

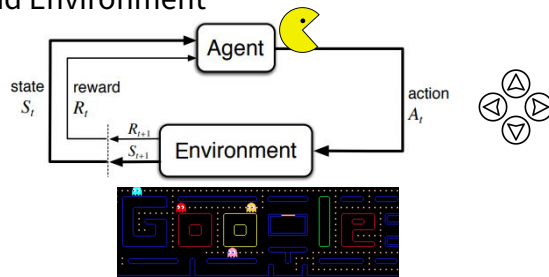
Cornell Bowers CIS

## Logistics

- Will release midterm grades by the end of the week
- Releasing a sign-up sheet for final presentation slots later today
  - First day of presentations is April 30th
- Announcement will include instructions for the final presentation

Cornell Bowers CIS

## Agent and Environment



Agent:

- Perceives environment.
- Makes decisions.
- Aims to maximize reward.

Environment:

- The external context in which an agent operates and interacts with
- Provides feedback to agent

Cornell Bowers CIS

## Markov Decision Process (MDP)

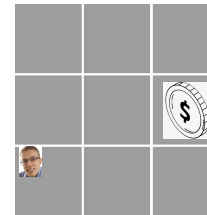
- MDPs provide a framework for modeling sequential decision-making problems.
- An MDP is defined by a tuple  $\langle \mathcal{S}, \mathcal{A}, \mathcal{P}, \mathcal{R}, \gamma \rangle$ :
  - $\mathcal{S}$ : Set of states representing the environment.
  - $\mathcal{A}$ : Set of actions the agent can take.
  - $\mathcal{P}$ : Transition probability function,  $\mathcal{P}(s'|s, a)$ .
  - $\mathcal{R}$ : Reward function,  $\mathcal{R}(s, a)$ .
  - $\gamma$ : Discount factor,  $\gamma \in [0, 1]$ .

## Q-Table

- In Q-Learning, the Q-function is typically represented using a Q-table.
  - A table that stores the estimated Q-values for state-action pairs.
- Q-table has dimensions  $|\mathcal{S}| \times |\mathcal{A}|$ , where:
  - $|\mathcal{S}|$  is the number of states in the state space.
  - $|\mathcal{A}|$  is the number of actions in the action space.
- The Q-table is initialized with arbitrary values and iteratively updated based on the agent's experiences during the learning process.

## Q-Table Example

- 9 states and 4 actions
- Initialize valid (s,a) tuples to 0s



	Up	Down	Right	Left
Bottom Left	0	-	0	-
Bottom Middle	0	-	0	0
Bottom Right	0	-	-	0
Mid Left	0	0	0	-
Mid Middle	0	0	0	0
Mid Right	0	0	-	0
Top Left	-	0	0	-
Top Middle	-	0	0	0
Top Right	-	0	-	0

## Q-Learning Update Rule

$$Q^*(s, a) = \mathcal{R}(s, a) + \gamma \sum_{s' \in \mathcal{S}} \mathcal{P}(s' | s, a) \max_{a' \in \mathcal{A}} Q^*(s', a')$$

- The Q-Learning update rule can be expanded as:

$$Q(s, a) \leftarrow Q(s, a) + \alpha [\mathcal{R}(s, a) + \gamma \max_{a'} Q(s', a') - Q(s, a)]$$

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

- The update rule adjusts the current Q-value estimate  $Q(s, a)$  in the direction of the target Q-value based on the Bellman error.
  - Q-values are gradually improved and converge towards the optimal Q-function.

## Q-Learning Algorithm

Initialize  $Q(s, a), \forall s \in \mathcal{S}, a \in \mathcal{A}(s)$ , arbitrarily, and  $Q(\text{terminal-state}, \cdot) = 0$

Repeat (for each episode):

  Initialize  $S$

  Repeat (for each step of episode):

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

    Take action  $A$ , observe  $R, S'$

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

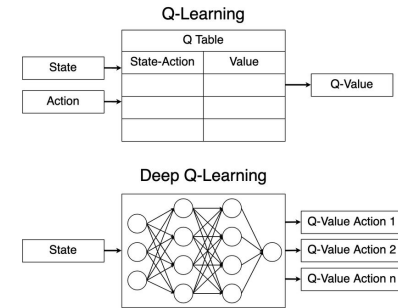
$S \leftarrow S'$ ;

  until  $S$  is terminal

## $\epsilon$ -Greedy Policy

- Simple solution to balance exploration and exploitation:  $\epsilon$ -greedy policy
- $\epsilon$ -greedy policy:
  - With probability  $1 - \epsilon$ , choose the optimal action according to the learned Q-values.
  - With probability  $\epsilon$ , choose a random action.
- The  $\epsilon$ -greedy policy ensures that the agent explores the environment while still exploiting the learned knowledge.
- Despite its simplicity,  $\epsilon$ -greedy is still widely used in practice and often yields good results.

