

Artificial Intelligence Foundations and Applications

Partha Pratim Chakrabarti

Plaban Bhowmik

Sudeshna Sarkar

Centre of Excellence in Artificial Intelligence

IIT Kharagpur

Planning

Stochastic Planning

Reinforcement Learning

Many slides adapted from

CS 188:University of California, Berkeley by Pieter Abbeel

CS221 : Stanford University by Percy Liang

Markov decision processes

states S

Start state s_0

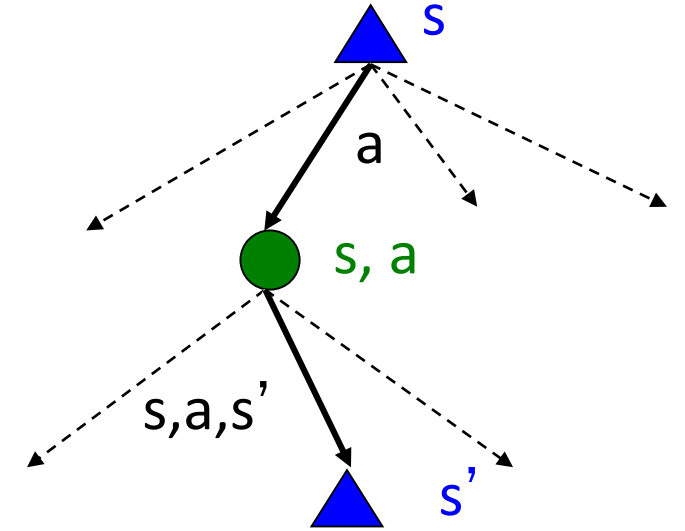
actions A

Transitions $P(s' | s, a)$ or $T(s, a, s')$

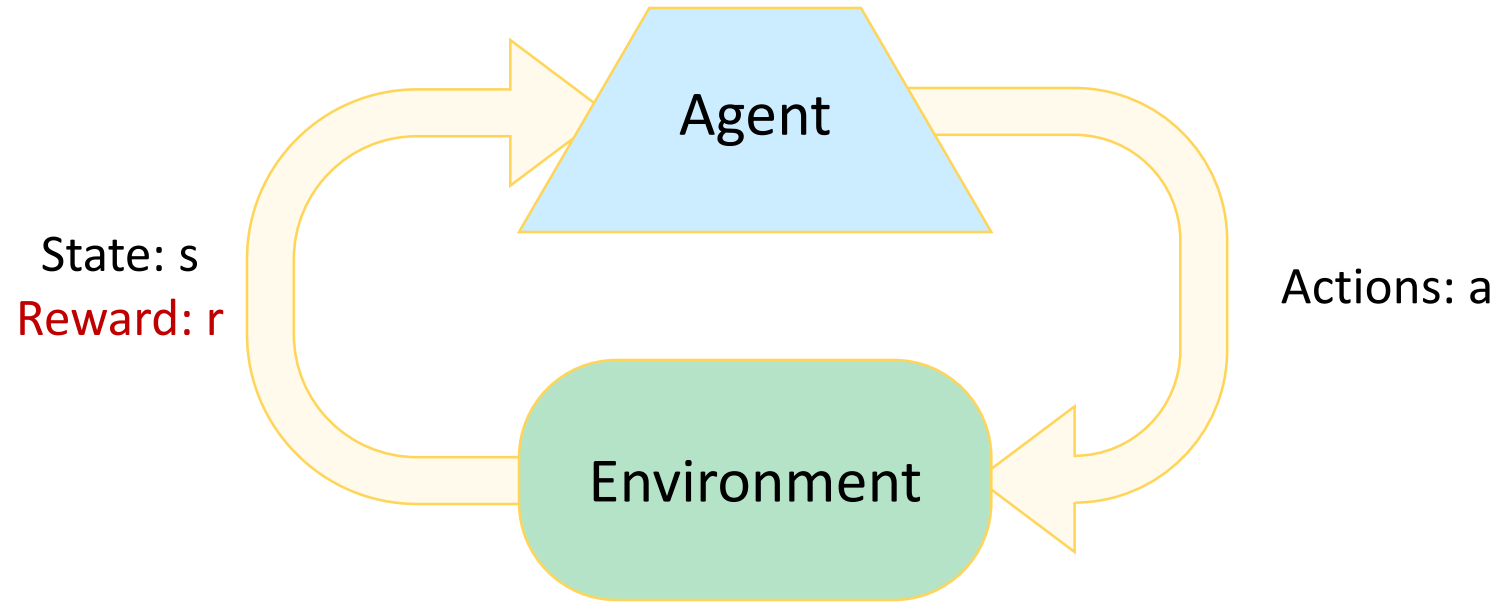
Rewards $R(s, a, s')$ (and discount γ)

