## COSC 689: Deep Reinforcement Learning Spring 2019

Prof. Grace Hui Yang Department of Computer Science Georgetown University

## Midterm

- Apr 3, in class, open book
- Written questions, covering:
  - MDP
  - Dynamic Programming (Value iteration, policy iteration, policy improvement, policy evaluation)
  - PG and AC (REINFORCE, REINFORCE w/baseline, AC, A3C, TRPO, PPO, GAE)
  - TD (SARSA, Q-learning, expected SARSA, double Q-learning, DQN, DDPG)

## Outline

- Deep Q-learning
- Deep Q-Network (DQN)
- Deep Deterministic Policy Gradient (DDPG)
- What are working?
  - Experience replay
  - Target network
  - Double Q-Learning
  - N-Step Q-learning
- Practical considerations for Deep Q-learning

# Q-learning

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

# Deep Q-Learning

- Use a deep neural network to learn the Q-function
- How:
  - Parameterize Q function with  $\phi$
  - Use SGD to find the optimal  $\phi$  and the optimal  $\, Q_{\phi} \,$

# Deep Q-Learning



Image from UC Berkeley CS294-112

• Also known as Online Q-iteration algorithm:

Repeat:

Take some action  $a_i$  and observe  $(S_i, a_i, S'_i, r_i)$ Calculate TD error's target  $y_i = r(s_i, a_i) + \gamma \max_{a'} Q_{\phi}(s'_i, a'_i)$ Update  $\phi \leftarrow \phi - \alpha \frac{dQ_{\phi}}{d\phi}(s_i, a_i)[Q_{\phi}(s_i, a_i) - y_i]$ 

• Also known as Online Q-iteration algorithm:

Repeat:

Take some action  $a_i$  and observe  $(s_i, a_i, s'_i, r_i)$ Calculate TD error's target  $y_i = r(s_i, a_i) + \gamma \max_{a'} Q_{\phi}(s'_i, a'_i)$ Update  $\phi \leftarrow \phi - \alpha \frac{dQ_{\phi}}{d\phi}(s_i, a_i)[Q_{\phi}(s_i, a_i) - y_i]$ Target TD error

- Suffers from two major issues:
  - Correlated samples in sequential states
    - The neural network used to learn Q is too easily to overfit to with current episodes
    - Once it is overfitted, it won't be able to generate varying experiences
    - It gets stuck!

Image from UC Berkeley CS294-112

• A moving target

• It is because the target  $y_i = r(s_i, a_i) + \gamma \max_{a'} Q_{\phi}(s'_i, a'_i)$  depends on the same parameters that we are trying to learn  $-\phi$ 

- It becomes a chicken-and-egg problem and the learning becomes super unstable
- And when we use SGD to seek for  $\phi$  , the gradient is not though the target function part (see next slide)

• Also known as Online Q-iteration algorithm:

Repeat:

Take some action  $a_i$  and observe  $(s_i, a_i, s_i', r_i)$ Calculate TD error's target  $y_i = r(s_i, a_i) + \gamma \max_{a'} Q_{\phi}(s_i', a_i')$ Update  $\phi \leftarrow \phi - \alpha \frac{dQ_{\phi}}{d\phi}(s_i, a_i)[Q_{\phi}(s_i, a_i) - y_i]$  $= \phi - \alpha \frac{dQ_{\phi}}{d\phi}(s_i, a_i)[Q_{\phi}(s_i, a_i) - (r(s_i, a_i) + \gamma \max_{a'} Q_{\phi}(s_i', a_i'))]$ 

gradient is not through this part

### Solutions to Correlated samples - asynchronous learning

 Solution 1: asynchronous learning to break the correlations

synchronized parallel Q-learning



asynchronous parallel Q-learning



Image from UC Berkeley CS294-112

### Solutions to Correlated samples - Experience Replay

- Solution 2: sample the recent states from a buffer
- This technique is called "experience replay"



Image from UC Berkeley CS294-112

# **Experience Replay**

- Idea:
  - Stores the experiences, including state transitions, rewards, actions, into a buffer
    - What are stored are those that are essential to learn a Q-function
      - state transitions (from s to s'; note it is not the entire trajectory, but pairs of Markov transitions)
      - rewards
      - actions
  - Then, random samples (mini-batches) from this buffer to update the neural network

