

# Value-based Reinforcement Learning

## Some Discussions

Kan Ren

Apex Data and Knowledge Management Lab  
Shanghai Jiao Tong University

Aug. 3 2017



# Outline

- 1 Revision of Value-based RL
  - Dynamic Programming(omitted)
  - Monte-carlo Method(omitted)
  - TD: Sarsa and Q-learning
- 2 Deep Q-network
  - Nature DQN
  - Several Improvements
- 3 Issues in Q-learning
  - Overestimation
  - Double Q-learning
  - Averaged Q-learning
- 4 Convergence of Tabular TD
  - Sarsa
  - Q-learning (TBE)



# Outline

- 1 Revision of Value-based RL
  - Dynamic Programming(omitted)
  - Monte-carlo Method(omitted)
  - TD: Sarsa and Q-learning
- 2 Deep Q-network
  - Nature DQN
  - Several Improvements
- 3 Issues in Q-learning
  - Overestimation
  - Double Q-learning
  - Averaged Q-learning
- 4 Convergence of Tabular TD
  - Sarsa
  - Q-learning (TBE)



# Outline

- 1 Revision of Value-based RL
  - Dynamic Programming(omitted)
  - **Monte-carlo Method(omitted)**
  - TD: Sarsa and Q-learning
- 2 Deep Q-network
  - Nature DQN
  - Several Improvements
- 3 Issues in Q-learning
  - Overestimation
  - Double Q-learning
  - Averaged Q-learning
- 4 Convergence of Tabular TD
  - Sarsa
  - Q-learning (TBE)



# Outline

- 1 Revision of Value-based RL
  - Dynamic Programming(omitted)
  - Monte-carlo Method(omitted)
  - TD: Sarsa and Q-learning
- 2 Deep Q-network
  - Nature DQN
  - Several Improvements
- 3 Issues in Q-learning
  - Overestimation
  - Double Q-learning
  - Averaged Q-learning
- 4 Convergence of Tabular TD
  - Sarsa
  - Q-learning (TBE)



# Sarsa & Q-learning

## Algorithm

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

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  $S$   
   Choose  $A$  from  $S$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)  
   Repeat (for each step of episode):  
     Take action  $A$ , observe  $R, S'$   
     Choose  $A'$  from  $S'$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -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

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

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



# Difference

- Exploration
  - Sarsa: on-policy
  - Q-learning: off-policy

- Update Rule

- Sarsa

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

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

- Q-learning

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



# Outline

- 1 Revision of Value-based RL
  - Dynamic Programming(omitted)
  - Monte-carlo Method(omitted)
  - TD: Sarsa and Q-learning
- 2 Deep Q-network
  - Nature DQN
  - Several Improvements
- 3 Issues in Q-learning
  - Overestimation
  - Double Q-learning
  - Averaged Q-learning
- 4 Convergence of Tabular TD
  - Sarsa
  - Q-learning (TBE)





# Outline

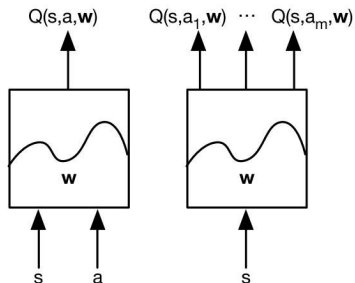
- 1 Revision of Value-based RL
  - Dynamic Programming(omitted)
  - Monte-carlo Method(omitted)
  - TD: Sarsa and Q-learning
- 2 Deep Q-network
  - Nature DQN
  - **Several Improvements**
- 3 Issues in Q-learning
  - Overestimation
  - Double Q-learning
  - Averaged Q-learning
- 4 Convergence of Tabular TD
  - Sarsa
  - Q-learning (TBE)



# Q-networks

Represent value function by **Q-network** with weights  $w$

$$Q(s, a; w) \approx Q^*(s, a) . \quad (1)$$



# Deep Q-network

Refer to D. Silver's slides P31 - P45.