MDP quantities :

- **Policy** = Choice of action for each state
- **Utility** = sum of (discounted) rewards
- **Values** = expected future utility from a state
- **Q-Values** = expected future utility from a q-state



Reinforcement Learning

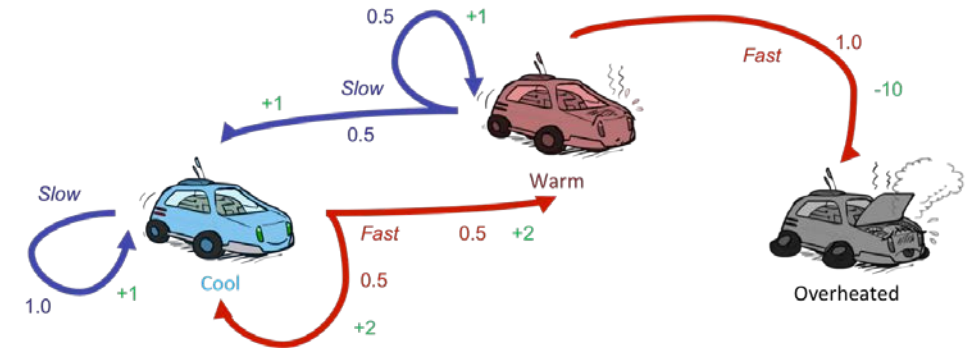


Basic idea:

- 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!

Reinforcement Learning

- Assume a Markov decision process (MDP):
- Looking for a policy $\pi(s)$
- BUT: don't know T or R



From MDPs to reinforcement learning

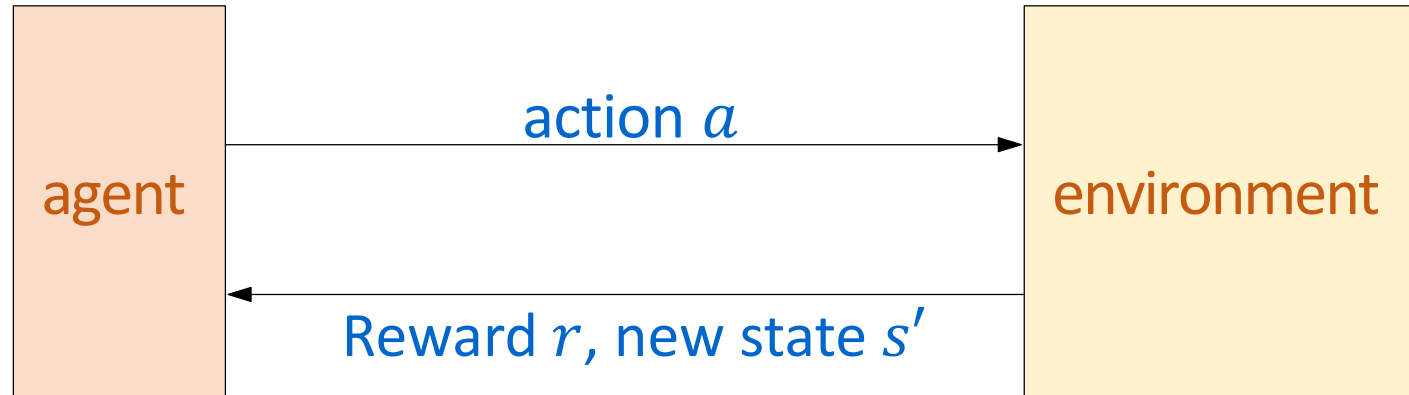
Markov decision process (offline)

- Have mental model of how the world works.
- Find policy to collect maximum rewards.

