# CS 471/571 (Fall 2023): Introduction to Artificial Intelligence

## Lecture 14: Reinforcement Learning (Part 3)
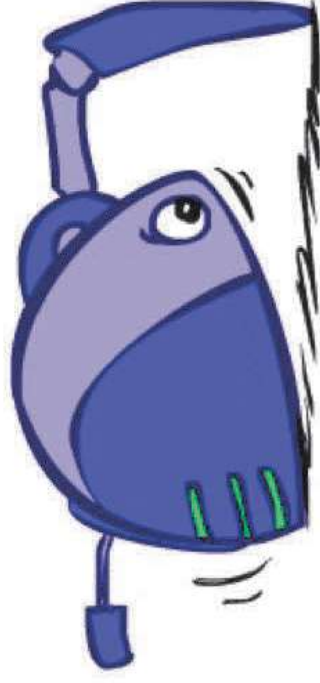
Thanh H. Nguyen

# Reminder

- Written assignment 3: MDPs and Reinforcement Learning
  - Deadline: Nov 08th, 2023

# Reinforcement Learning

- We still assume an MDP:
  - A set of states s ∈ S
  - A set of actions (per state) A
  - A model T(s,a,s')
  - A reward function R(s,a,s')
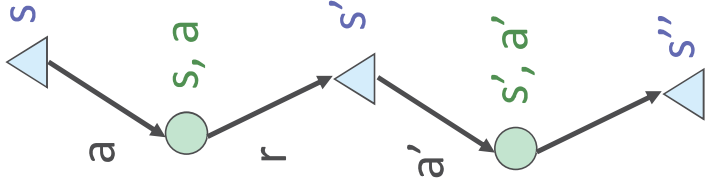- Still looking for a policy π(s)

- New twist: don't know T or R, so must try out actions

- Big idea: Compute all averages over T using sample outcomes

# Model-Free Learning

- Model-free (temporal difference) learning
  - Experience world through episodes

$$(s, a, r, s', a', r', s'', a'', r'', s''' \ldots)$$

  - Update estimates each transition $(s, a, r, s')$

  - Over time, updates will mimic Bellman updates

# Q-Learning

- We'd like to do Q-value updates to each Q-state:

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

  - But can't compute this update without knowing T, R

- Instead, compute average as we go
  - Receive a sample transition (s,a,r,s')
  - This sample suggests

$$Q(s, a) \approx r + \gamma \max_{a'} Q(s', a')$$

  - But we want to average over results from (s,a)   (Why?)
  - So keep a running average

$$Q(s, a) \leftarrow (1 - \alpha) Q(s, a) + (\alpha) \left[ r + \gamma \max_{a'} Q(s', a') \right]$$
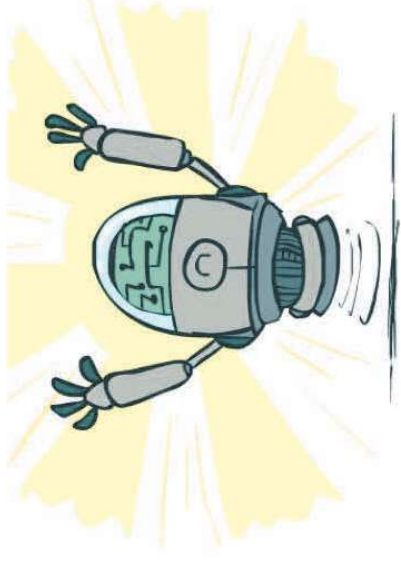
# Example

- Two states: A, B
- Two actions: Up, Down
- Discount factor: $\gamma = 0.5$
- Learning rate: $\alpha = 0.5$

- Q(A, Down) = ?
- Q(B, Up) = ?

