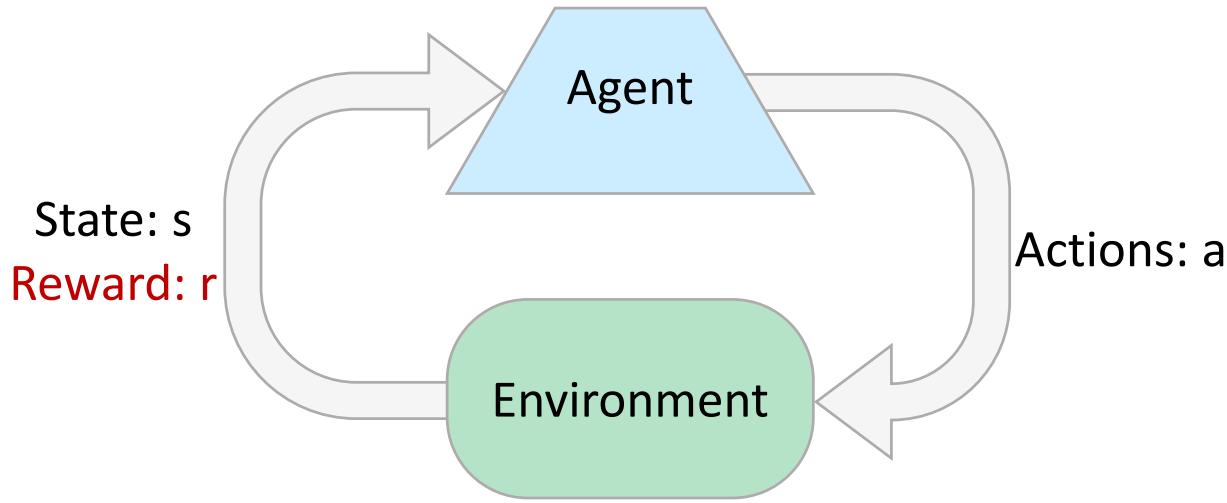


Search with Uncertainty



Reinforcement Learning

Review: Reinforcement Learning



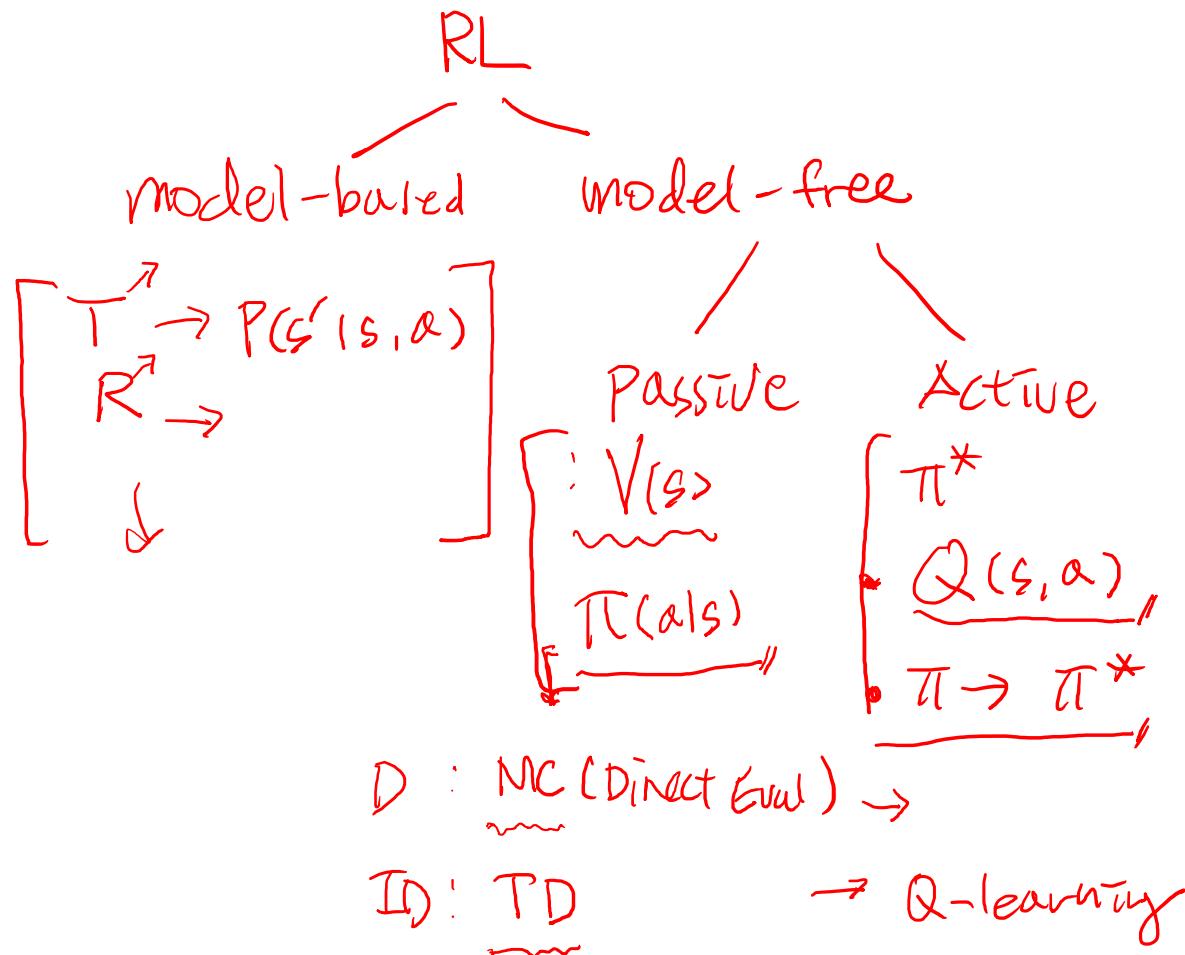
- Basic idea:
 - Receive feedback in the form of **rewards**
 - Agent's utility is defined by the reward function
 - Must (learn to) act so as to **maximize expected rewards**
 - All learning is based on observed samples of outcomes!