# Duelling network

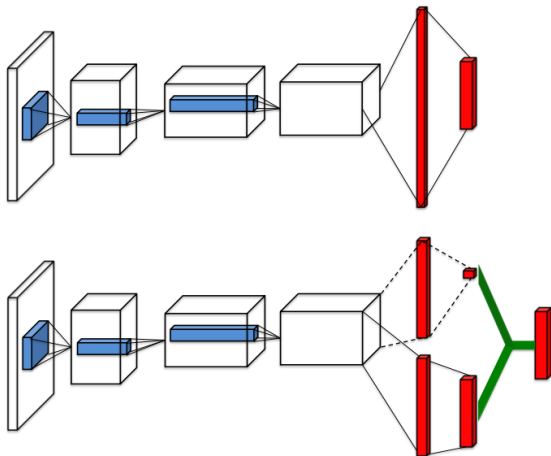


Figure: Duelling network: split Q-network into two channels

# Outline

- 1 Revision of Value-based RL
  - Dynamic Programming(omitted)
  - Monte-carlo Method(omitted)
  - TD: Sarsa and Q-learning
- 2 Deep Q-network
  - Nature DQN
  - Several Improvements
- 3 Issues in Q-learning
  - Overestimation
  - Double Q-learning
  - Averaged Q-learning
- 4 Convergence of Tabular TD
  - Sarsa
  - Q-learning (TBE)



# Overestimation

## Preliminaries

Recall that

$$Q(s, a) \leftarrow r_s^a + \gamma \max_{\hat{a}} Q(s', \hat{a}) \quad (2)$$

Repeated application of this update equation eventually yields Q-values that give rise to a policy which maximizes the expected cumulative discounted reward<sup>1</sup> in the look-up table case.

The **max** operation may cause some problems under the approximation scenario.

---

<sup>1</sup>C. J. C. H. Watkins, Learning from Delayed Rewards. PhD thesis, Kings College, Cambridge, England, 1989.



# Overestimation

Assume  $Q^{approx}(\cdot)$  representing implicit target values  $Q^{target}$ , corrupted by a noise term  $Y$  such that

$$Q^{approx}(s', \hat{a}) = Q^{target}(s', \hat{a}) + Y_{s'}^{\hat{a}}$$

$$\begin{aligned} Z_s &\stackrel{def}{=} r_s^a + \gamma \max_{\hat{a}} Q^{approx}(s', \hat{a}) - \left( r_s^a + \gamma \max_{\hat{a}} Q^{target}(s', \hat{a}) \right) \\ &= \gamma \left( \max_{\hat{a}} Q^{approx}(s', \hat{a}) - \max_{\hat{a}} Q^{target}(s', \hat{a}) \right) \end{aligned} \quad (3)$$

The key observation is

$$E[Y_{s'}^{\hat{a}}] = 0, \quad \forall \hat{a} \xrightarrow{\text{often}} E[Z_s] > 0.$$



# Expectation of $Z$

## Lemma

Let  $n$  denote the number of actions applicable at state  $s'$ . If all  $n$  actions share the same target  $Q$ -value, i.e.,  $\exists q : \forall \hat{a} : q = Q^{\text{target}}(s', \hat{a})$ , then the average overestimation  $E[Z_s]$  is  $\gamma c$  with  $c \stackrel{\text{def}}{=} \epsilon \frac{n-1}{n+1}$ .

The proof can be referred to the paper<sup>2</sup>.

## Corollary

$0 \leq E[Z_s] \leq \gamma c$  with  $c = \epsilon \frac{n-1}{n+1}$ .

<sup>2</sup>Thrun S, Schwartz A. Issues in using function approximation for reinforcement learning[C] Proceedings of the 1993 Connectionist Models Summer School Hillsdale, NJ. Lawrence Erlbaum. 1993.



# Bounds for Expected Failure of Q-learning

## Simple Assumptions

- There is a set of goal states;
- Positive reward  $r_{goal}$  is only received upon entering a goal state;
- $r_{goal} = 1$ ;
- The state transition function is deterministic.

One *necessary* condition for the success of Q-learning is that the sequence of Q-values  $Q(s_i, a_i)$  is monotonically increasing in  $i$ :

$$Q(s_i, a_i) \leq Q(s_{i+1}, a_{i+1}), \text{ for all } i \in \{0, \dots, L-1\} \quad (4)$$



# Bounds for Expected Failure of Q-learning

## Simple Assumptions

Case 1: the learner *always* overestimates Q-values by  $\gamma c$ .