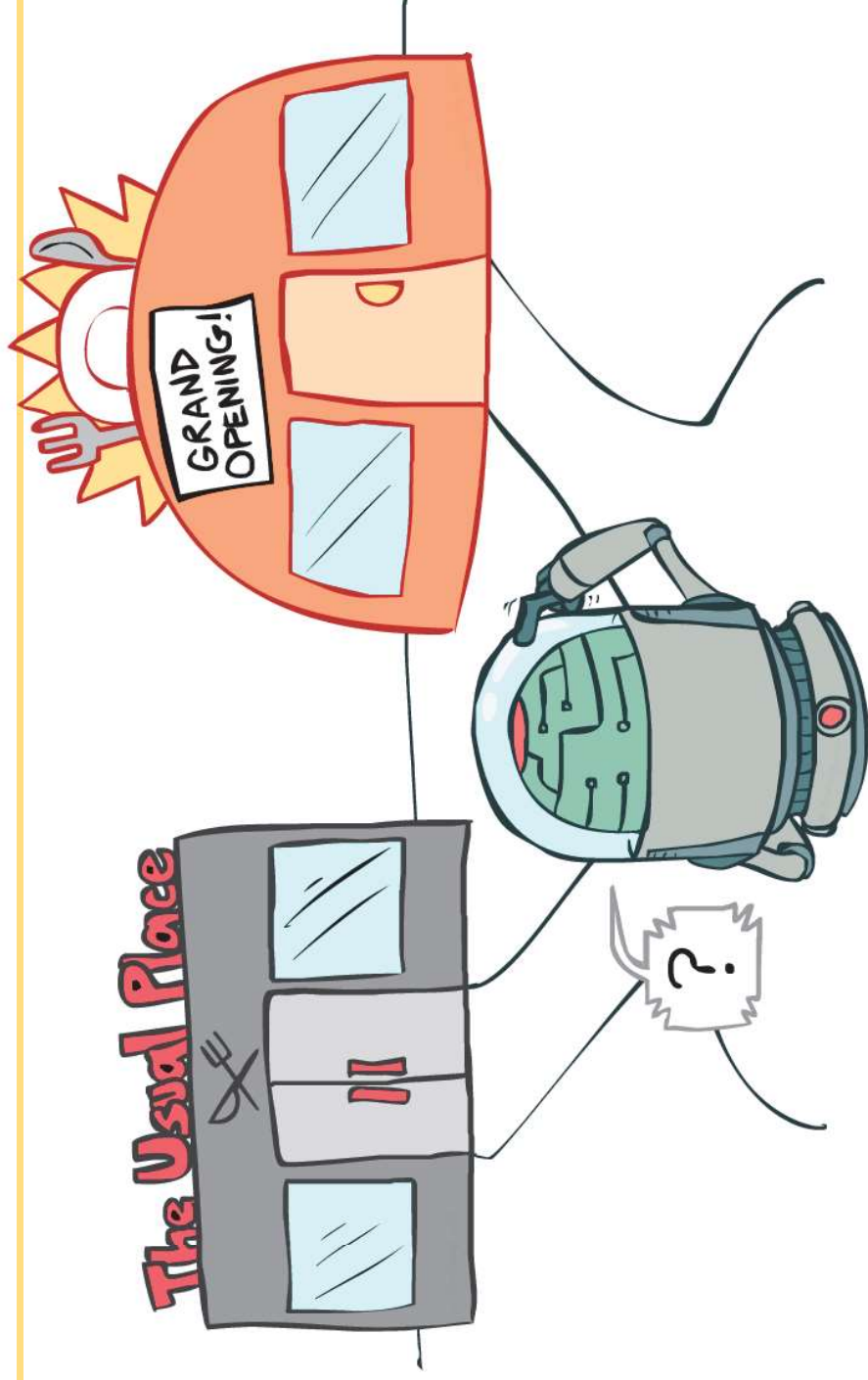| $t$ | $s_t$ | $a_t$ | $s_{t+1}$ | $r_t$ |
|---|---|---|---|---|
| 0 | A | Down | B | 2 |
| 1 | B | Down | B | -4 |
| 2 | B | Up | B | 0 |
| 3 | B | Up | A | 3 |
| 4 | A | Up | A | -1 |

$$Q(s,a) \leftarrow (1-\alpha)Q(s,a) + (\alpha)\left[ r + \gamma \max_{a'} Q(s',a') \right]$$

# Q-Learning Properties

- Amazing result: Q-learning converges to optimal policy -- even if you're acting suboptimally!

- 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 (!)

