Manipulation

**Training time**

**Reward**: Grasp success determined by subtracting pre and post-drop images

**Distributed RL**

State, Action, Reward

Learned weights

**Inference time**

Camera

Robot

**State**: 472x472 Image and gripper aperture

**Action**: Gripper displacement and aperture

Action proposals

**Critic Function**
$Q(\text{State}, \text{Action})$

Q-Values

**Cross-Entropy Method**
$\arg\max_{\text{Action}} Q(\text{State}, \text{Action})$

# Q-learning: Learning off-policy



$(s, a, s', c)$

dataset of transitions
("replay buffer")

off-policy
Q-learning

$\pi(\mathbf{a}|\mathbf{s})$ (e.g., $\epsilon$-greedy)

For every $\left( s_t, a_t, c_t, s_{t+1} \right)$

Can learn from any data!

$$Q^*(s_t, a_t) = Q^*(s_t, a_t) + \alpha(c(s_t, a_t) + \gamma \min_{a'} Q^*(s_{t+1}, a') - Q^*(s_t, a_t))$$

# Large-scale Q-learning with continuous actions (QT-Opt)



stored data from all past experiments
$$\{(\mathbf{s}_i, \mathbf{a}_i, \mathbf{s}'_i)\}_i$$

**training buffers**

off-policy $(\mathbf{s}, \mathbf{a}, \mathbf{s}', r)$

on-policy $(\mathbf{s}, \mathbf{a}, \mathbf{s}', r)$

labeled $(\mathbf{s}, \mathbf{a}, Q_T(\mathbf{s}, \mathbf{a}))$

**Bellman updaters**

compute $Q_T(\mathbf{s}, \mathbf{a}) = r + \max_{\mathbf{a}'} Q_\theta(\mathbf{s}', \mathbf{a}')$

**training threads**

$$\min_\theta ||Q_\theta(\mathbf{s}, \mathbf{a}) - Q_T(\mathbf{s}, \mathbf{a})||^2$$

$Q_\theta$

$\mathbf{s}$

$\mathbf{a}$

live data collection

Kalashnikov, Irpan, Pastor, Ibarz, Herzong, Jang, Quillen, Holly, Kalakrishnan, Vanhoucke, Levine. **QT-Opt: Scalable Deep Reinforcement Learning of Vision-Based Robotic Manipulation Skills**

# Making Q-learning better!

Problem: Q-learning suffers from an estimation bias $\min_{a'} Q^*(s_{t+1}, a')$

Solution: Double Q-learning $Q^*(s_{t+1}, \arg\min_{a'} \tilde{Q}(s_{t+1}, a'))$
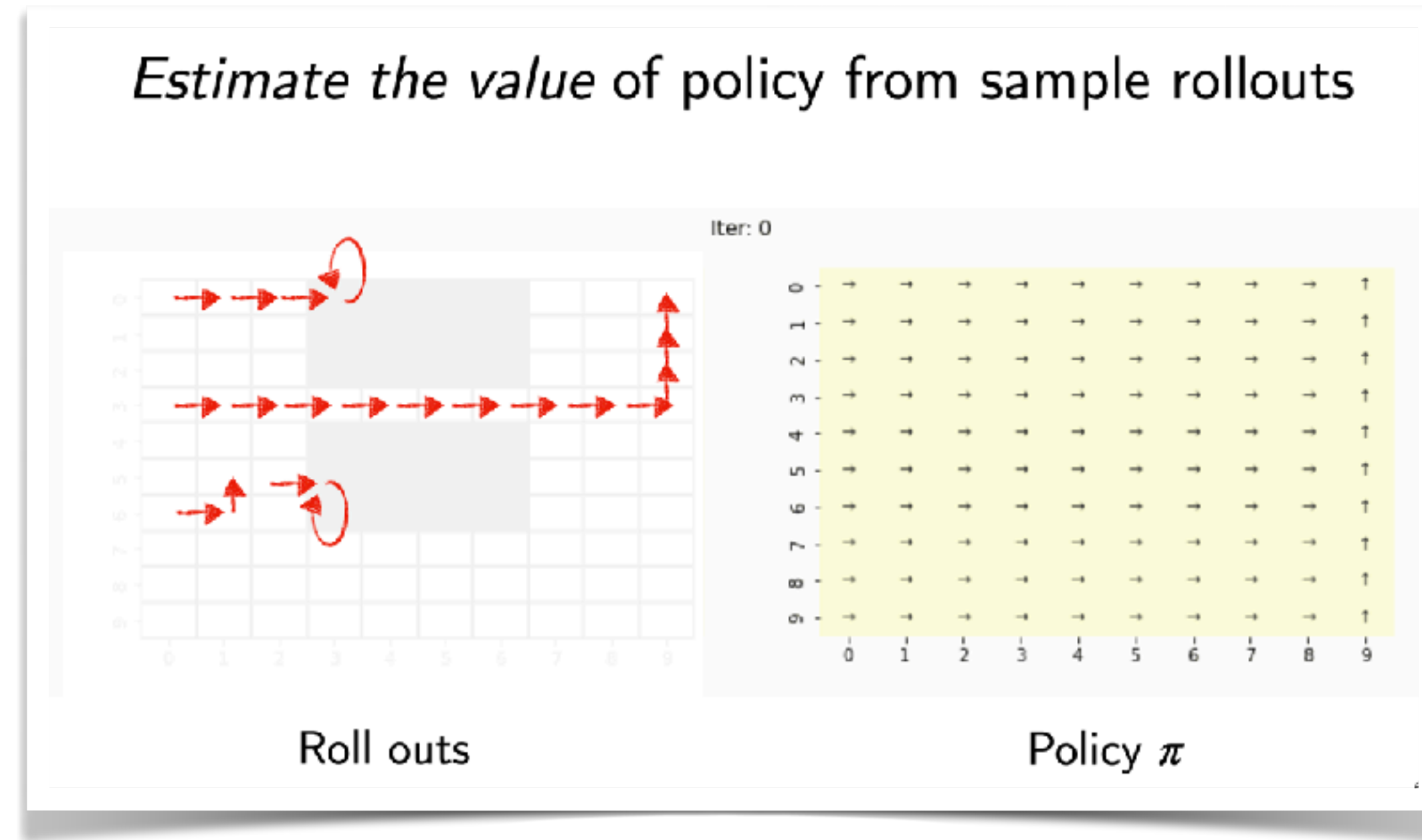
Problem: Q-learning samples uniformly from replay buffer

Solution: Prioritized DQN - samples states with higher bellman error

Problem: Q-learning doesn't seem to learn ….

Solution: Start with high exploration + learning rate, anneal!

# tl;dr

Estimate the value of policy from sample rollouts



Iter: 0

Roll outs | Policy $\pi$

| Monte-Carlo | Temporal Difference |
|---|---|
| $V(s) \leftarrow V(s) + \alpha(G_t - V(s))$ | $V(s) \leftarrow V(s) + \alpha(c + \gamma V(s') - V(s))$ |
| Zero Bias | Can have bias |
| High Variance | Low Variance |
| Always convergence<br>(Just have to wait till heat death of the universe) | May *not* converge if<br>using function approximation |

## Q-learning: Learning off-policy

For every $(s_t, a_t, c_t, s_{t+1})$

$$Q^*(s_t, a_t) = Q^*(s_t, a_t) + \alpha(c(s_t, a_t) + \gamma \min_{a'} Q^*(s_{t+1}, a') - Q^*(s_t, a_t))$$

Notice we are *not* approximating $Q^\pi(s_t, a_t)$

We don't even care about $\pi$

We can learn from any data!

28