Reinforcement learning (online)

- Don't know how the world works.
- Perform actions in the world to find out and collect rewards.

Reinforcement learning framework



Algorithm: reinforcement learning template

For $t = 1, 2, 3, \dots$

Choose action $a_t = \pi_{act}(s_{t-1})$ (how?)

Receive reward r_t and observe new state s_t

Update parameters (how?)

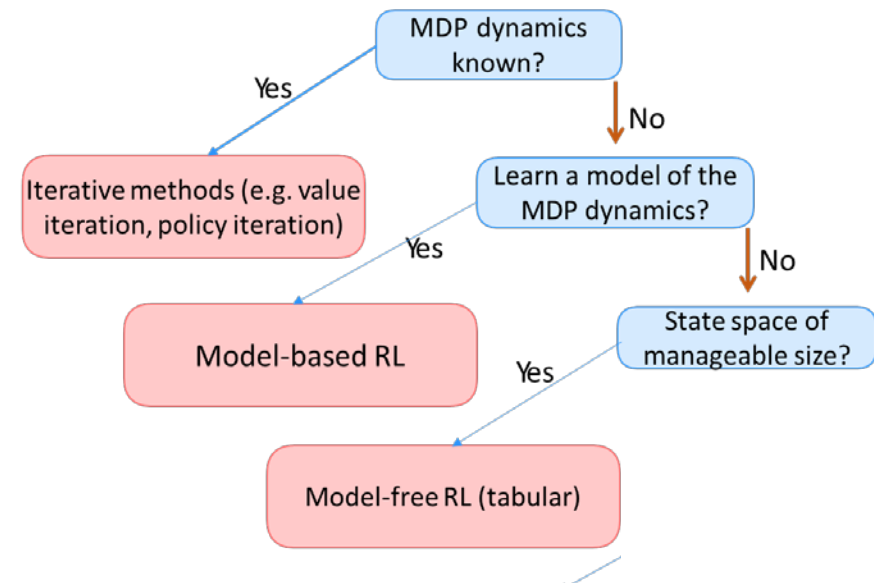
Reinforcement Learning -- Overview

Passive Reinforcement Learning (= how to learn from experiences)

- Model-based Passive RL
 - Learn the MDP model from experiences, then solve the MDP
- Model-free Passive RL
 - Forego learning the MDP model, directly learn V or Q:
Q learning – learns Q values of the optimal policy (uses a Q version of TD Learning)

Active Reinforcement Learning (= agent also needs to decide how to collect experiences)

- Key challenges:
 - How to efficiently explore?
 - How to trade off exploration <> exploitation



Reinforcement Learning -- Overview

- **Passive Reinforcement Learning (= how to learn from experiences)**
 - **Model-based Passive RL**
 - Learn the MDP model from experiences, then solve the MDP
 - **Model-free Passive RL**
 - Forego learning the MDP model, directly learn V or Q:
 - Value learning – learns value of a fixed policy; 2 approaches: Direct Evaluation & TD Learning
 - Q learning – learns Q values of the optimal policy (uses a Q version of TD Learning)
- **Active Reinforcement Learning (= agent also needs to decide how to collect experiences)**
 - Key challenges:
 - How to efficiently explore?
 - How to trade off exploration <> exploitation
 - Applies to both model-based and model-free.