# Exploration vs. Exploitation

# How to Explore?

- Several schemes for forcing exploration
  - Simplest: random actions ($\varepsilon$-greedy)
    - Every time step, flip a coin
    - With (small) probability $\varepsilon$, act randomly
    - With (large) probability $1-\varepsilon$, 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 $\varepsilon$ 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.
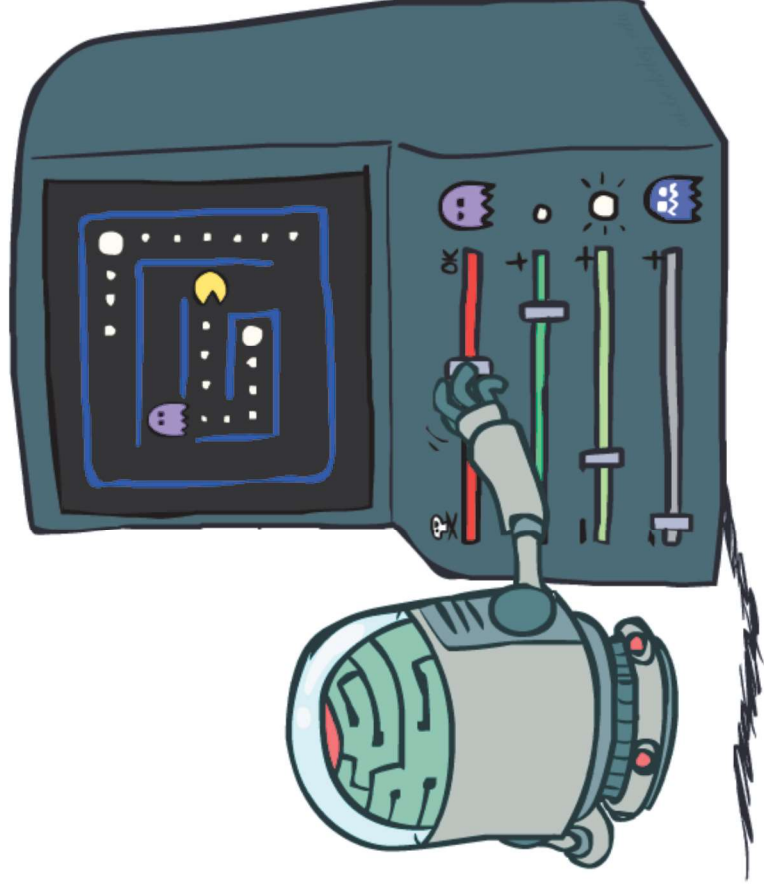
$$f(u, n) = u + k/n$$

Regular Q-Update: $\quad Q(s, a) \leftarrow_\alpha R(s, a, s') + \gamma \max_{a'} Q(s', a')$

Modified Q-Update: $Q(s, a) \leftarrow_\alpha R(s, a, s') + \gamma \max_{a'} f(Q(s', a'), N(s', a'))$

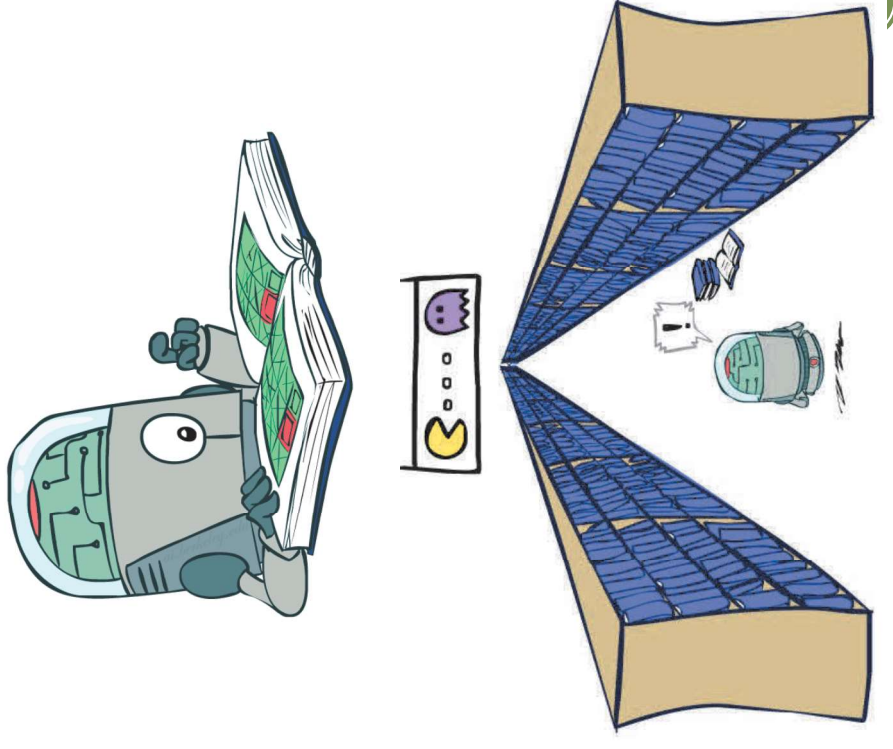- Note: this propagates the "bonus" back to states that lead to unknown states as well!
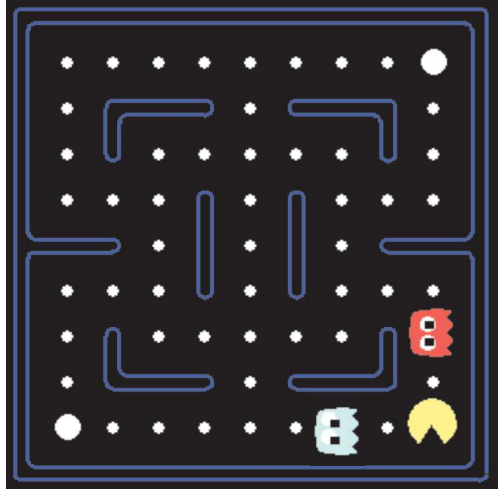
# Approximate Q-Learning