### Theorem

*If there is maximal, repeated overestimation of magnitude  $\gamma c$  along an optimal path, Q-learning is expected to fail to learn an optimal policy if*

$$\gamma > \frac{1}{1+c}.$$



Case 2: Assume that Q-learning managed to learn the *last*  $L-1$  Q-values of this optimal path correctly.

- Q-values are given by iteratively discounting the final reward with the *distance* to the goal state, i.e.,  $Q(s_{L-i}, a_{L-i}) = \gamma^i$  for  $i \in \{1, \dots, L-1\}$ .
- Correct Q-value  $Q^{correct}(s_0, a_0)$  is  $\gamma^L$ .
- In order to maintain monotonicity of Q, we need to make sure that

$$\gamma^{L-1} - \gamma^L \geq \gamma c . \quad (5)$$

## Theorem

*Under the conditions above, Q-learning is expected to fail if*

$$\gamma^{L-1} - \gamma^L < \gamma c . \quad (6)$$

## Theorem

*Under the conditions above, Q-learning is expected to fail if*

$$\epsilon > \frac{n+1}{n-1} \cdot \frac{(L-2)^{L-2}}{(L-1)^{L-1}} . \quad (7)$$



# Outline

- 1 Revision of Value-based RL
  - Dynamic Programming(omitted)
  - Monte-carlo Method(omitted)
  - TD: Sarsa and Q-learning
- 2 Deep Q-network
  - Nature DQN
  - Several Improvements
- 3 Issues in Q-learning
  - Overestimation
  - **Double Q-learning**
  - Averaged Q-learning
- 4 Convergence of Tabular TD
  - Sarsa
  - Q-learning (TBE)



# Double Q-learning

## Preliminaries

- a set of random variables  $X = \{X_i, \dots, X_M\}$  Our interest is that

$$\max_i E[X_i] , \quad (8)$$

which is in the Q-learning update rule.

- $S = \cup_{i=1}^M S_i$  where  $S_i$  is the subset contains samples for the variable  $X_i$  and each  $s \in S_i$  is i.i.d.
- $E[X_i] = E[\mu_i] \approx \mu_i(S) \stackrel{\text{def}}{=} \frac{1}{|S_i|} \sum_{s \in S_i} s$  , where  $\mu_i$  is an unbiased estimate for the value of  $E[X_i]$ .
- $f_i^\mu$  is PDF and  $F_i^\mu$  is CDF of  $X_i$ .

$$\max_i E[X_i] = \max_i \int_{-\infty}^{\infty} x f_i^\mu(x) dx .$$



# Double Q-learning

## Single Estimator

An obvious way to approximate the value in Eq. (8) is

$$\max_i E[X_i] = \max_i E[\mu_i] \approx \max_i \mu_i(S) . \quad (9)$$

- Assume the maximal estimator  $\max_i \mu_i(S)$  is distributed as PDF  $f_{max}^\mu$ .
- $f_{max}^\mu \neq f_i^\mu$  but  $f_{max}^\mu$  is dependent on  $f_i^\mu$ .



# Double Q-learning

## Single Estimator

An obvious way to approximate the value in Eq. (8) is

$$\max_i E[X_i] = \max_i E[\mu_i] \approx \max_i \mu_i(S). \quad (9)$$

- Assume the maximal estimator  $\max_i \mu_i(S)$  is distributed as PDF  $f_{max}^\mu$ .
- $f_{max}^\mu \neq f_i^\mu$  but  $f_{max}^\mu$  is dependent on  $f_i^\mu$ .
- CDF  $F_{max}^\mu(x) \stackrel{def}{=} P(\max_i \mu_i \leq x) = \prod_{i=1}^M P(\mu_i \leq x) \stackrel{def}{=} \prod_{i=1}^M F_i^\mu(x)$ .





# Double Q-learning

## Biased Estimation of $E[X_i]$

- The value  $\max_i \mu_i(S)$  is an unbiased estimate for  $E[\max_j \mu_j]$ .

$$\begin{aligned}
 E[\max_i \mu_i] &= \int_{-\infty}^{\infty} x f_{\max}^{\mu}(x) \\
 &= \int_{-\infty}^{\infty} x \frac{d}{dx} \prod_{i=1}^M F_i^{\mu}(x) dx \\
 &= \sum_j^M \int_{-\infty}^{\infty} x f_j^{\mu}(x) \prod_{i \neq j} F_i^{\mu}(x) dx .
 \end{aligned} \tag{10}$$

- $E[\max_i \mu_i]$  is not the same as  $\max_j E[X_j]$ .