# **Experience Replay**

- Benefits:
  - It reduces the correlation between recent experiences when updating the NN
  - With a good size of mini-batches, it can speed up the learning
  - reuses past experience to avoid sudden change/ forgetting

# **Experience Replay**

- How big the buffer should be?
  - the replay buffer should be big enough to contain a wide range of experiences
  - however, not keep everything. It would be very inefficient
- How recent the experiences should be?
  - If you use very recent date only, then it again easily overfit to recent samples
  - If you use a lot of past experiences, which means to keep more data in the buffer, it will slow you down
- Tuning is thus important

## Naive DQN + Experience Replay

Repeat:

collect experiences  $\{(s_i, a_i, s_i', r_i)\}$ , add them to replay buffer B

Repeat:

sample a batch of  $(s_i, a_i, s'_i, r_i)$  from **B** 

calculate TD error's target  $y_i = r(s_i, a_i) + \gamma \max_{a'} Q_{\phi}(s'_i, a'_i)$ update  $\phi \leftarrow \phi - \alpha \frac{dQ_{\phi}}{d\phi}(s_i, a_i)[Q_{\phi}(s_i, a_i) - y_i]$ 

# Issue - Moving Target

Repeat:

Take some action  $a_i$  and observe  $(s_i, a_i, s'_i, r_i)$ 

Calculate TD error's target  $y_i = r(s_i, a_i) + \gamma \max_{a'} Q_{\phi}(s'_i, a'_i)$ Update  $\phi \leftarrow \phi - \alpha \frac{dQ_{\phi}}{d\phi}(s_i, a_i)[Q_{\phi}(s_i, a_i) - (r(s_i, a_i) + \gamma \max_{a'} Q_{\phi}(s'_i, a'_i))]$ 

gradient is not through this part

- The targets don't change in the inner loop
- Q-learning (Naive DQN) is not gradient descent
- The learning is super unstable

## Solution to moving target -Target Network

- Idea:
  - Use a separate target network (parameterized by  $\phi'$ ) in Q-value fitting
  - In the new target network, we use a set of parameters that is close to  $\phi$  , but not exactly it, with a time delay
  - Ways to do it:
    - Directly copy over the learned  $\phi$  from Q-network to the target network (as in DQN):

$$\phi' \leftarrow \phi$$

• Use Polyak averaging (as in DDPG)

 $\phi' \leftarrow \tau \phi' + (1 - \tau) \phi$   $\tau = 0.999$  for instance

## Naive DQN + Experience Replay + Target Network

Repeat:

Update target network parameter by copying:  $\phi' \leftarrow \phi$ 

Repeat (N times):

collect experiences  $\{(s_i, a_i, s_i', r_i)\}$ , add them to  $\boldsymbol{B}$ 

Repeat (K times):

sample a batch of  $\{(s_j, a_j, s_j', r_j)\}$  from B

compute TD error's target using target network parameters  $y_j = r(s_j, a_j) + \gamma \max_{a'} Q'_{\phi}(s'_j, a'_j)$ update  $\phi \leftarrow \phi - \alpha \frac{dQ_{\phi}}{d\phi}(s_j, a_j)[Q_{\phi}(s_j, a_j) - y_j]$ 

## Naive DQN + Experience Replay + Target Network

Repeat:

Update target network parameter by copying:  $\phi' \leftarrow \phi$ Repeat (N times): collect experiences  $\{(s_i, a_i, s'_i, r_i)\}$ , add them to BRepeat (K times): sample a batch of  $\{(s_j, a_j, s'_j, r_j)\}$  from BNote the index changes from i to j

parameters

compute TD error's target using target network ers  $y_j = r(s_j, a_j) + \gamma \max_{a'} Q'_{\phi}(s'_j, a'_j)$ update  $\phi \leftarrow \phi - \alpha \frac{dQ_{\phi}}{d\phi}(s_j, a_j)[Q_{\phi}(s_j, a_j) - y_j]$ 

## Illustration



- Online Q-learning (last lecture): evict immediately, process 1, process 2, and process 3 all run at the same speed •
- DQN: process 1 and process 3 run at the same speed, process 2 is slow

Image from UC Berkeley CS294-112

# Deep Q-Network (DQN)

- It is Naive DQN + Experience Replay + Target Network, when N=1 and K=1.
- DQN Paper: Human-level control through deep reinforcement learning: Q-learning with convolutional networks for playing Atari. <u>https://www.nature.com/</u> <u>articles/nature14236</u>

## DQN

Repeat:

Update target network parameter by copying:  $\phi' \leftarrow \phi$ 

Repeat (N times): N=1 and K=1

collect experiences  $\{(s_i, a_i, s_i', r_i)\}$ , add them to B

Repeat (K times):

sample a batch of  $\{(s_j, a_j, s'_j, r_j)\}$  from B

compute TD error's target using target network parameters  $y_j = r(s_j, a_j) + \gamma \max_{a'} Q'_{\phi}(s'_j, a'_j)$ update  $\phi \leftarrow \phi - \alpha \frac{dQ_{\phi}}{d\phi}(s_j, a_j)[Q_{\phi}(s_j, a_j) - y_j]$ 

#### It is your Homework 3.

# DQN Algo in the Paper

Algorithm 1: deep Q-learning with experience replay. Initialize replay memory D to capacity N Initialize action-value function Q with random weights  $\theta$ Initialize target action-value function Q with weights  $\theta^- = \theta$ For episode = 1, M do Initialize sequence  $s_1 = \{x_1\}$  and preprocessed sequence  $\phi_1 = \phi(s_1)$ For t = 1,T do With probability  $\varepsilon$  select a random action  $a_t$ otherwise select  $a_t = \operatorname{argmax}_a Q(\phi(s_t), a; \theta)$ Execute action  $a_t$  in emulator and observe reward  $r_t$  and image  $x_{t+1}$ Set  $s_{t+1} = s_t, a_t, x_{t+1}$  and preprocess  $\phi_{t+1} = \phi(s_{t+1})$ Store transition  $(\phi_t, a_t, r_t, \phi_{t+1})$  in D Sample random minibatch of transitions  $(\phi_j, a_j, r_j, \phi_{j+1})$  from D  $\operatorname{Set} y_{j} = \begin{cases} r_{j} & \text{if episode terminates at step } j+1 \\ r_{j} + \gamma \max_{a'} \hat{Q}(\phi_{j+1}, a'; \theta^{-}) & \text{otherwise} \end{cases}$ Perform a gradient descent step on  $(y_j - Q(\phi_j, a_j; \theta))^2$  with respect to the network parameters  $\theta$ Every C steps reset  $\bar{Q} = Q$ End For End For

## Further Improvement: Double Q-learning

- Overestimation in Q-learning
- (We mentioned in last lecture) It is because Q-learning is a greedy algorithm
  - Or at least \epsilon-greedy
  - It can be biased towards overestimation of the next value

target value 
$$y_j = r_j + \gamma \max_{\mathbf{a}'_j} Q_{\phi'}(\mathbf{s}'_j, \mathbf{a}'_j)$$

### Overestimation in Q-Learning



# **Double Q-Learning**

- Idea:
  - Do not use the same network for choosing actions and evaluating values
  - Instead, use two different networks to do the above (with 0.5 probably for each)

$$Q_{\phi_A}(\mathbf{s}, \mathbf{a}) \leftarrow r + \gamma Q_{\phi_B}(\mathbf{s}', \arg\max_{\mathbf{a}'} Q_{\phi_A}(\mathbf{s}'))$$
$$Q_{\phi_B}(\mathbf{s}, \mathbf{a}) \leftarrow r + \gamma Q_{\phi_A}(\mathbf{s}', \arg\max_{\mathbf{a}'} Q_{\phi_B}(\mathbf{s}'))$$

 In practice, we use target network to evaluate Q-value and current network to choose action

# Double Q-Learning

 In practice, we use target network to evaluate Q-value and <u>use current network to choose action</u>:

$$y_j = r(s_j, a_j) + \gamma Q_{\phi'}(s'_j, \arg\max_{a'} Q_{\phi}(s'_j, a'_j))$$

# Double Q-Learning

 In practice, we use target network to evaluate Q-value and <u>use current network to choose action</u>:

$$y_j = r(s_j, a_j) + \gamma Q_{\phi'}(s'_j, \arg \max_{a'} Q_{\phi}(s'_j, a'_j))$$
Not  $Q_{\phi'}$ 

## Further Improvement: N-Step Return

- Q-learning's target is high in bias and low in variance.
- It is because it only looks ahead 1 step and 1 reward

$$y_j = r(s_j, a_j) + \gamma Q_{\phi'}(s'_j, \arg\max_{a'} Q_{\phi}(s'_j, a'_j))$$

- How to improve (to make it low in both)?
  - We can use N-step returns as in actor-critic