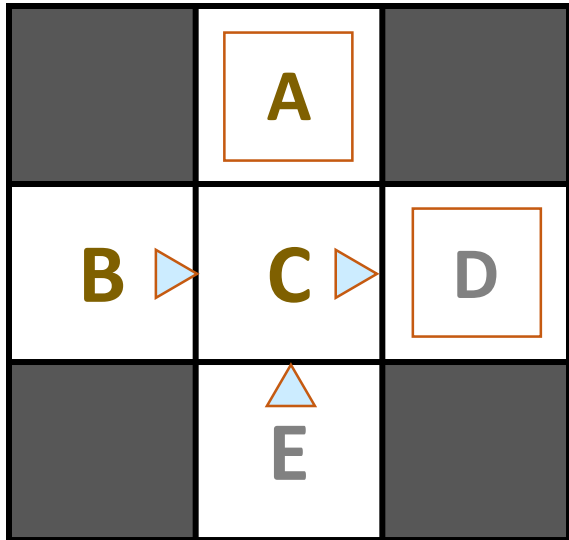
Model-Based Reinforcement Learning

- Model-Based Idea:
 - Learn an approximate model based on experiences
 - Solve for values as if the learned model were correct
- Step 1: Learn empirical MDP model
 - Count outcomes s' for each s, a
 - Normalize to give an estimate of $\hat{T}(s, a, s')$
 - Discover each $\hat{R}(s, a, s')$ when we experience (s, a, s')
- Step 2: Solve the learned MDP
 - For example, use value iteration, as before

(and repeat as needed)

Example: Model-Based RL

Input Policy π



Assume: $\gamma = 1$

Observed Episodes (Training)

Episode 1

B, east, C, -1
C, east, D, -1
D, exit, x, +10

Episode 2

B, east, C, -1
C, east, D, -1
D, exit, x, +10

Episode 3

E, north, C, -1
C, east, D, -1
D, exit, x, +10

Episode 4

E, north, C, -1
C, east, A, -1
A, exit, x, -10

Learned Model

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

T(B, east, C) = 1.00
T(C, east, D) = 0.75
T(C, east, A) = 0.25
...

$$\hat{R}(s, a, s')$$

R(B, east, C) = -1
R(C, east, D) = -1
R(D, exit, x) = +10
...

Analogy: Expected Age

Goal: Compute expected age of the students in the class

Known $P(A)$

$$E[A] = \sum_a P(a) \cdot a = 0.35 \times 20 + \dots$$

Without $P(A)$, instead collect samples $[a_1, a_2, \dots, a_N]$

Unknown $P(A)$: “Model Based”

$$\hat{P}(a) = \frac{\text{num}(a)}{N}$$
$$E[A] \approx \sum_a \hat{P}(a) \cdot a$$

Why does this work? Because eventually you learn the right model.

Unknown $P(A)$: “Model Free”

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

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

Reinforcement Learning -- Overview

- Passive Reinforcement Learning (= how to learn from experiences)
 - Model-based Passive RL
 - Learn the MDP model from experiences, then solve the MDP
 - **Model-free Passive RL**
 - **Forego learning the MDP model, directly learn V or Q:**
 - **Value learning – learns value of a fixed policy; 2 approaches: Direct Evaluation & TD Learning**
 - Q learning – learns Q values of the optimal policy (uses a Q version of TD Learning)
- Active Reinforcement Learning (= agent also needs to decide how to collect experiences)
 - Key challenges:
 - How to efficiently explore?
 - How to trade off exploration <> exploitation
 - Applies to both model-based and model-free.

Passive Model-Free Reinforcement Learning

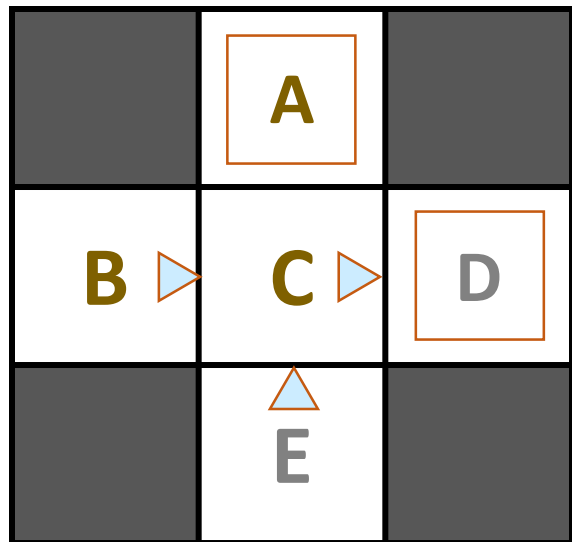
- Simplified task: policy evaluation
 - Input: a fixed policy $\pi(s)$
 - You don't know the transitions $T(s,a,s')$
 - You don't know the rewards $R(s,a,s')$
 - Goal: learn the state values
- In this case:
 - Learner is “along for the ride”
 - No choice about what actions to take
 - Just execute the policy and learn from experience
 - This is NOT offline planning! You actually take actions in the world.