# Double Q-learning

## Double Estimators

- Two sets of estimators:  $\mu^A = \mu_1^A, \dots, \mu_M^A$ ,  $\mu^B = \mu_1^B, \dots, \mu_M^B$ .
- Two subsets of samples:  $S = S^A \cup S^B$ ,  $S^A \cap S^B = \emptyset$
- $\mu_i^A(S) \stackrel{\text{def}}{=} \frac{1}{|S_i^A|} \sum_{s \in S_i^A} s$ ,  $\mu_i^B(S) \stackrel{\text{def}}{=} \frac{1}{|S_i^B|} \sum_{s \in S_i^B} s$ .
- Both  $\mu_i^A$  and  $\mu_i^B$  are unbiased if we assume proper split on the sample set  $S$ .
- $Max^A(S) \stackrel{\text{def}}{=} \{j | \max_i \mu_i^A(S)\}$ .
- Since  $\mu_i^B(S)$  is an independent, unbiased set of estimators, we have  $E[\mu_j^B(S)] = E[X_j]$  for all  $j$  including  $j \in Max^A$ . We can pick  $a^*$  such that  $\mu_{a^*}^A \stackrel{\text{def}}{=} \max_i \mu_i^A(S)$ . So that

$$\max_i E[X_i] = \max_i E[\mu_i^B] \approx \mu_{a^*}^B .$$



# Double Q-learning

## Difference Between Single/Double Estimators

$$\begin{aligned}
 P(j = a^*) &= \int_{-\infty}^{\infty} P(\mu_j^A = x) \prod_{i \neq j}^M P(\mu_j^A < x) dx \\
 &\stackrel{\text{def}}{=} \int_{-\infty}^{\infty} f_j^A(x) \prod_{i \neq j}^M F_i^A(x) dx
 \end{aligned} \tag{12}$$

$$\sum_j^M P(j = a^*) E[\mu_j^B] = \sum_j^M E[\mu_j^B] \int_{-\infty}^{\infty} f_j^A(x) \prod_{i \neq j}^M F_i^A(x) dx . \tag{13}$$

Recall Eq. (10) of single estimator that

$$E[\max_i \mu_i] = \int_{-\infty}^{\infty} x f_{\max}^{\mu}(x) = \sum_j^M \int_{-\infty}^{\infty} x f_j^{\mu}(x) \prod_{i \neq j}^M F_i^{\mu}(x) dx .$$



# Double Q-learning

## Algorithm<sup>34</sup>

---

### Algorithm 1 Double Q-learning

---

```

1: Initialize  $Q^A, Q^B, s$ 
2: repeat
3:   Choose  $a$ , based on  $Q^A(s, \cdot)$  and  $Q^B(s, \cdot)$ , observe  $r, s'$ 
4:   Choose (e.g. random) either UPDATE(A) or UPDATE(B)
5:   if UPDATE(A) then
6:     Define  $a^* = \arg \max_a Q^A(s', a)$ 
7:      $Q^A(s, a) \leftarrow Q^A(s, a) + \alpha(s, a) (r + \gamma Q^B(s', a^*) - Q^A(s, a))$ 
8:   else if UPDATE(B) then
9:     Define  $b^* = \arg \max_a Q^B(s', a)$ 
10:     $Q^B(s, a) \leftarrow Q^B(s, a) + \alpha(s, a) (r + \gamma Q^A(s', b^*) - Q^B(s, a))$ 
11:   end if
12:    $s \leftarrow s'$ 
13: until end

```

---

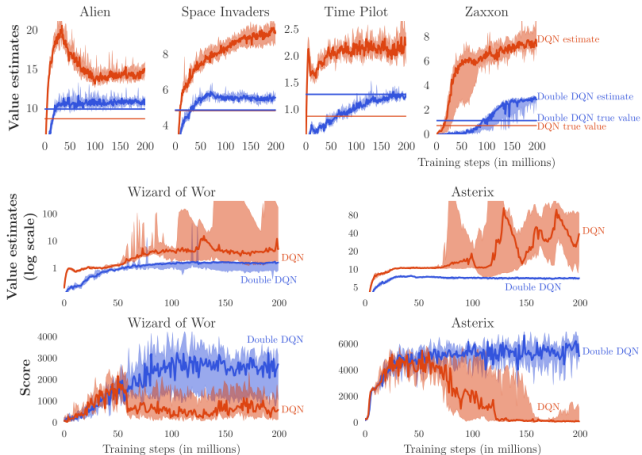
<sup>3</sup> Hasselt H V. Double Q-learning[C] Advances in Neural Information Processing Systems. 2010: 2613-2621.

