

# Deep Q-Learning

# Deep Q-Learning

- Q-Learning + Redes Neurais



# Deep Q-Learning

- Q-Learning + Redes Neurais
- Estado  $s$  como um input, valores  $\hat{q}(s, a)$  como outputs
- Como no Deep SARSA, a alvo é lidar com estados contínuos
- *Off Policy*



# Deep Q-Learning

## Algorithm 1 Deep Q-Learning

```
1: Input:  $\alpha$  learning rate,  $\epsilon$  random action probability,  
2:    $\gamma$  discount factor,  
3: Initialize q-value parameters  $\theta$  and target parameters  $\theta_{targ} \leftarrow \theta$   
4:  $b \leftarrow \epsilon$ -greedy policy w.r.t  $\hat{q}(s, a|\theta)$   
5:  $\pi \leftarrow$  greedy policy w.r.t  $\hat{q}(s, a|\theta)$   
6: Initialize replay buffer  $B$   
7: for episode  $\in 1..N$  do  
8:   Restart environment and observe the initial state  $S_0$   
9:   for  $t \in 0..T - 1$  do  
10:    Select action  $A_t \sim b(S_t)$   
11:    Execute action  $A_t$  and observe  $S_{t+1}, R_{t+1}$   
12:    Insert transition  $(S_t, A_t, R_{t+1}, S_{t+1})$  into the buffer  $B$   
13:     $K = (S, A, R, S') \sim B$   
14:    Select actions  $A' \sim \pi(S')$   
15:    Compute loss function over the batch of experiences:
```

$$L(\theta) = \frac{1}{|K|} \sum_{i=1}^{|K|} [R_i + \gamma \hat{q}(S'_i, A'_i|\theta_{targ}) - \hat{q}(S_i, A_i|\theta)]^2 \quad (1)$$

```
16:   end for  
17:   Every  $k$  episodes synchronize  $\theta_{targ} \leftarrow \theta$   
18: end for  
19: Output: Near optimal policy  $\pi$  and q-value approximations  $\hat{q}(s, a|\theta)$ 
```



# Deep Q-Learning

## Algorithm 1 Deep Q-Learning

```
1: Input:  $\alpha$  learning rate,  $\epsilon$  random action probability,  
2:    $\gamma$  discount factor,  
3: Initialize q-value parameters  $\theta$  and target parameters  $\theta_{targ} \leftarrow \theta$   
4:  $b \leftarrow \epsilon$ -greedy policy w.r.t  $\hat{q}(s, a|\theta)$   
5:  $\pi \leftarrow$  greedy policy w.r.t  $\hat{q}(s, a|\theta)$   
6: Initialize replay buffer  $B$   
7: for episode  $\in 1..N$  do  
8:   Restart environment and observe the initial state  $S_0$   
9:   for  $t \in 0..T - 1$  do  
10:    Select action  $A_t \sim b(S_t)$   
11:    Execute action  $A_t$  and observe  $S_{t+1}, R_{t+1}$   
12:    Insert transition  $(S_t, A_t, R_{t+1}, S_{t+1})$  into the buffer  $B$   
13:     $K = (S, A, R, S') \sim B$   
14:    Select actions  $A' \sim \pi(S')$   
15:    Compute loss function over the batch of experiences:
```