Review: 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

Review: Reinforcement Learning approaches



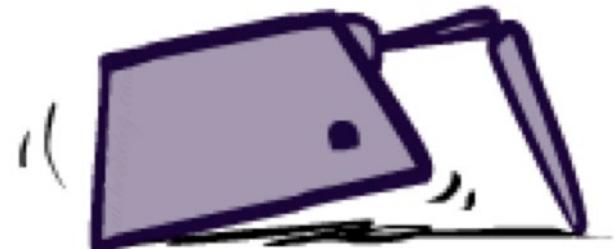
Review: Model-Based Learning

- 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 $\hat{T}(s, a, s')$ $P(s' | s, a)$
 - Discover each $\hat{R}(s, a, s')$ when we experience (s, a, s')

T, R



- Step 2: Solve the learned MDP
 - For example, use value iteration, as before

Review: Model-Free Learning

Passive

: To evaluate \mathbb{V}, π

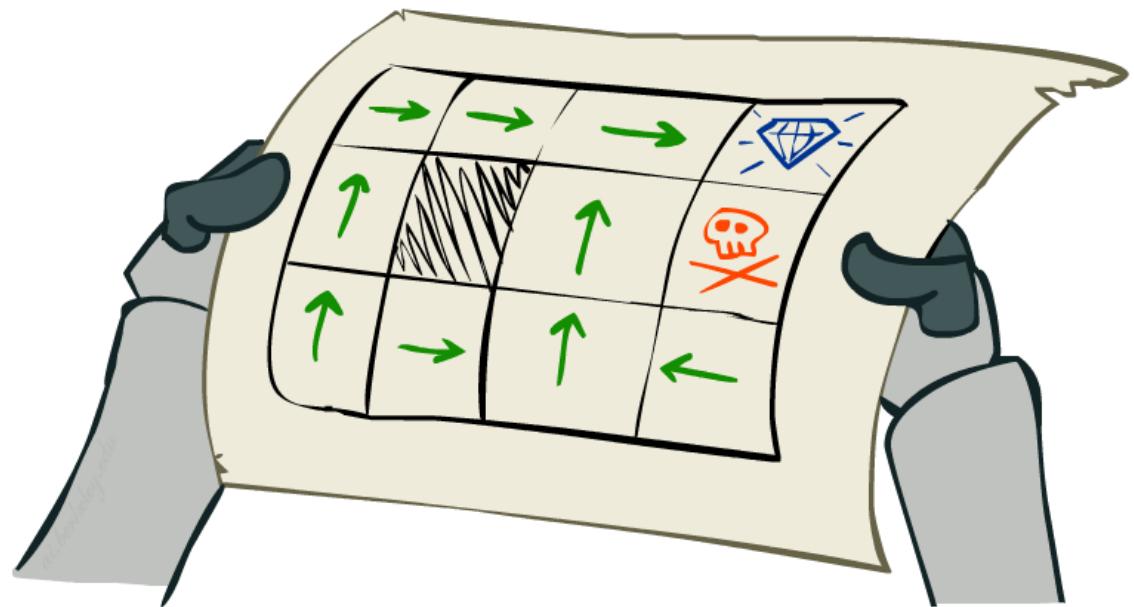
π : given, fixed

{ Direct eval
Temporal-Diff

Active

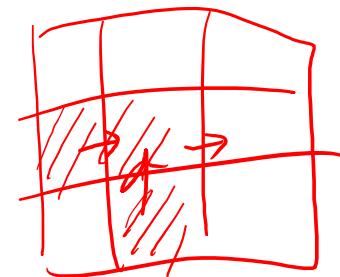
Review: Passive Reinforcement Learning

- 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.



Review: 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



Review: Sample-Based Policy Evaluation

- We want to improve our estimate of V by computing these averages:

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

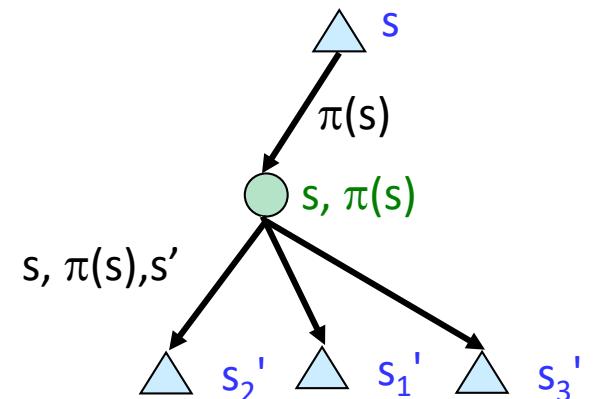
$s \rightarrow \text{Ter}$
 $s \rightarrow s'$

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

$$\left\{ \begin{array}{l} sample_1 = R(s, \pi(s), s'_1) + \gamma V_k^{\pi}(s'_1) \\ sample_2 = R(s, \pi(s), s'_2) + \gamma V_k^{\pi}(s'_2) \\ \dots \\ sample_n = R(s, \pi(s), s'_n) + \gamma V_k^{\pi}(s'_n) \end{array} \right.$$

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

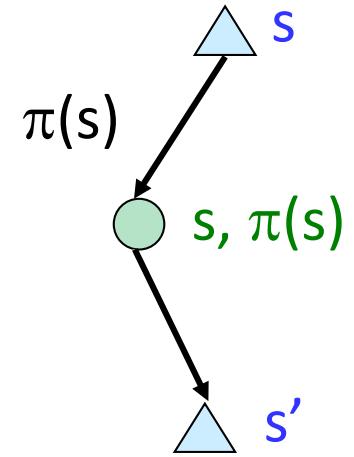
vs $\frac{1}{n} \sum_i P(a_i) A_i$



Almost! But we can't
rewind time to get sample
after sample from state s .