Direct Evaluation

- Goal: Compute values for each state under π
- Idea: Average together observed sample values
 - Act according to π
 - Every time you visit a state, write down what the sum of discounted rewards turned out to be
 - Average those samples
- This is called direct evaluation

Example: Direct Evaluation

Input Policy π



Assume: $\gamma = 1$

Observed Episodes (Training)

Episode 1

B, east, C, -1
C, east, D, -1
D, exit, x, +10

Episode 2

B, east, C, -1
C, east, D, -1
D, exit, x, +10

Episode 3

E, north, C, -1
C, east, D, -1
D, exit, x, +10

Episode 4

E, north, C, -1
C, east, A, -1
A, exit, x, -10

Output Values

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

If B and E both go to C under this policy, how can their values be different?

Problems with Direct Evaluation

- What's good about direct evaluation?
 - It's easy to understand
 - It doesn't require any knowledge of T , R
 - It eventually computes the correct average values, using just sample transitions
- What bad about it?
 - It wastes information about state connections
 - Each state must be learned separately
 - So, it takes a long time to learn

Output Values

	<div>-10 A</div>	
<div>+8 B</div>	<div>+4 C</div>	<div>+10 D</div>
	<div>-2 E</div>	

If B and E both go to C under this policy, how can their values be different?

Reinforcement Learning -- Overview

- Passive Reinforcement Learning (= how to learn from experiences)
 - Model-based Passive RL
 - Learn the MDP model from experiences, then solve the MDP
 - **Model-free Passive RL**
 - **Forego learning the MDP model, directly learn V or Q:**
 - **Value learning – learns value of a fixed policy; 2 approaches: Direct Evaluation & TD Learning**
 - Q learning – learns Q values of the optimal policy (uses a Q version of TD Learning)
- Active Reinforcement Learning (= agent also needs to decide how to collect experiences)
 - Key challenges:
 - How to efficiently explore?
 - How to trade off exploration <> exploitation

Temporal Difference Value Learning

Sample-Based Policy Evaluation?

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

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

$V_k(s)$

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

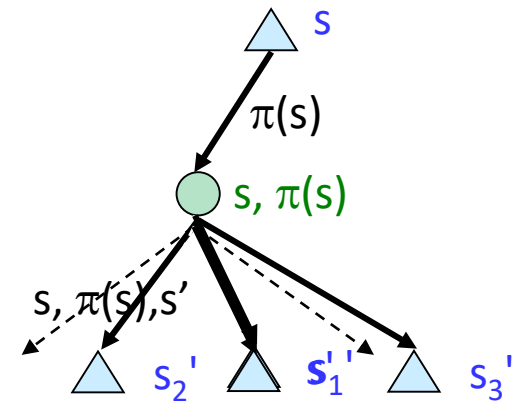
$$\text{sample}_1 = R(s, \pi(s), s'_1) + \gamma V_k^\pi(s'_1)$$

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

...

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

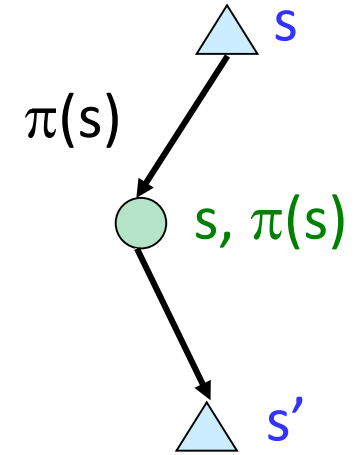
$$\underline{V_{k+1}^\pi(s)} \leftarrow \frac{1}{n} \sum_i \text{sample}_i$$



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

Temporal Difference Value Learning

- Idea: Learn from every experience!
 - Update $V(s)$ each time we experience a transition (s, a, s', r)
 - Likely outcomes s' will contribute updates more often
- Temporal difference learning of values
 - Policy still fixed, still doing evaluation!
 - Move values toward value of whatever successor occurs: running average