<sup>4</sup> Van Hasselt H, Guez A, Silver D. Deep Reinforcement Learning with Double Q-Learning[C] AAAI. 2016: 2094-2100.



# Double Q-learning

## Performance



# Outline

- 1 Revision of Value-based RL
  - Dynamic Programming(omitted)
  - Monte-carlo Method(omitted)
  - TD: Sarsa and Q-learning
- 2 Deep Q-network
  - Nature DQN
  - Several Improvements
- 3 Issues in Q-learning
  - Overestimation
  - Double Q-learning
  - **Averaged Q-learning**
- 4 Convergence of Tabular TD
  - Sarsa
  - Q-learning (TBE)



# Averaged Deep Q-Network

- Double Q-learning aims to correct the *overestimation* of natural Q-learning.
- Averaged DQN focus on variance reduction and stabilization.



# Averaged Deep Q-Network

Revision of DQN<sup>5</sup>

---

## Algorithm 1 DQN

---

- 1: Initialize  $Q(s, a; \theta)$  with random weights  $\theta_0$
  - 2: Initialize Experience Replay (ER) buffer  $\mathcal{B}$
  - 3: Initialize exploration procedure  $Explore(\cdot)$
  - 4: **for**  $i = 1, 2, \dots, N$  **do**
  - 5:    $y_{s,a}^i = \mathbb{E}_{\mathcal{B}} [r + \gamma \max_{a'} Q(s', a'; \theta_{i-1}) | s, a]$
  - 6:    $\theta_i \approx \operatorname{argmin}_{\theta} \mathbb{E}_{\mathcal{B}} [(y_{s,a}^i - Q(s, a; \theta))^2]$
  - 7:    $Explore(\cdot)$ , update  $\mathcal{B}$
  - 8: **end for**
- output**  $Q^{\text{DQN}}(s, a; \theta_N)$
- 

<sup>5</sup>Mnih V, Kavukcuoglu K, Silver D, et al. Human-level control through deep reinforcement learning[J]. Nature, 2015, 518(7540): 529-533.





# Averaged Deep Q-Network

## Algorithm<sup>6</sup>

---

### Algorithm 2 Averaged DQN

---

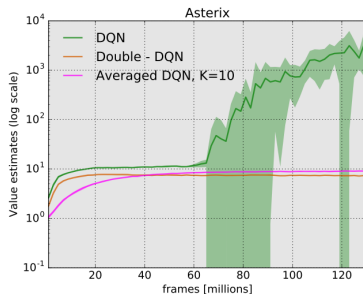
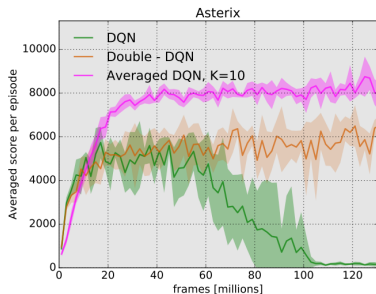
- 1: Initialize  $Q(s, a; \theta)$  with random weights  $\theta_0$
  - 2: Initialize Experience Replay (ER) buffer  $\mathcal{B}$
  - 3: Initialize exploration procedure  $Explore(\cdot)$
  - 4: **for**  $i = 1, 2, \dots, N$  **do**
  - 5:    $Q_{i-1}^A(s, a) = \frac{1}{K} \sum_{k=1}^K Q(s, a; \theta_{i-k})$
  - 6:    $y_{s,a}^i = \mathbb{E}_{\mathcal{B}} [r + \gamma \max_{a'} Q_{i-1}^A(s', a') | s, a]$
  - 7:    $\theta_i \approx \operatorname{argmin}_{\theta} \mathbb{E}_{\mathcal{B}} [(y_{s,a}^i - Q(s, a; \theta))^2]$
  - 8:    $Explore(\cdot)$ , update  $\mathcal{B}$
  - 9: **end for**
- output**  $Q_N^A(s, a) = \frac{1}{K} \sum_{k=0}^{K-1} Q(s, a; \theta_{N-k})$
- 

<sup>6</sup>Anschel O, Baram N, Shimkin N. Averaged-DQN: Variance Reduction and Stabilization for Deep Reinforcement Learning[C] International Conference on Machine Learning. 2017: 176-185.