Review: 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

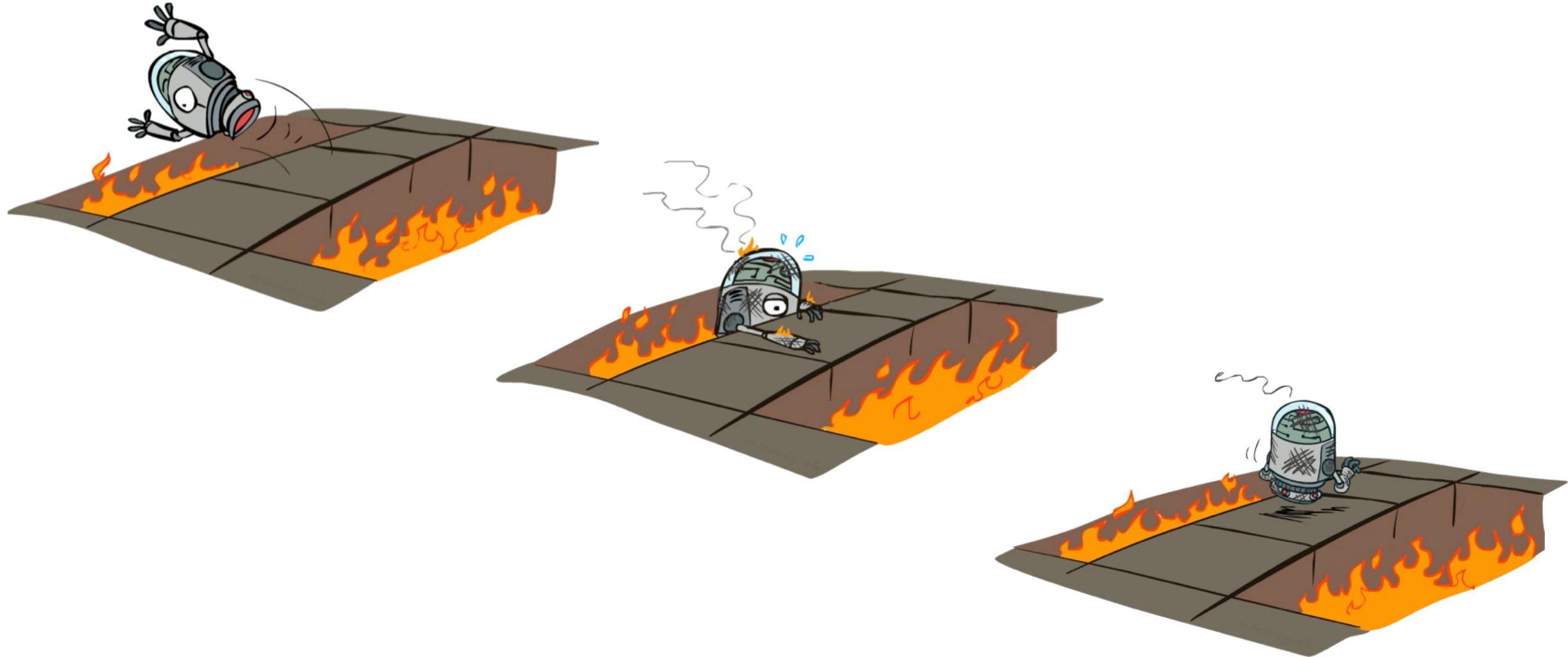


Sample of $V(s)$: $sample = \underbrace{R(s, \pi(s), s')}_{\text{Reward}} + \gamma \underbrace{V^\pi(s')}_{\text{Future Value}} \checkmark$

Update to $V(s)$: $V^\pi(s) \leftarrow \underbrace{(1 - \alpha)V^\pi(s)}_{\text{Old Value}} + \underbrace{(\alpha)sample}_{\text{New Value}}$

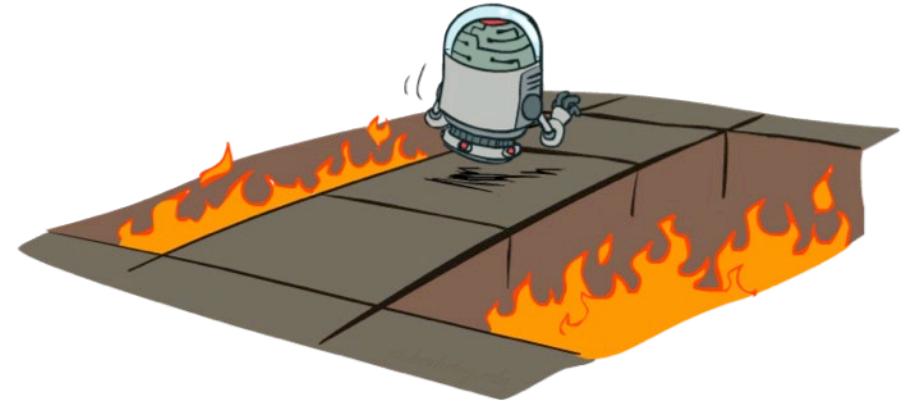
Same update: $V^\pi(s) \leftarrow V^\pi(s) + \alpha(sample - V^\pi(s))$

Active Reinforcement Learning



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:

$$V_{k+1}(s) \leftarrow \max_a \left[\sum_{s'} T(s, a, s') \left[R(s, a, s') + \gamma V_k(s') \right] \right]$$