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

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

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

Exponential Moving Average

- Exponential moving average
 - The running interpolation update: $\bar{x}_n = (1 - \alpha) \cdot \bar{x}_{n-1} + \alpha \cdot x_n$
 - Makes recent samples more important
 - Forgets about the past (distant past values were wrong anyway)
- Decreasing learning rate (alpha) can give converging averages

Example: Temporal Difference Value Learning

States

	A	
B	C	D
	E	

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

Observed Transitions

B, east, C, -2

	0	
0	0	8
	0	

C, east, D, -2

	0	
-1	0	8
	0	

	0	
-1	3	8
	0	

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

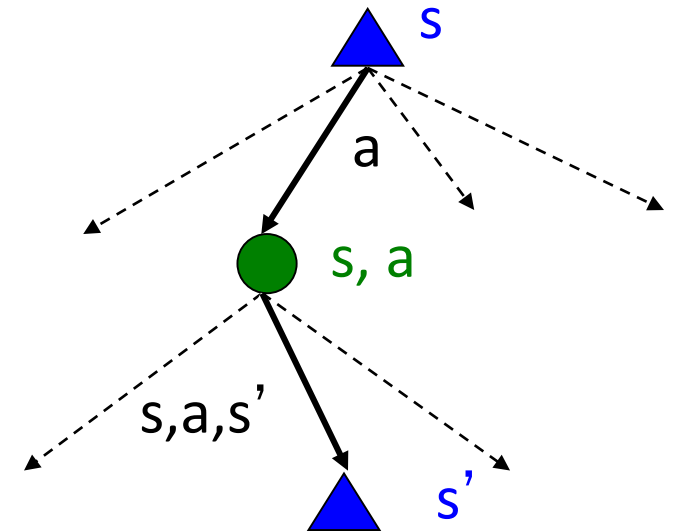
Problems with TD Value Learning

- TD value learning is a model-free way to do policy evaluation, mimicking Bellman updates with running sample averages
- However, if we want to turn values into a (new) policy, we're sunk:

$$\pi(s) = \arg \max_a Q(s, a)$$

$$Q(s, a) = \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V(s')]$$

- Idea: learn Q-values, not values
- Makes action selection model-free too!



Reinforcement Learning -- Overview

- Passive Reinforcement Learning (= how to learn from experiences)
 - Model-based Passive RL
 - Learn the MDP model from experiences, then solve the MDP
 - **Model-free Passive RL**
 - Forego learning the MDP model, directly learn V or Q:
 - Value learning – learns value of a fixed policy; 2 approaches: Direct Evaluation & TD Learning
 - **Q learning – learns Q values of the optimal policy (uses a Q version of TD Learning)**
- Active Reinforcement Learning (= agent also needs to decide how to collect experiences)
 - Key challenges:
 - How to efficiently explore?
 - How to trade off exploration <> exploitation

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:

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

- 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:

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

Q-Learning

- Q-value iteration

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

- Q-Learning: learn $Q(s,a)$ values as you go

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

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

- Incorporate the new estimate into a running average:

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

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:
 $sample = R(s, a, s') + \gamma \max_{a'} Q_k(s', a')$
 - Incorporate the new estimate into a running average:
 $Q(s, a) \leftarrow (1 - \alpha)Q(s, a) + (\alpha)[sample]$



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 (!)