## Q-learning + N-Step Return

• Q-learning target becomes:

$$y_{j,t} = \sum_{t'=t}^{t'+N-1} r_{j,t'} + \gamma^N \max_{\mathbf{a}_{j,t+N}} Q_{\phi'}(\mathbf{s}_{j,t+N}, \mathbf{a}_{j,t+N})$$

• When we do action selection,

$$\pi(\mathbf{a}_t | \mathbf{s}_t) = \begin{cases} 1 \text{ if } \mathbf{a}_t = \arg \max_{\mathbf{a}_t} Q_\phi(\mathbf{s}_t, \mathbf{a}_t) \\ 0 \text{ otherwise} \end{cases}$$

we will need transitions  $\mathbf{s}_{j,t'}, \mathbf{a}_{j,t'}, \mathbf{s}_{j,t'+1}$  all to come from the same policy  $\pi$  for all t' - t < N - 1

It is hard to do though - so we sometimes ignore the problem

# Summary: DQN

- What are working?
  - Experience replay
  - Use a Target network
  - Double Q-Learning
  - N-Step Q-learning

## Deep Q-learning with Continuous Actions

• DQN select actions based on the following or \epsilon-greedy:

$$\pi(\mathbf{a}_t | \mathbf{s}_t) = \begin{cases} 1 \text{ if } \mathbf{a}_t = \arg \max_{\mathbf{a}_t} Q_\phi(\mathbf{s}_t, \mathbf{a}_t) \\ 0 \text{ otherwise} \end{cases}$$

- This works well with discrete actions
- However, when use DQN with continuous actions, we cannot exhaustively evaluate all actions in the space to find out the max
- Moreover, solving this optimization (the max) is in the inner loop of calculating the target, which makes the optimization highly nontrivial

## Deep Q-learning with Continuous Actions

• DQN select actions based on the following or \epsilon-greedy:

$$\pi(\mathbf{a}_t | \mathbf{s}_t) = \begin{cases} 1 \text{ if } \mathbf{a}_t = \arg \max_{\mathbf{a}_t} Q_\phi(\mathbf{s}_t, \mathbf{a}_t) \\ 0 \text{ otherwise} \end{cases}$$

This works well with discrete actions

- a painfully expensive subroutine
- However, when use DQN with continuous actions, we cannot exhaustively evaluate all actions in the space to find out the max
- Moreover, solving this optimization (the max) is in the inner loop of calculating the target, which makes the optimization highly nontrivial

## Deep Q-learning with Continuous Actions

- Solution:
  - Parameterize the entire argmax part  $\frac{rg \max_{\mathbf{a}_t} Q_{\phi}(\mathbf{s}_t, \mathbf{a}_t)}{\mathbf{b}_t}$  by heta
  - Then, learn a maximizer (another network)  $\mu_{\theta}$  to approximate it  $\arg \max_{\mathbf{a}_t} Q_{\phi}(\mathbf{s}_t, \mathbf{a}_t)$
  - Assumption: the Q-function is differentiable w.r.t the action  $\frac{dQ_{\phi}}{d\theta} = \frac{d\mathbf{a}}{d\theta} \frac{dQ_{\phi}}{d\mathbf{a}}$
  - This leads to DDPG

## Deep Deterministic Policy Gradient (DDPG)

- train another network  $\mu_{\theta}(\mathbf{s})$  such that  $\mu_{\theta}(\mathbf{s}) \approx \arg \max_{\mathbf{a}} Q_{\phi}(\mathbf{s}, \mathbf{a})$
- Then, solve  $\theta \leftarrow \arg \max_{\theta} Q_{\phi}(\mathbf{s}, \mu_{\theta}(\mathbf{s}))$
- After this action selection, the new target becomes

$$y_j = r_j + \gamma \max_{\mathbf{a}'_j} Q_{\phi'}(\mathbf{s}'_j, \mu_{\theta}(\mathbf{s}'_j))$$

 DDPG paper: Continuous control with deep reinforcement learning: continuous Q-learning with actor network for approximate maximization.