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

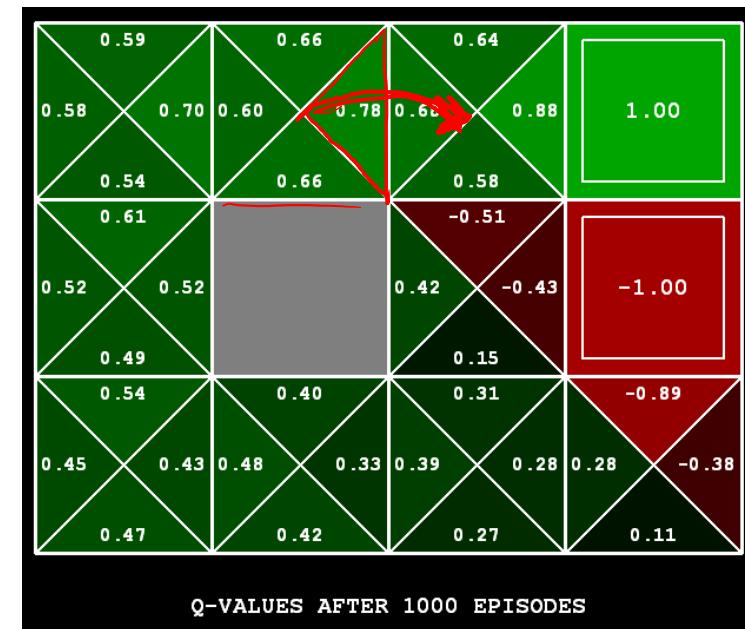
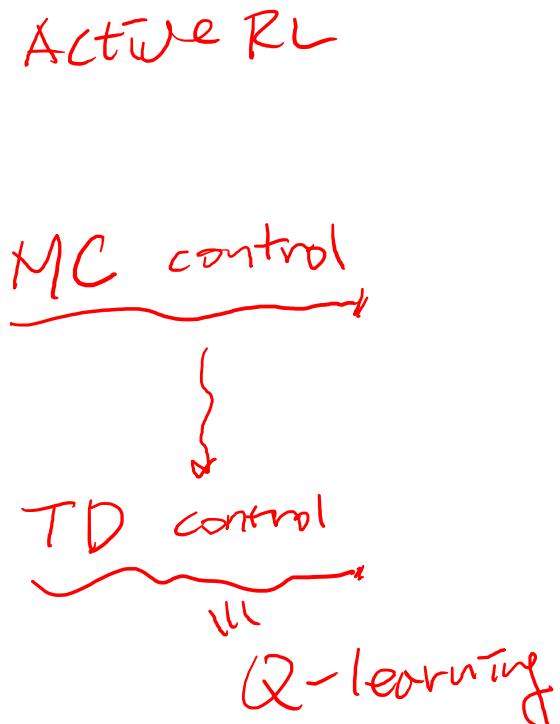
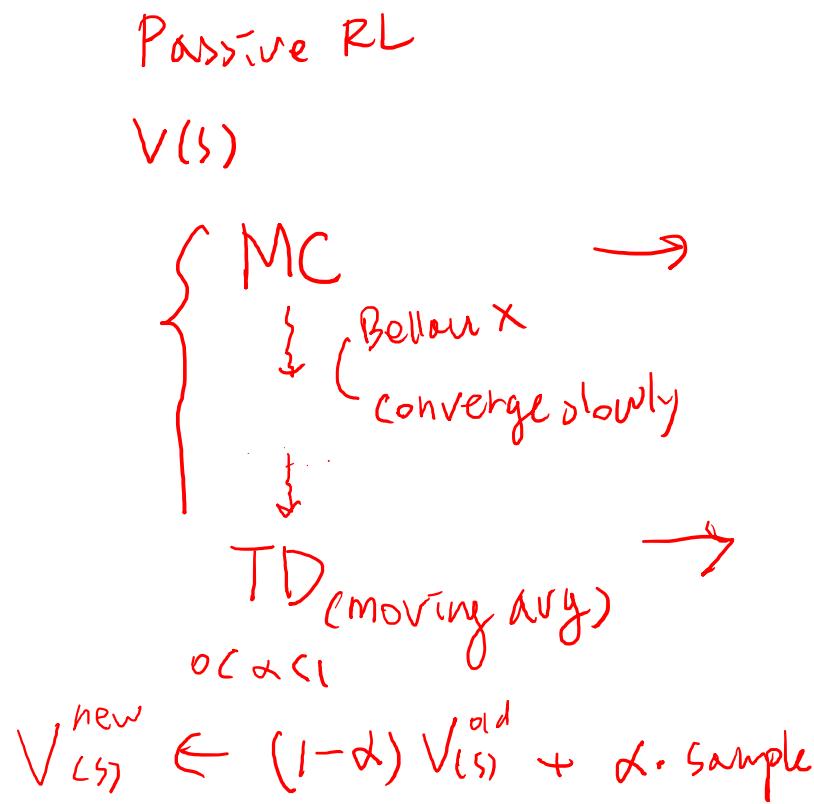
- Start with $Q_0(s, a) = 0$, which we know is right
- Given Q_k , calculate the depth $k+1$ q-values for all q-states:

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

$$Q = \sum_{s'} P(s' | s, a) V$$

Learning the Q-function

- Analogy to the value function learning in Passive RL



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

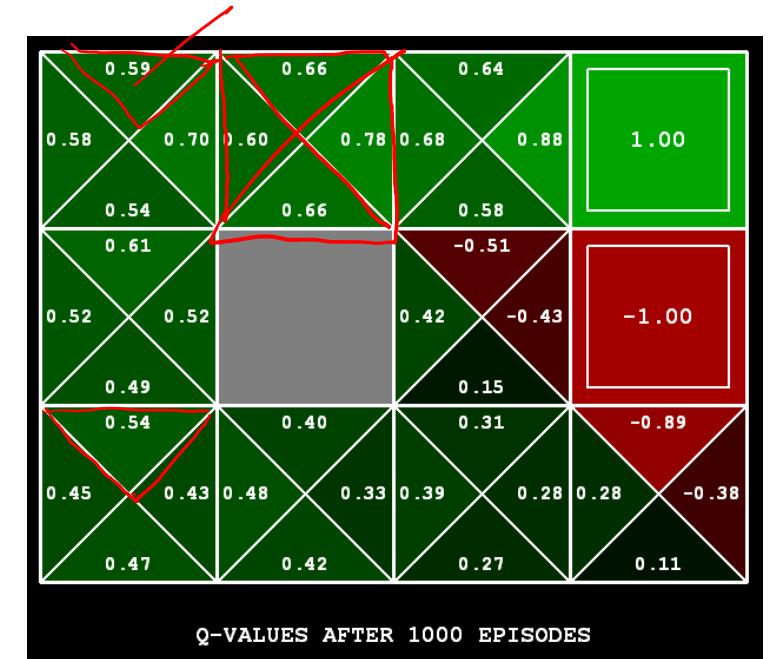
- Receive a sample (s, a, s', r)
- Consider your old estimate: $Q(s, a)$
- Consider your new sample estimate:

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

- Incorporate the new estimate into a running average:

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

π^*

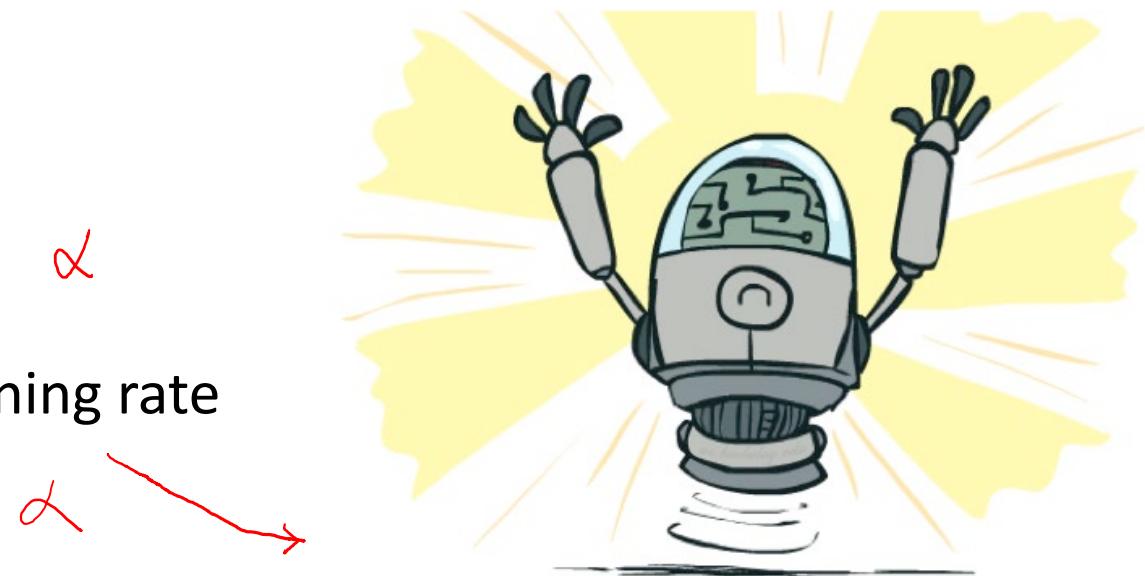


Video of Demo Q-Learning Auto Cliff Grid



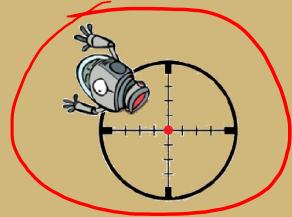
Q-Learning Properties

- 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 (!)



Note on Q-learning and Active RL

- Q-learning is actually a TD-learning with control
 $\pi \rightarrow \pi^*$
- There are more Active RL methods
(Q-learning is not the only one) $MC + control$
- Active RL allows the policy to be changed (improved)



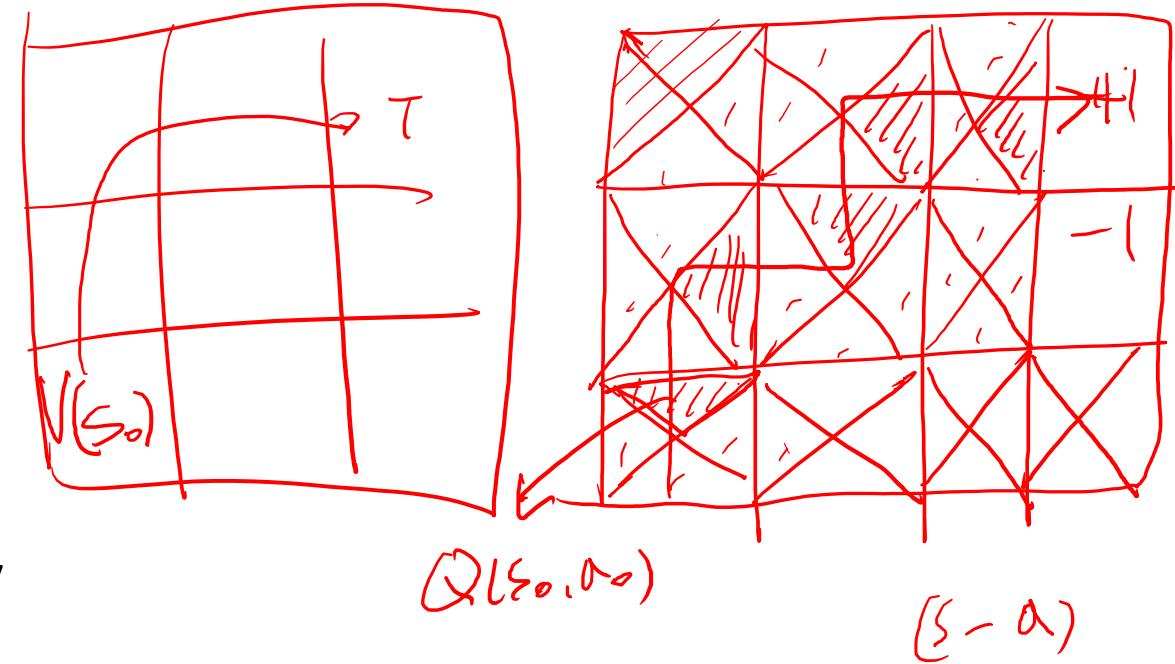
More Active RL: How about MC?