# 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
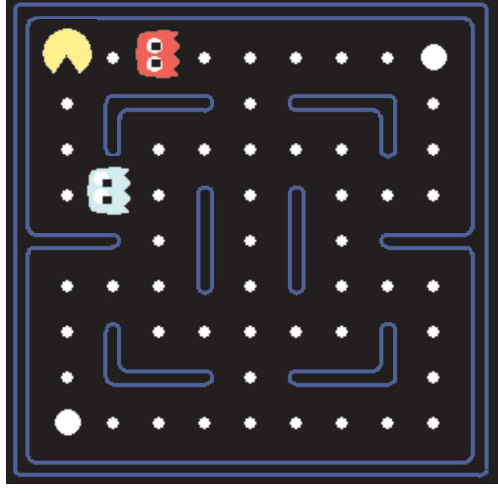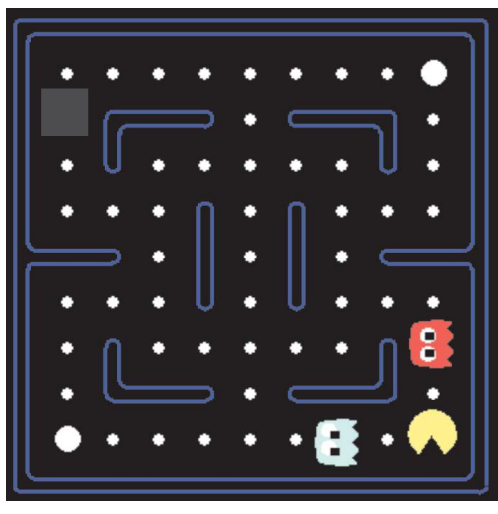
# Example: Pacman

Let's say we discover through experience that this state is bad:

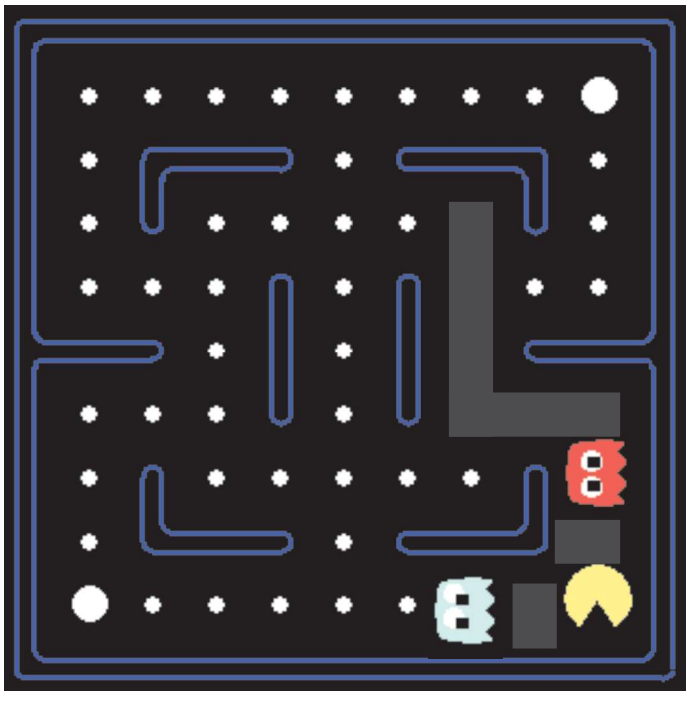In naïve q-learning, we know nothing about this state:

Or even this one!