# Averaged Deep Q-Network

## Performance



# Averaged Deep Q-Network

## Error Analysis

Let  $Q(s, a; \theta_i)$  be the value function of DQN at iteration  $i$ ,

$$\begin{aligned}
 \Delta_i &= Q(s, a; \theta_i) - Q^*(s, a) \\
 &= \underbrace{Q(s, a; \theta_i) - y_{s,a}^i}_{\text{Target Approximation Error}} + \underbrace{y_{s,a}^i - \hat{y}_{s,a}^i}_{\text{Overestimation Error}} \\
 &\quad + \underbrace{\hat{y}_{s,a}^i - Q^*(s, a)}_{\text{Optimality Difference}} .
 \end{aligned} \tag{14}$$

Here  $y_{s,a}^i$  is the *DQN target*, and  $\hat{y}_{s,a}^i$  is the *true target*, such that

$$\begin{aligned}
 y_{s,a}^i &= E_{\mathcal{B}} \left[ r + \gamma \max_{a'} Q(s', a'; \theta_{i-1}) | s, a \right] , \\
 \hat{y}_{s,a}^i &= E_{\mathcal{B}} \left[ r + \gamma \max_{a'} (\hat{y}_{s', a'}^{i-1}) | s, a \right] .
 \end{aligned} \tag{15}$$



# Averaged Deep Q-Network

## Background and Related Work

Define  $Z_{s,a}^i$  as TAE (Target Approximation Error) and  $R_{s,a}^i$  as overestimation error.

$$\begin{aligned} Z_{s,a}^i &= Q(s, a; \theta_i) - y_{s,a}^i, \\ R_{s,a}^i &= y_{s,a}^i - \hat{y}_{s,a}^i. \end{aligned} \quad (16)$$

In Thrun & Schwartz (1993),  $Z_{s,a}^i$  is considered as a random variable uniformly distributed error in  $[-\epsilon, \epsilon]$  and

$$E_{\mathbf{Z}}[R_{s,a}^i] = \gamma E_{\mathbf{Z}}[\max_{a'} [Z_{s',a'}^{i-1}]] = \gamma \epsilon \frac{n-1}{n+1}. \quad (17)$$

In Double Q-learning paper, the author replaces *positive* bias with a *negative* one.



# Averaged Deep Q-Network

## TAE Variance Reduction

Assume that

$$\begin{aligned} E[Z_{s,a}^i] &= 0, \quad \text{Var}[Z_{s,a}^i] = \sigma_s^2, \\ \text{for } i \neq j, \text{Cov}[Z_{s,a}^i, Z_{s',a'}^j] &= 0. \end{aligned} \tag{18}$$

We consider a fixed policy for updating the target values, and conveniently consider a zero reward  $r = 0$  everywhere since it has no effect on variance calculations.



# Averaged Deep Q-Network

## TAE Variance Reduction (cont.)

Consider M-state unidirectional MDP as



$$\begin{aligned}
 Q^{DQN}(s_0, a; \theta_i) &= Z_{s_0, a}^i + y_{s_0, a}^i \\
 &= Z_{s_0, a}^i + \gamma Q(s_1, a; \theta_{i-1}) \\
 &= Z_{s_0, a}^i + \gamma [Z_{s_1, a}^{i-1} + y_{s_1, a}^{i-1}] = \dots = \\
 &= Z_{s_0, a}^i + \gamma Z_{s_1, a}^{i-1} + \dots + \gamma^{M-1} Z_{s_{M-1}, a}^{i-(M-1)}
 \end{aligned} \tag{19}$$