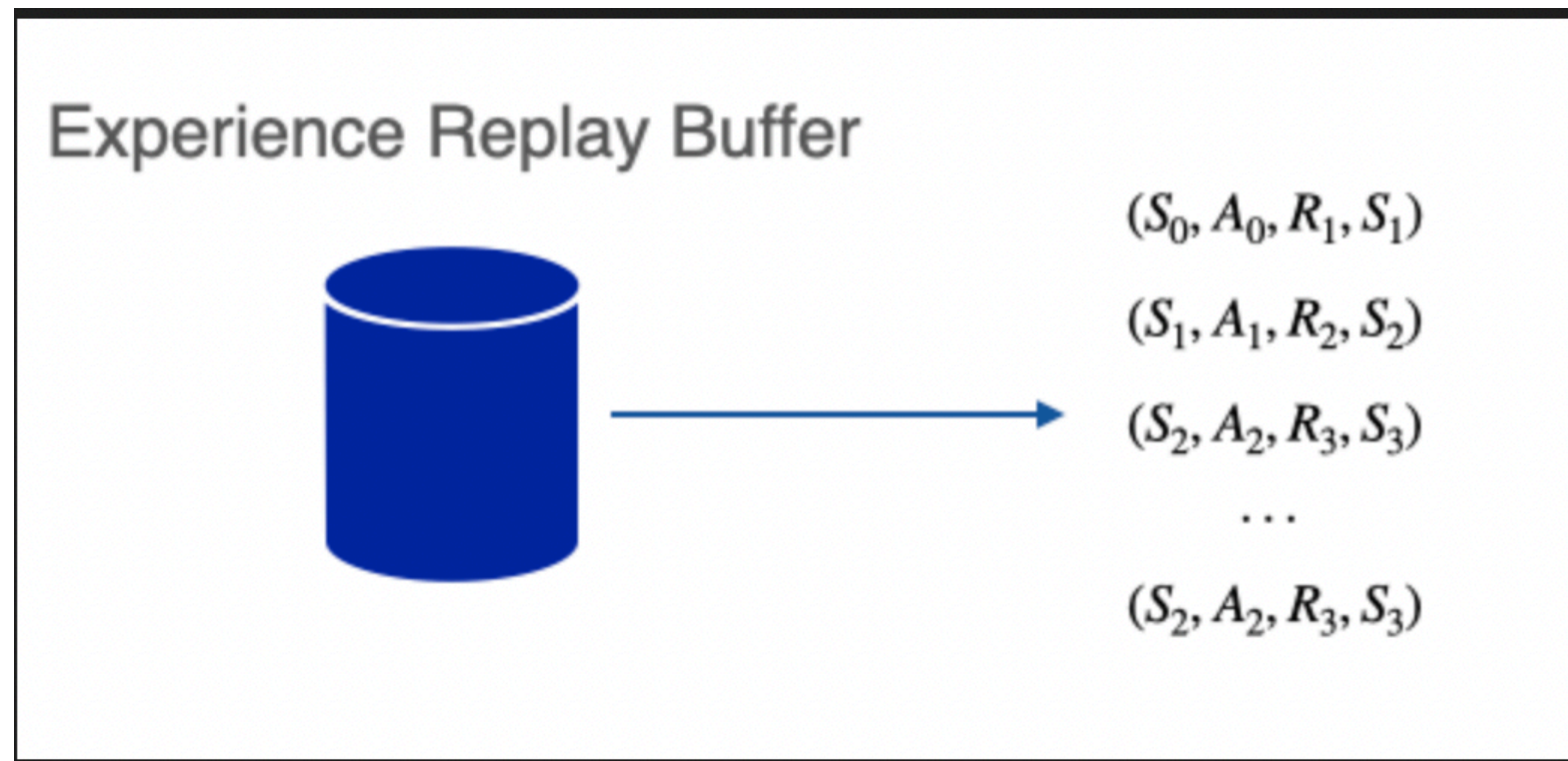
$$L(\theta) = \frac{1}{|K|} \sum_{i=1}^{|K|} [R_i + \gamma \hat{q}(S'_i, A'_i|\theta_{targ}) - \hat{q}(S_i, A_i|\theta)]^2 \quad (1)$$

```
16:   end for  
17:   Every  $k$  episodes synchronize  $\theta_{targ} \leftarrow \theta$   
18: end for  
19: Output: Near optimal policy  $\pi$  and q-value approximations  $\hat{q}(s, a|\theta)$ 
```



# Memória de replay

- Para atualizar a rede neural, escolhe-se aleatoriamente um lote de transições da memória de replay



# Rede Neural Alvo

- Fazemos uma copia da rede neural para calcular o alvo, essa rede não é atualizada com o gradiente.

$$\theta_{targ} \leftarrow \theta$$

$$L(\theta) = \frac{1}{|K|} \sum_{i=1}^{|K|} [R_i + \gamma \hat{q}(S'_i, A'_i | \theta_{targ}) - \hat{q}(S_i, A_i | \theta)]^2$$

$$\theta_{targ} \leftarrow \theta_k$$