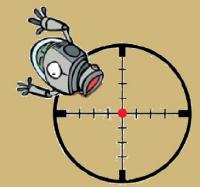
- Direct Evaluation (Monte Carlo evaluation of V)
 T_R)^x
- We can evaluate Q using MC
- Monte Carlo Control

Monte Carlo Control



- Exploring Starts
- Using ϵ -greedy policy





Monte Carlo Control Exploring Starts

Monte Carlo ES (Exploring Starts), for estimating $\pi \approx \pi_*$

Initialize:

$\pi(s) \in \mathcal{A}(s)$ (arbitrarily), for all $s \in \mathcal{S}$

$Q(s, a) \in \mathbb{R}$ (arbitrarily), for all $s \in \mathcal{S}, a \in \mathcal{A}(s)$

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

(ξ, α)

Loop forever (for each episode):

Choose $S_0 \in \mathcal{S}, A_0 \in \mathcal{A}(S_0)$ randomly such that all pairs have probability > 0

Generate an episode from S_0, A_0 , following π : $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T$

$G \leftarrow 0$

Loop for each step of episode, $t = T-1, T-2, \dots, 0$:

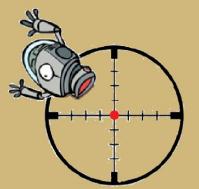
$G \leftarrow \gamma G + R_{t+1}$

Unless the pair S_t, A_t appears in $S_0, A_0, S_1, A_1, \dots, S_{t-1}, A_{t-1}$:

Append G to $Returns(S_t, A_t)$

$Q(S_t, A_t) \leftarrow \text{average}(Returns(S_t, A_t))$

$\pi(S_t) \leftarrow \arg \max_a Q(S_t, a)$



On-policy Monte Carlo Control (ε -greedy)

On-policy first-visit MC control (for ε -soft policies), estimates $\pi \approx \pi_*$

Algorithm parameter: small $\varepsilon > 0$

Initialize:

$\pi \leftarrow$ an arbitrary ε -soft policy

$Q(s, a) \in \mathbb{R}$ (arbitrarily), for all $s \in \mathcal{S}$, $a \in \mathcal{A}(s)$

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

Repeat forever (for each episode):

Generate an episode following π : $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T$

$G \leftarrow 0$

Loop for each step of episode, $t = T-1, T-2, \dots, 0$:

$$G \leftarrow \gamma G + R_{t+1}$$

Unless the pair S_t, A_t appears in $S_0, A_0, S_1, A_1, \dots, S_{t-1}, A_{t-1}$:

Append G to $\text{Returns}(S_t, A_t)$

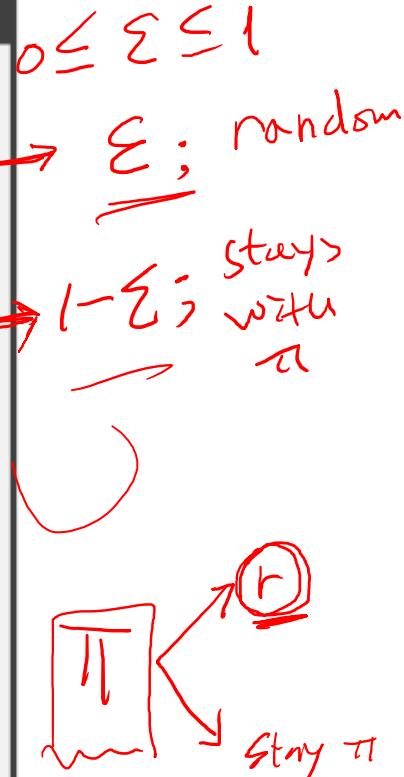
$Q(S_t, A_t) \leftarrow \text{average}(\text{Returns}(S_t, A_t))$

$A^* \leftarrow \arg \max_a Q(S_t, a)$

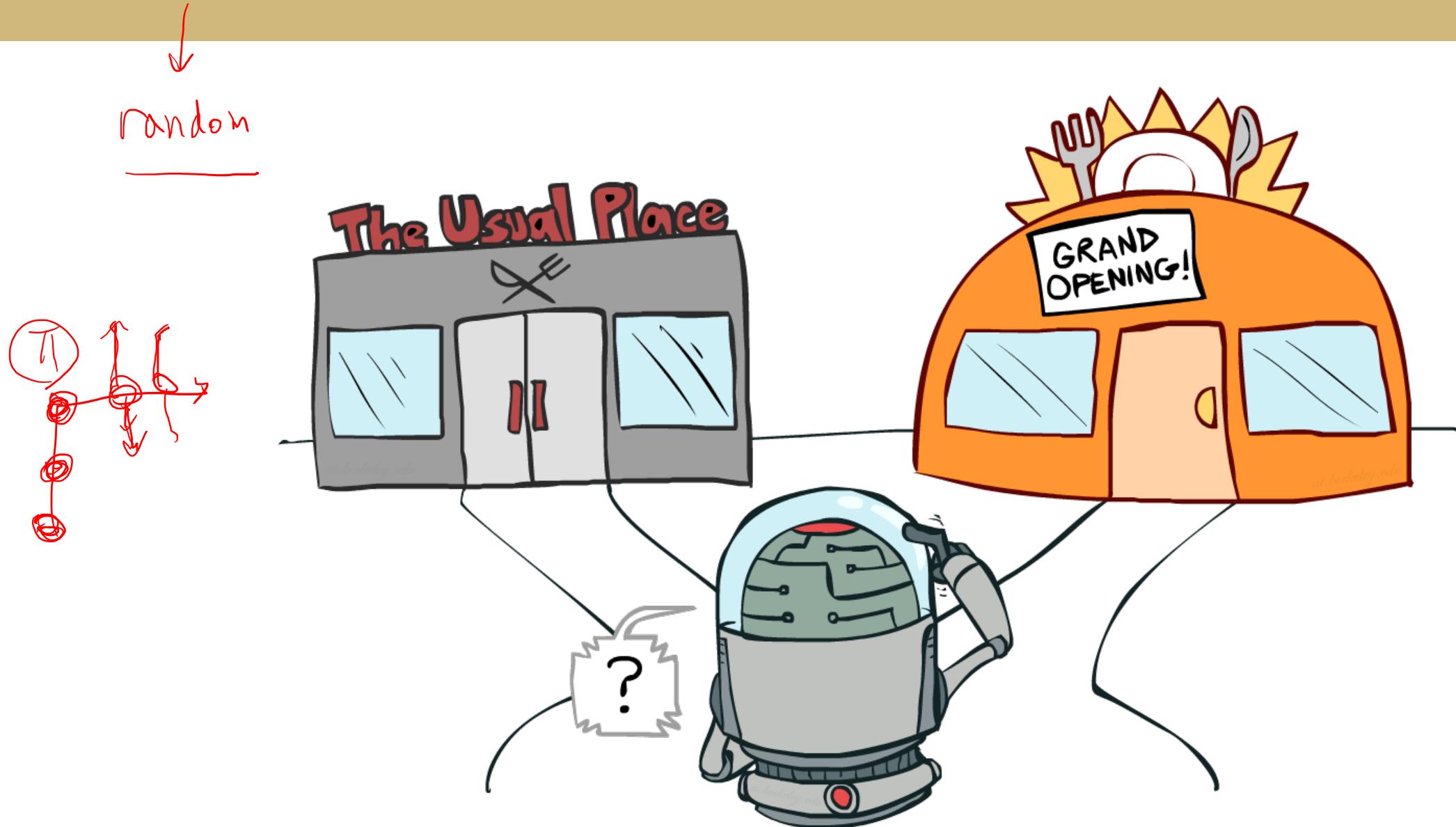
(with ties broken arbitrarily)