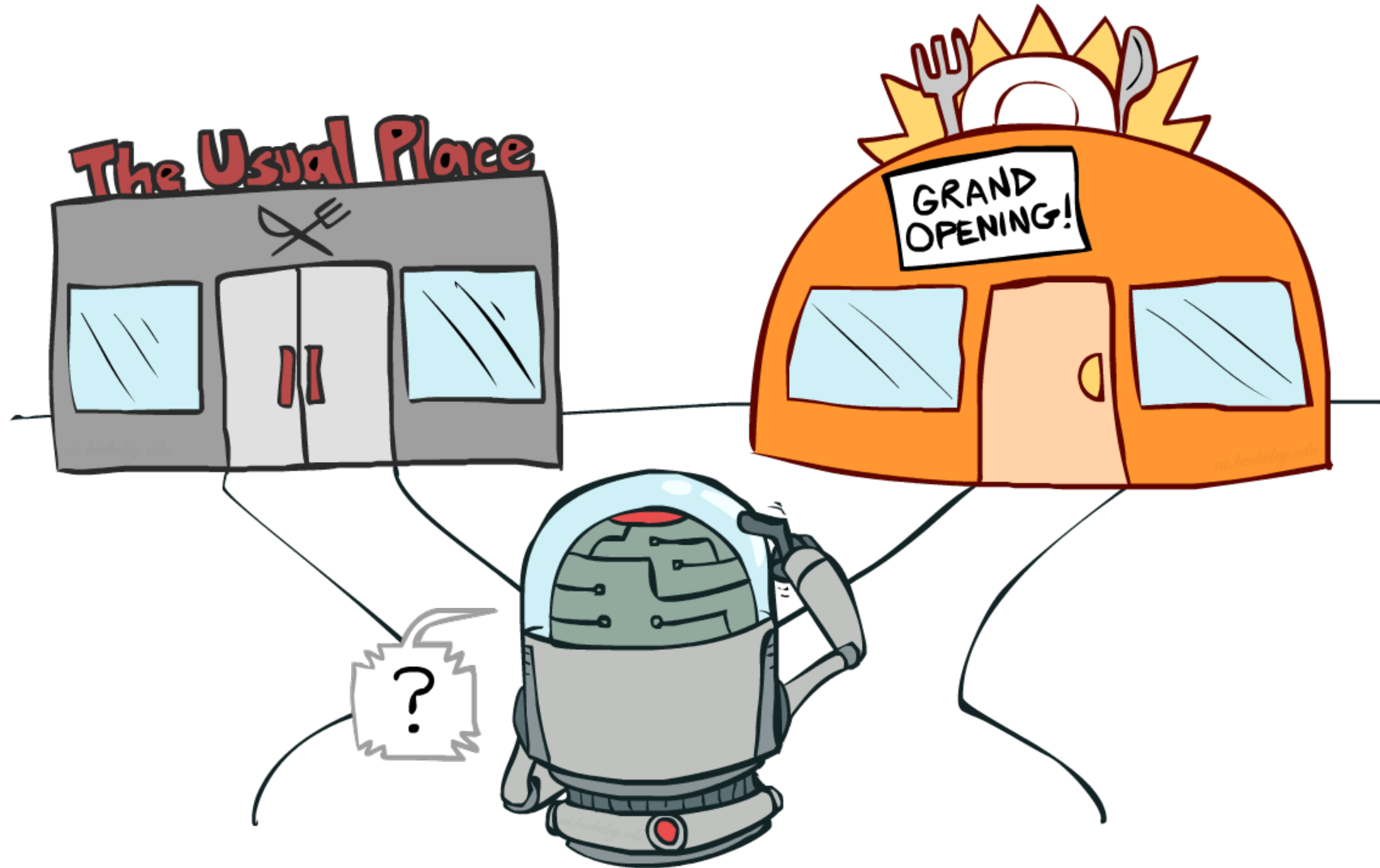
Reinforcement Learning -- Overview

- Passive Reinforcement Learning (= how to learn from experiences)
 - Model-based Passive RL
 - Learn the MDP model from experiences, then solve the MDP
 - Model-free Passive RL
 - Forego learning the MDP model, directly learn V or Q:
 - Value learning – learns value of a fixed policy; 2 approaches: Direct Evaluation & TD Learning
 - Q learning – learns Q values of the optimal policy (uses a Q version of TD Learning)
- **Active Reinforcement Learning (= agent also needs to decide how to collect experiences)**
 - Key challenges:
 - How to efficiently explore?
 - How to trade off exploration <> exploitation

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...

Exploration vs. Exploitation



How to Explore?