Since for  $i \neq j$ ,  $\text{Cov}[Z_{s, a}^i, Z_{s', a'}^j] = 0$ , we have

$$\text{Var}[Q^{DQN}(s_0, a; \theta_i)] = \sum_{m=0}^{M-1} \gamma^{2m} \sigma_{s_m}^2 .$$

(20)

# Averaged Deep Q-Network

## TAE Variance Reduction (cont.)

For Averaged DQN,

$$Q_i = Z_i + \gamma P \frac{1}{K} \sum_{k=1}^K Q_{i-k} , \quad (21)$$

where  $P \in \mathbb{R}_+^{S \times S}$  is the transition probabilities matrix for the given policy.  
Recall that  $Z_{s,a}^i = Q(s, a; \theta_i) - y_{s,a}^i$ .



# Averaged Deep Q-Network

## Ensemble DQN

---

### Algorithm 3 Ensemble DQN

---

- 1: Initialize  $K$  Q-networks  $Q(s, a; \theta^k)$  with random weights  $\theta_0^k$  for  $k \in \{1, \dots, K\}$
  - 2: Initialize Experience Replay (ER) buffer  $\mathcal{B}$
  - 3: Initialize exploration procedure  $Explore(\cdot)$
  - 4: **for**  $i = 1, 2, \dots, N$  **do**
  - 5:    $Q_{i-1}^E(s, a) = \frac{1}{K} \sum_{k=1}^K Q(s, a; \theta_{i-1}^k)$
  - 6:    $y_{s,a}^i = \mathbb{E}_{\mathcal{B}} [r + \gamma \max_{a'} Q_{i-1}^E(s', a') | s, a]$
  - 7:   **for**  $k = 1, 2, \dots, K$  **do**
  - 8:      $\theta_i^k \approx \operatorname{argmin}_{\theta} \mathbb{E}_{\mathcal{B}} [(y_{s,a}^i - Q(s, a; \theta))^2]$
  - 9:   **end for**
  - 10:    $Explore(\cdot)$ , update  $\mathcal{B}$
  - 11: **end for**
- output**  $Q_N^E(s, a) = \frac{1}{K} \sum_{k=1}^K Q(s, a; \theta_i^k)$
- 





# Averaged Deep Q-Network

## Ensemble DQN Variance

For  $i > M$ ,

$$\begin{aligned} Q_i^E(s_0, a) &= \sum_{m=0}^{M-1} \gamma^m \frac{1}{K} \sum_{k=1}^K Z_{s_m, a}^{k, i-m} \\ \text{Var}[Q_i^E(s_0, a)] &= \sum_{m=0}^{M-1} \frac{1}{K} \gamma^{2m} \sigma_{s_m}^2 \\ &= \frac{1}{K} \text{Var}[Q^{DQN}(s_0, a; \theta_i)] \end{aligned} \tag{22}$$



# Averaged Deep Q-Network

## Averaged DQN Variance

For  $i > KM$ ,

$$\text{Var}[Q_i^A(s_0, a)] = \sum_{m=0}^{M-1} D_{K,m} \gamma^{2m} \sigma_{s_m}^2, \quad (23)$$

where  $D_{K,m} = \frac{1}{N} \sum_{n=0}^{N-1} |U_n/K|^{2(m+1)}$  and  $U = (U_n)_{n=0}^{N-1}$  denoting a Discrete Fourier Transform of a rectangle pulse.

Furthermore,  $D_{K,m} < \frac{1}{K}$  and

$$\begin{aligned} \text{Var}[Q_i^A(s_0, a)] &< \text{Var}[Q_i^E(s_0, a)] \\ &= \frac{1}{K} \text{Var}[Q^{DQN}(s_0, a; \theta_i)]. \end{aligned} \quad (24)$$



# Outline

- 1 Revision of Value-based RL
  - Dynamic Programming(omitted)
  - Monte-carlo Method(omitted)
  - TD: Sarsa and Q-learning
- 2 Deep Q-network
  - Nature DQN
  - Several Improvements
- 3 Issues in Q-learning
  - Overestimation
  - Double Q-learning
  - Averaged Q-learning
- 4 Convergence of Tabular TD
  - Sarsa
  - Q-learning (TBE)



# Convergence of Sarsa(0)

## Convergence of Random Iterative Process

### Lemma

*A random iterative process*

$$\Delta_{t+1}(x) = (1 - \alpha_t(x))\Delta_t(x) + \alpha_t(x)F_t(x), \quad x \in X, \quad t = 0, 1, 2, \dots \quad (25)$$