For all $a \in \mathcal{A}(S_t)$:

$$\pi(a|S_t) \leftarrow \begin{cases} 1 - \varepsilon + \varepsilon / |\mathcal{A}(S_t)| & \text{if } a = A^* \\ \varepsilon / |\mathcal{A}(S_t)| & \text{if } a \neq A^* \end{cases}$$



Exploration vs. Exploitation



Note on Exploration vs. Exploitation

- ϵ -greedy policy allows random exploration
- General method: can be used with MC or TD (Q-learning)
- Alters the policy: called on-policy method

$$\pi \begin{cases} \epsilon \\ 1-\epsilon \end{cases}$$

Problems with random actions?

- You do eventually explore the space, but keep thrashing around once learning is done
- One solution: lower ϵ over time
- Another solution: exploration functions

$$\epsilon \rightarrow \downarrow$$



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



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

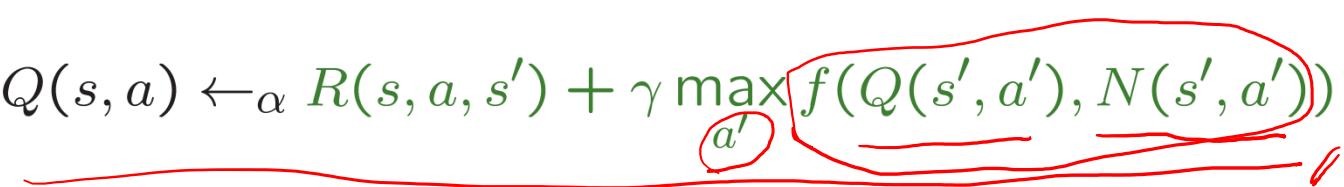
$$\text{Q-state } (S-a) \quad \text{incentive}$$

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



Video of Demo Q-learning – Exploration Function – Crawler

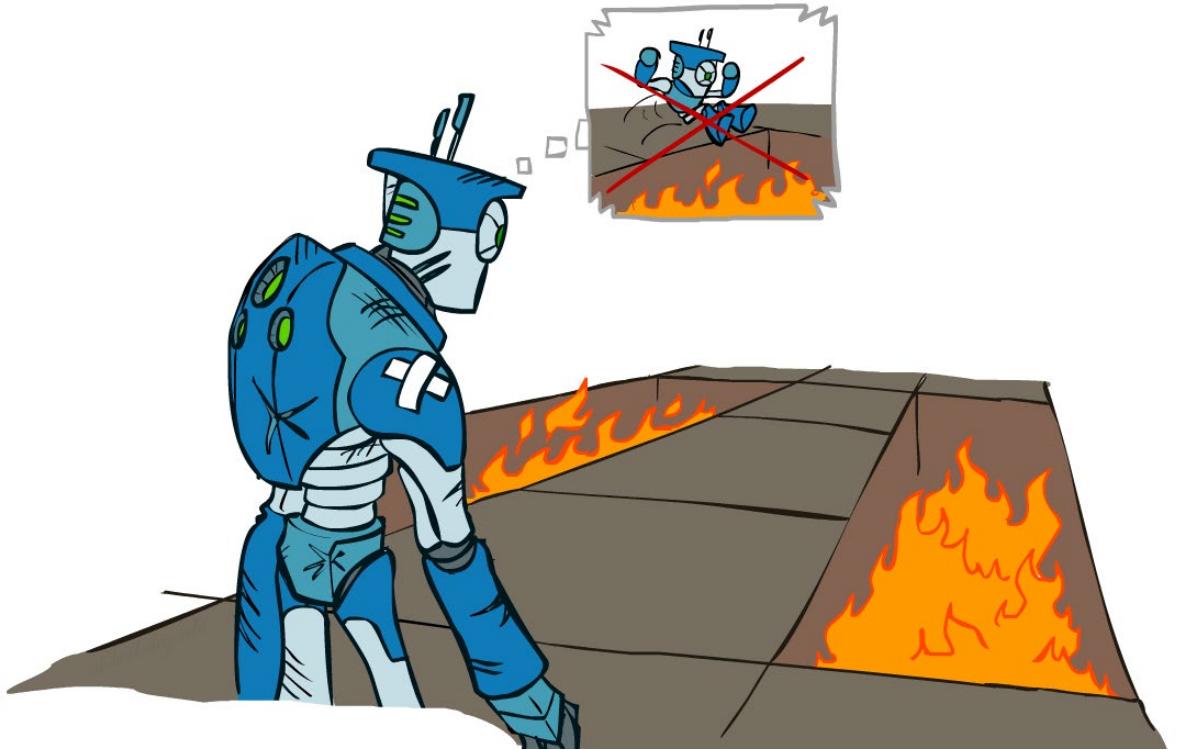


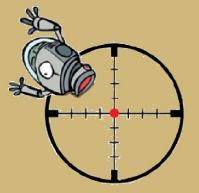
Regret

- 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, including youthful suboptimality, 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

ϵ -greedy $\epsilon \rightarrow \text{Random}$
 $1 - \epsilon$

Exploration Function





Off-policy method

- On-policy method compromises the optimal policy in order to include exploration.
- An ϵ -greedy policy is suboptimal after convergence

one π ↘
 $1-\epsilon$

- Solution: We'll have two policies. One we'll keep updating to be optimal, and the other for exploration.