# 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:

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

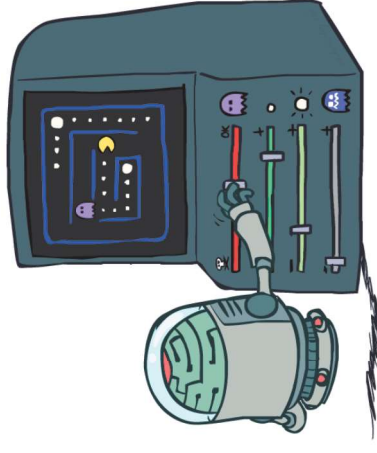$$Q(s,a) = w_1 f_1(s,a) + w_2 f_2(s,a) + \ldots + w_n f_n(s,a)$$

- Advantage: our experience is summed up in a few powerful numbers

- Disadvantage: states may share features but actually be very different in value!
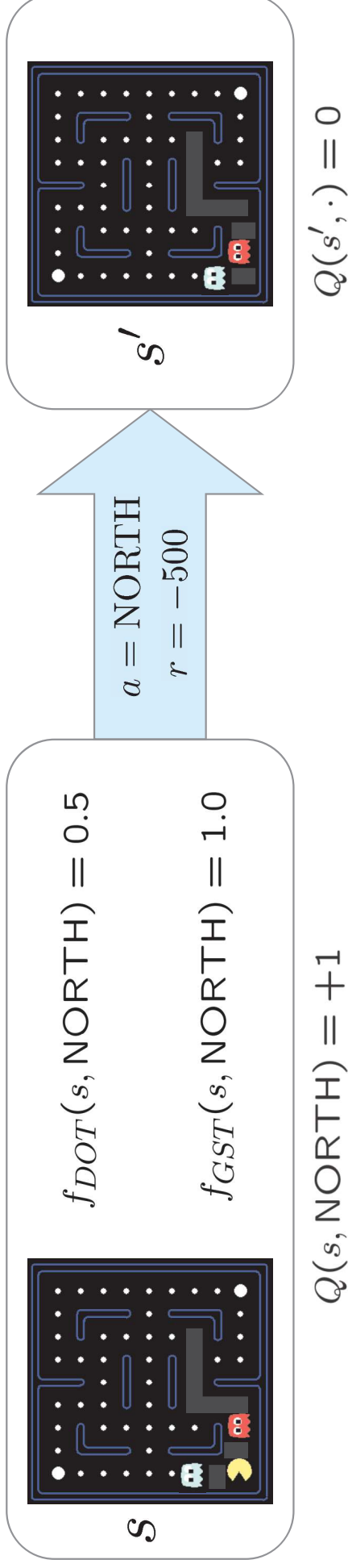
# Approximate Q-Learning

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

- Q-learning with linear Q-functions:

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

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

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

- Formal justification: online least squares

# Example: Q-Pacman

$$Q(s, a) = 4.0 f_{DOT}(s, a) - 1.0 f_{GST}(s, a)$$

$f_{DOT}(s, \text{NORTH}) = 0.5$

$f_{GST}(s, \text{NORTH}) = 1.0$

$a = \text{NORTH}$
$r = -500$

$s'$

$Q(s', \cdot) = 0$

$Q(s, \text{NORTH}) = +1$

$r + \gamma \max_{a'} Q(s', a') = -500 + 0$

difference $= -501$

$w_{DOT} \leftarrow 4.0 + \alpha \, [-501] \, 0.5$

$w_{GST} \leftarrow -1.0 + \alpha \, [-501] \, 1.0$

$$Q(s, a) = 3.0 f_{DOT}(s, a) - 3.0 f_{GST}(s, a)$$

$s$

# Q-learning with Linear Approximation

**Algorithm 4:** Q-learning with linear approximation.

1 Initialize q-value function $Q$ with random weights $w$: $Q(s, a; w) = \sum_m w_m f_m(s, a)$;

2 **for** $episode = 1 \to M$ **do**

3  Get initial state $s_0$;

4  **for** $t = 1 \to T$ **do**

5   With prob. $\epsilon$, select a random action $a_t$;

6   With prob. $1 - \epsilon$, select $a_t \in \mathrm{argmax}_a Q(s_t, a; w)$;

7   Execute selected action $a_t$ and observe reward $r_t$ and next state $s_{t+1}$;

8   Set target $y_t = \begin{cases} r_t & \text{if episode terminates at step } t+1 \\ r_t + \gamma \max_{a'} Q(s_{t+1}, a'; w) & \text{otherwise} \end{cases}$ ;

9   Perform a gradient descent step to update $w$: $w_m \leftarrow w_m + \alpha \left[ y_t - Q(s_t, a_t; w) \right] f_m(s, a)$;

Thanh H. Nguyen

18