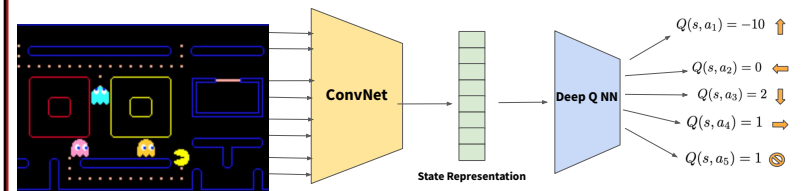
How can we use deep learning to improve classical RL techniques?



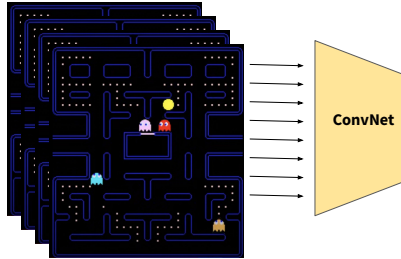
## Discussion: Which do we prefer?



## Play Pacman with DQN



We feed in multiple



Training DQN

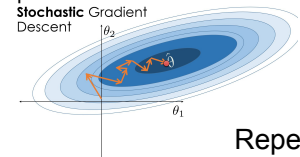
Naive Approach

Agent performs action



$$(s_t, a_t, r_t, s'_t)$$

Gradient step



Repeat

Experience Replay

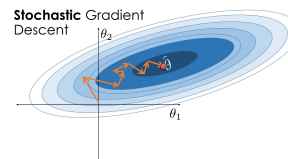
Store Transitions

- $(s_1, a_1, r_1, s'_1)$
- ...
- ...
- ...
- $(s_t, a_t, r_t, s'_t)$

Sample Minibatch

- $(s_1, a_1, r_1, s'_1)$
- ...
- $(s_j, a_j, r_j, s'_j)$

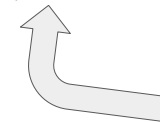
Training



- Efficient use of training data
- Avoid forgetting previous experiences
- Remove correlations

How do we compute loss?

$$\|y - Q(s, a)\|_2^2$$



"True" Q value

### Temporal Difference Target

- Consider the Bellman update

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

- When does this converge?

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

Discussion: Why might issues arise with this loss function?

$$\left\| \underbrace{r + \gamma \max_a Q(s', a)}_{\text{Target}} - \underbrace{Q(s, a)}_{\text{Prediction}} \right\|_2^2$$

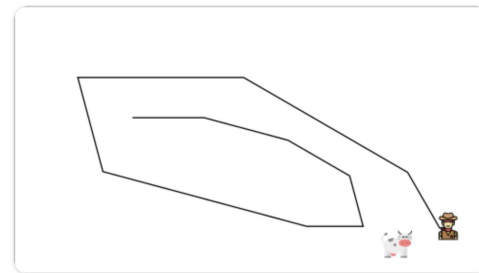
### Chasing a moving target!

When  $Q(s, a)$  changes, so does  $y$



### Chasing a moving target!

When  $Q(s, a)$  changes, so does  $y$



## Solution: Fix the target to stabilize training

$$\|r + \gamma \max_a \overset{\text{lock}}{\boxed{Q}}(s', a) - Q(s, a)\|_2^2$$

- Use a frozen network to compute the target
- Update the frozen network periodically

## Deep Q-learning

- Sampling
  - Perform actions and store the observed experience tuples in a replay memory
  - Tuples are of the form  $(s_t, a_t, r_t, s'_t)$
- Training
  - Select a batch from the replay buffer
  - Update neural network based on the batch

**Algorithm 1** Deep Q-Learning with Replay Memory and Target Network

---

```

1: Initialize replay memory  $D$  to capacity  $N$ 
2: Initialize action-value function  $Q$  with random weights  $\theta$ 
3: Initialize target action-value function  $\hat{Q}$  with weights  $\theta^- = \theta$ 
4: for episode = 1,  $M$  do
5:   Initialize state  $s_1$ 
6:   for  $t = 1, T$  do
7:     With probability  $\epsilon$  select a random action  $a_t$ 
8:     otherwise select  $a_t = \arg \max_a Q(s_t, a; \theta)$ 
9:     Execute action  $a_t$  in simulator and observe reward  $r_t$  and new state  $s_{t+1}$ 
10:    Store transition  $(s_t, a_t, r_t, s_{t+1})$  in  $D$ 
11:  end for
12:  Sample random minibatch of transitions  $(s_j, a_j, r_j, s_{j+1})$  from  $D$ 
13:  for each transition  $(s_j, a_j, r_j, s_{j+1})$  in the minibatch do
14:    if episode terminates at step  $j + 1$  then
15:      Set  $y_j = r_j$ 
16:    else
17:      Set  $y_j = r_j + \gamma \max_{a'} \hat{Q}(s_{j+1}, a'; \theta^-)$ 
18:    end if
19:  end for
20:  Perform a gradient descent step on  $\frac{1}{|B|} \sum_j (y_j - Q(s_j, a_j; \theta))^2$  with respect to  $\theta$ 
21:   $\theta^- = \tau \theta + (1 - \tau) \theta^-$ 
22: end for