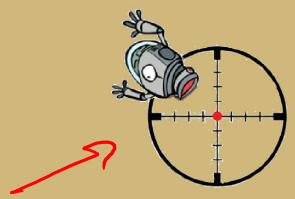
1

2

- Target policy ; π to be improved $\rightarrow \pi^*$ to optimize

- Behavior policy : $b \models \epsilon$ -greedy to explore

- This is called off-policy method



Off-policy MC control

Off-policy MC control, for estimating $\pi \approx \pi_*$

Initialize, for all $s \in \mathcal{S}$, $a \in \mathcal{A}(s)$:

$$Q(s, a) \in \mathbb{R} \text{ (arbitrarily)}$$

$$C(s, a) \leftarrow 0$$

$$\pi(s) \leftarrow \arg \max_a Q(s, a) \quad (\text{with ties broken consistently})$$

Loop forever (for each episode):

$b \leftarrow$ any soft policy (ϵ -greedy)

Generate an episode using b : $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T$

$$G \leftarrow 0$$

$$W \leftarrow 1$$

Loop for each step of episode, $t = T-1, T-2, \dots, 0$:

$$G \leftarrow \gamma G + R_{t+1}$$

$$C(S_t, A_t) \leftarrow C(S_t, A_t) + W$$

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \frac{W}{C(S_t, A_t)} [G - Q(S_t, A_t)]$$

$$\pi(S_t) \leftarrow \arg \max_a Q(S_t, a) \quad (\text{with ties broken consistently})$$

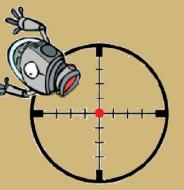
If $A_t \neq \pi(S_t)$ then exit inner Loop (proceed to next episode)

$$W \leftarrow W \frac{1}{b(A_t | S_t)}$$

Handwritten notes on the left side:

- $Q^{\text{new}} = (1-\lambda)Q^{\text{old}} + \lambda \text{ sample}$ (with ties broken consistently)
- $T \Delta \text{ learning}$
- $R \times \frac{1}{C}$
- $\pi(s) \leftarrow \arg \max_a Q(s, a)$ (with ties broken consistently)
- $\pi(s) \leftarrow \arg \max_a Q(s, a)$ (with ties broken consistently)
- $\pi(s) \leftarrow \arg \max_a Q(s, a)$ (with ties broken consistently)

$$W = \frac{\pi(a_t | s_t)}{b(a_t | s_t)}$$



On-policy vs. Off-policy TD control

Sarsa (on-policy TD control) for estimating $Q \approx q_*$

Algorithm parameters: step size $\alpha \in (0, 1]$, small $\varepsilon > 0$

Initialize $Q(s, a)$, for all $s \in \mathcal{S}^+$, $a \in \mathcal{A}(s)$, arbitrarily except that $Q(\text{terminal}, \cdot) = 0$

Loop for each episode:

 Initialize S

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

 Loop for each step of episode:

 Take action A , observe R, S'

 Choose A' from S' using policy derived from Q (e.g., ε -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

TD

π ε random
 $(1-\varepsilon)$ best policy

Q-learning (off-policy TD control) for estimating $\pi \approx \pi_*$

Algorithm parameters: step size $\alpha \in (0, 1]$, small $\varepsilon > 0$

Initialize $Q(s, a)$, for all $s \in \mathcal{S}^+$, $a \in \mathcal{A}(s)$, arbitrarily except that $Q(\text{terminal}, \cdot) = 0$

Loop for each episode:

 Initialize S

 Loop for each step of episode:

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

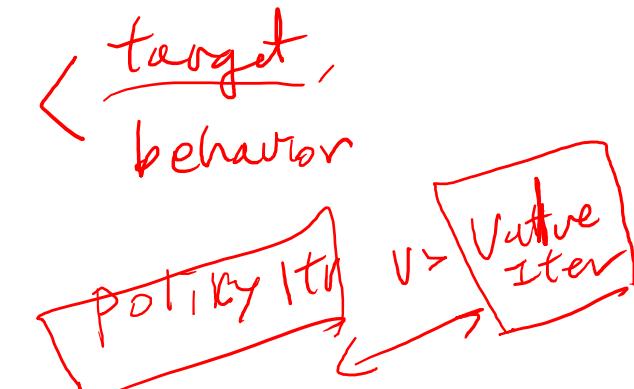
 Take action A , observe R, S'

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

$S \leftarrow S'$

 until S is terminal

Target π



The Story So Far: MDPs and RL

Known MDP: Offline Solution

Goal

Compute V^*, Q^*, π^*

V, Q
 τ^l

Technique

Value / policy iteration

Evaluate a fixed policy π

Policy evaluation

Unknown MDP: Model-Based

Goal

Compute V^*, Q^*, π^*

Technique

P, R
VI/PI on approx. MDP

Evaluate a fixed policy π

PE on approx. MDP

Unknown MDP: Model-Free

Goal

Compute V^*, Q^*, π^*

Technique

Q-learning

P, R

$MC \rightarrow$

$TD \rightarrow$

Evaluate a fixed policy π

Value Learning

$Q\text{-learny}$

Recap- concepts in RL

→ Exact Solution (S, A)