## DDPG

1. take some action a<sub>i</sub> and observe (s<sub>i</sub>, a<sub>i</sub>, s'<sub>i</sub>, r<sub>i</sub>), add it to β
2. sample mini-batch {s<sub>j</sub>, a<sub>j</sub>, s'<sub>j</sub>, r<sub>j</sub>} from β uniformly
3. compute y<sub>j</sub> = r<sub>j</sub> + γ max<sub>a'<sub>j</sub></sub> Q<sub>φ'</sub>(s'<sub>j</sub>, μ<sub>θ'</sub>(s'<sub>j</sub>)) using target nets Q<sub>φ'</sub> and μ<sub>θ'</sub>
4. φ ← φ − α ∑<sub>j</sub> dQ<sub>φ</sub>/dφ (s<sub>j</sub>, a<sub>j</sub>)(Q<sub>φ</sub>(s<sub>j</sub>, a<sub>j</sub>) − y<sub>j</sub>)
5. θ ← θ + β ∑<sub>j</sub> dμ/dθ (s<sub>j</sub>) dQ<sub>φ</sub>/da (s<sub>j</sub>, a)
6. update φ' and θ' (e.g., Polyak averaging)

 DDPG paper: Continuous control with deep reinforcement learning: continuous Q-learning with actor network for approximate maximization.

Adapted from UC Berkeley CS294-112

## **DDPG** in the Paper

#### Algorithm 1 DDPG algorithm

Randomly initialize critic network  $Q(s, a|\theta^Q)$  and actor  $\mu(s|\theta^\mu)$  with weights  $\theta^Q$  and  $\theta^\mu$ . Initialize target network Q' and  $\mu'$  with weights  $\theta^{Q'} \leftarrow \theta^Q$ ,  $\theta^{\mu'} \leftarrow \theta^\mu$ Initialize replay buffer Rfor episode = 1, M do Initialize a random process  $\mathcal{N}$  for action exploration Receive initial observation state  $s_1$ for t = 1, T do Select action  $a_t = \mu(s_t|\theta^\mu) + \mathcal{N}_t$  according to the current policy and exploration noise Execute action  $a_t$  and observe reward  $r_t$  and observe new state  $s_{t+1}$ Store transition  $(s_t, a_t, r_t, s_{t+1})$  in RSample a random minibatch of N transitions  $(s_i, a_i, r_i, s_{i+1})$  from RSet  $y_i = r_i + \gamma Q'(s_{i+1}, \mu'(s_{i+1}|\theta^{\mu'})|\theta^{Q'})$ Update critic by minimizing the loss:  $L = \frac{1}{N} \sum_i (y_i - Q(s_i, a_i|\theta^Q))^2$ Update the actor policy using the sampled policy gradient:

$$abla_{ heta^{\mu}} J pprox rac{1}{N} \sum_{i} 
abla_{a} Q(s, a | heta^{Q})|_{s=s_{i}, a=\mu(s_{i})} 
abla_{ heta^{\mu}} \mu(s | heta^{\mu})|_{s_{i}}$$

Update the target networks:

$$\begin{aligned} \theta^{Q'} &\leftarrow \tau \theta^Q + (1 - \tau) \theta^{Q'} \\ \theta^{\mu'} &\leftarrow \tau \theta^\mu + (1 - \tau) \theta^{\mu'} \end{aligned}$$

end for end for

## DDPG

- It is an actor-critic method
- The value learning part (critic) is similar to DQN
- The action section (actor) is using Policy Gradient
- It also uses tricks including
  - experience replay
  - target network

## **Practical Tips for Q-learning**

- Test on easy, reliable tasks first, make sure your implementation is correct
- Large replay buffers help improve stability
- Start with high exploration (\epsilon) and gradually reduce
- Start with high learning rate (\alpha) and gradually reduce
- Q-learning could be slow, so be patient
- Double Q-learning helps a lot in practice, simple and no downsides
- N-step returns also help a lot, but have some downsides
- Run multiple random seeds, it could be inconsistent between runs

## Summary: Deep TD-Learning

- Widely used value-based method family
- Connects well with psychology and neuroscience
  - TD error
- DQN is a state-of-art RL method
- DDPG can be used for continuous action space; and works as an actor-critic method
- Key techniques to increase training stability:
  - experience replay
  - target network

## References

- Main Textbook: RL book chapters 6.
  - Note that some images or formulas from the main textbook are not individually cited. You should assume they are from the textbook.
- DQN Paper: Human-level control through deep reinforcement learning: Q-learning with convolutional networks for playing Atari.
- DDPG paper: Continuous control with deep reinforcement learning: continuous Q-learning with actor network for approximate maximization.
- UC Berkeley CS294-112 Deep Reinforcement Learning