- Several schemes for forcing exploration
 - Simplest: random actions (ϵ -greedy)
 - Every time step, flip a coin
 - With (small) probability ϵ , act randomly
 - With (large) probability $1-\epsilon$, act on current policy
 - 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

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.
$$f(u, n) = u + k/n$$
- Regular Q-Update: $Q(s, a) \leftarrow_{\alpha} R(s, a, s') + \gamma \max_{a'} Q(s', a')$
- Modified Q-Update: $Q(s, a) \leftarrow_{\alpha} R(s, a, s') + \gamma \max_{a'} f(Q(s', a'), N(s', a'))$
- Note: this propagates the “bonus” back to states that lead to unknown states as well!

Generalizing Across States

- Basic Q-Learning keeps a table of all q-values
- In realistic situations, we cannot possibly learn about every single state!
 - Too many states to visit them all in training
 - Too many states to hold the q-tables in memory
- Instead, we want to generalize:
 - Learn about some small number of training states from experience
 - Generalize that experience to new, similar situations
 - This is a fundamental idea in machine learning, and we'll see it over and over again

Feature-Based Representations

- Solution: describe a state using a vector of features (properties)
- Features are functions from states to real numbers (often 0/1) that capture important properties of the state
- Example features:
 - Distance to closest ghost
 - Distance to closest dot
 - Number of ghosts
 - $1 / (\text{dist to dot})^2$
 - Is Pacman in a tunnel? (0/1)
 - etc.
 - Is it the exact state on this slide?
- Can also describe a q-state (s, a) with features (e.g. action moves closer to food)

Linear Value Functions

- Using a feature representation, we can write a q function (or value function) for any state using a few weights:

$$V(s) = w_1 f_1(s) + w_2 f_2(s) + \dots + w_n f_n(s)$$

$$Q(s, a) = w_1 f_1(s, a) + w_2 f_2(s, a) + \dots + w_n f_n(s, a)$$

- Advantage: our experience is summed up in a few powerful numbers
- Disadvantage: states may share features but actually be very different in value!

Approximate Q-Learning

$$Q(s, a) = w_1 f_1(s, a) + w_2 f_2(s, a) + \dots + w_n f_n(s, a)$$

- Q-learning with linear Q-functions:

transition = (s, a, r, s')

difference = $\left[r + \gamma \max_{a'} Q(s', a') \right] - Q(s, a)$

$Q(s, a) \leftarrow Q(s, a) + \alpha [\text{difference}]$ Exact Q's

$w_i \leftarrow w_i + \alpha [\text{difference}] f_i(s, a)$ Approximate Q's

- Intuitive interpretation:

- Adjust weights of active features
- E.g., if something unexpectedly bad happens, blame the features that were on: disprefer all states with that state's features

Reinforcement Learning -- Overview

- Passive Reinforcement Learning (= how to learn from experiences)
 - Model-based Passive RL
 - Learn the MDP model from experiences, then solve the MDP
 - Model-free Passive RL
 - Forego learning the MDP model, directly learn V or Q:
 - Value learning – learns value of a fixed policy; 2 approaches: Direct Evaluation & TD Learning
 - Q learning – learns Q values of the optimal policy (uses a Q version of TD Learning)
- Active Reinforcement Learning (= agent also needs to decide how to collect experiences)
 - Key challenges:
 - How to efficiently explore?
 - How to trade off exploration <> exploitation
 - Applies to both model-based and model-free.

Deep Reinforcement Learning

- Deep Reinforcement Learning = Deep Learning + Reinforcement Learning
- Value-based RL
 - Estimate optimal value function $Q_*(s, a)$
 - Find maximum value achievable under policy
- Policy-based RL
 - Search directly for optimal policy
 - Policy for achieving maximum future reward
- Model-based RL
 - Build an environment model and plan by using look-ahead

Deep Q-learning

- Represent value function using a Q -network with weights w :
- Look for $Q(s, a, w) \approx Q_*(s, a)$

$$Q^*(s, a) = \mathbb{E} \left[r + \gamma \max_{a'} Q(s', a') \mid s, a \right]$$

- Instead treat RHS $(r + \gamma \max_{a'} Q(s', a', w))$ as a target, and
- Minimize MSE loss using stochastic gradient descent

$$I = \left(r + \gamma \max_{a'} Q(s', a', w) - Q(s, a, w) \right)^2$$

- Intuition: For optimal Q^* , above term will be zero, neural network will approximate the Q function