```

---

## Overestimation in DQN

$$y = r + \max_a Q(s', a; \theta^-)$$



Taking the max of noisy  
random variables

### Overestimation Intuition

- There are 300 people with the same weight: 150 pounds
- There is a weight scale that measures with an error of +/- 5
- Suppose someone (who doesn't know the original weights) and wants to estimate the max weight
  - Method 1:
    - sample n people and weigh them to obtain a list of weights  $x_1, \dots, x_n$
    - Output  $\max(x_1, \dots, x_n)$
  - Method 2:
    - Weigh each of the n people twice to obtain  $x_1, \dots, x_n$  and  $x'_1, \dots, x'_n$
    - Find  $i = \operatorname{argmax}(x_1, \dots, x_n)$  and then output  $x'_i$

### Overestimation in DQN

$$y = r + \max_a Q(s', a; \theta^-)$$



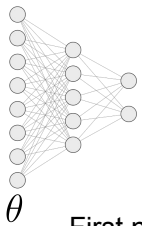
Taking the max of noisy random variables

Idea: Use different networks to select action & evaluate action

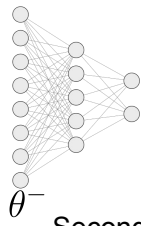
### Double DQN

$$a^* = \operatorname{argmax}_a Q(s', a; \theta)$$

$$y = r + \gamma Q(s', a^*; \theta^-)$$



First network selects best action



Second network evaluates action

Discuss: What's the difference between Fixed Q-Targets & Double DQN?

## Cornell Bowers CIS

### Algorithm 2 Double Deep Q-Learning

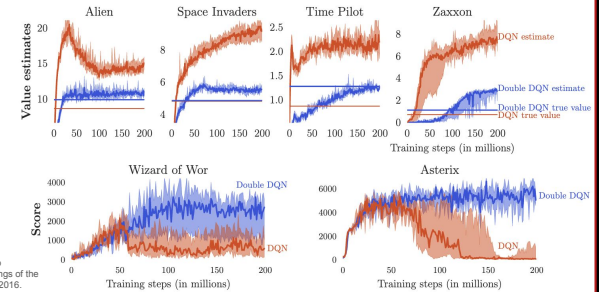
```

1: Initialize replay memory  $D$  to capacity  $N$ 
2: Initialize action-value function  $Q$  with random weights  $\theta$ 
3: Initialize target action-value function  $Q^-$  with weights  $\theta^- = \theta$ 
4: for episode = 1,  $M$  do
5:   Initialize state  $s_1$ 
6:   for  $t = 1, T$  do
7:     With probability  $\epsilon$  select a random action  $a_t$ 
8:     otherwise select  $a_t = \arg \max_a Q(s_t, a; \theta)$ 
9:     Execute action  $a_t$  in simulator and observe reward  $r_t$  and new state  $s_{t+1}$ 
10:    Store transition  $(s_t, a_t, r_t, s_{t+1})$  in  $D$ 
11:  end for
12:  Sample random minibatch of transitions  $(s_j, a_j, r_j, s_{j+1})$  from  $D$ 
13:  for each transition  $(s_j, a_j, r_j, s_{j+1})$  in the minibatch do
14:    if episode terminates at step  $j + 1$  then
15:      Set  $y_j = r_j$ 
16:    else
17:       $a^* = \max_{a'} Q(s_{j+1}, a'; \theta)$ 
18:      Set  $y_j = r_j + \gamma Q(s_{j+1}, a^*; \theta^-)$ 
19:    end if
20:  end for
21:  Perform a gradient descent step on  $\frac{1}{|B|} \sum_j (y_j - Q(s_j, a_j; \theta))^2$  with respect to  $\theta$ 
22:   $\theta^- = \tau \theta + (1 - \tau) \theta^-$ 
23: end for
  
```

## Cornell Bowers CIS

### Impact of Double DQN

- Double DQN reduces over-estimation
- Stabilizes training



## Cornell Bowers CIS

### Summary of Models

Q Learning	Vanilla DQN	Double DQN
Basic reinforcement learning algorithm	Use ConvNet to represent environment	Fixes the overestimation problem by using different networks to find and evaluate the best action
Tabular storage of q-values	Neural net for Q-values	

## Cornell Bowers CIS

### Review

- Tabular q-learning does not scale well when there a lot of states
- Deep q-learning uses a neural networks and is able to handle problems with large state-action spaces
- Issues with vanilla deep q-learning:
  - Correlation between subsequent time steps
  - Moving target
- Vanilla deep q-learning can be further improved by fixing the overestimation problem and using double deep q-learning