*converges to zero w.p.1 if the following properties hold:*

- 1. *the set of possible states  $X$  is finite.*
- 2.  *$0 \leq \alpha_t(x) \leq 1$ ,  $\sum_t \alpha_t(x) = \infty$ ,  $\sum_t \alpha_t^2(x) < \infty$  w.p.1, where the probability is over the learning rates  $\alpha_t$ .*
- 3.  *$\|E[F_t(\cdot)|P_t]\|_W \leq \kappa\|\Delta_t\|_W + c_t$ , where  $\kappa \in [0, 1)$  and  $c_t$  converges to zero w.p.1.*
- 4.  *$\text{Var}[F_t(x)] \leq K(1 + \|\Delta_t\|_W)^2$ , where  $K$  is some constant.*

*Here  $P_t$  is an increasing sequence of  $\sigma$ -fields that includes the past of the process. In particular we assume that  $\alpha_t, \Delta_t, F_{t-1} \in P_t$ . The notation  $\|\cdot\|_W$  refers to some (fixed) weighted maximum norm.*

# Convergence of Sarsa(0)

## Theorem

*In finite state-action MDPs, the  $Q_t$  values computed by the Sarsa(0) rule*

$$\begin{aligned} Q_{t+1}(s_t, a_t) &= Q_t(s_t, a_t) + \alpha_t(s_t, a_t)[r_t + \gamma Q_t(s_{t+1}, a_{t+1}) - Q_t(s_t, a_t)] \\ &= (1 - \alpha(s_t, a_t))Q_t(s_t, a_t) + \alpha_t(s_t, a_t)[r_t + \gamma Q_t(s_{t+1}, a_{t+1})] . \end{aligned}$$

*converges to  $Q^*$  and the learning policy  $\pi_t$  converges to an optimal policy  $\pi^*$  if the learning policy is GLIE with these additional conditions are satisfied*

- 1. *The  $Q$  values are stored in a lookup table.*
- 2. *The learning rates satisfy*  
 $0 \leq \alpha_t(s_t, a_t) \leq 1$ ,  $\sum_t \alpha_t(s_t, a_t) = \infty$ ,  $\sum_t \alpha_t^2(s_t, a_t) < \infty$  *and*  
 $\alpha_t(s_t, a_t) = 0$  *unless*  $(s, a) = (s_t, a_t)$ .
- 3.  $\text{Var}[r(s, a)] < \infty$  .

# Convergence of Sarsa(0)

- $x \stackrel{\text{def}}{=} (s_t, a_t)$ .
- $\Delta_t \stackrel{\text{def}}{=} Q_t(s, a) - Q^*(s, a)$ .

So we get

$$\begin{aligned} \Delta_{t+1}(s_t, a_t) &= Q_{t+1}(s_t, a_t) - Q^*(s, a) \\ &= (1 - \alpha(s_t, a_t))\Delta_t(s_t, a_t) + \alpha_t(s_t, a_t)F_t(s_t, a_t). \end{aligned} \quad (26)$$

where

$$\begin{aligned} F_t(s_t, a_t) &= r_t + \gamma \max_{a'} Q_t(s_{t+1}, a') - Q^*(s_t, a_t) \\ &+ \gamma \left[ Q_t(s_{t+1}, a_{t+1}) - \max_{a'} Q_t(s_{t+1}, a') \right] \\ &\stackrel{\text{def}}{=} r_t + \gamma \max_{a'} Q_t(s_{t+1}, a') - Q^*(s_t, a_t) + C_t(Q) \\ &\stackrel{\text{def}}{=} F_t^Q(s_t, a_t) + C_t(s_t, a_t) \end{aligned} \quad (27)$$



# Outline

- 1 Revision of Value-based RL
  - Dynamic Programming(omitted)
  - Monte-carlo Method(omitted)
  - TD: Sarsa and Q-learning
- 2 Deep Q-network
  - Nature DQN
  - Several Improvements
- 3 Issues in Q-learning
  - Overestimation
  - Double Q-learning
  - Averaged Q-learning
- 4 Convergence of Tabular TD
  - Sarsa
  - Q-learning (TBE)

