

Deep SARSA

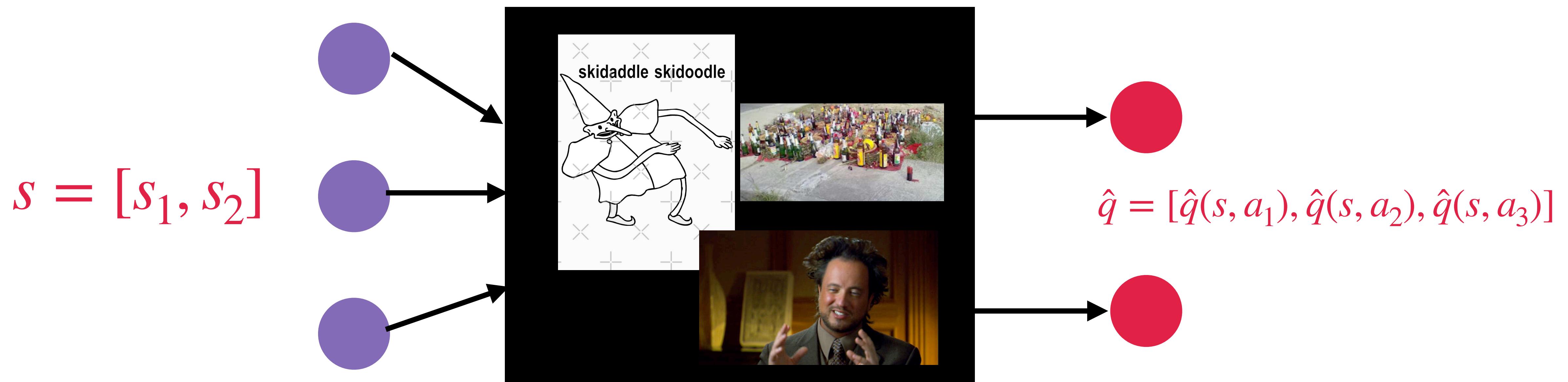
Introdução a redes Neurais

Informações sobre
O ambiente

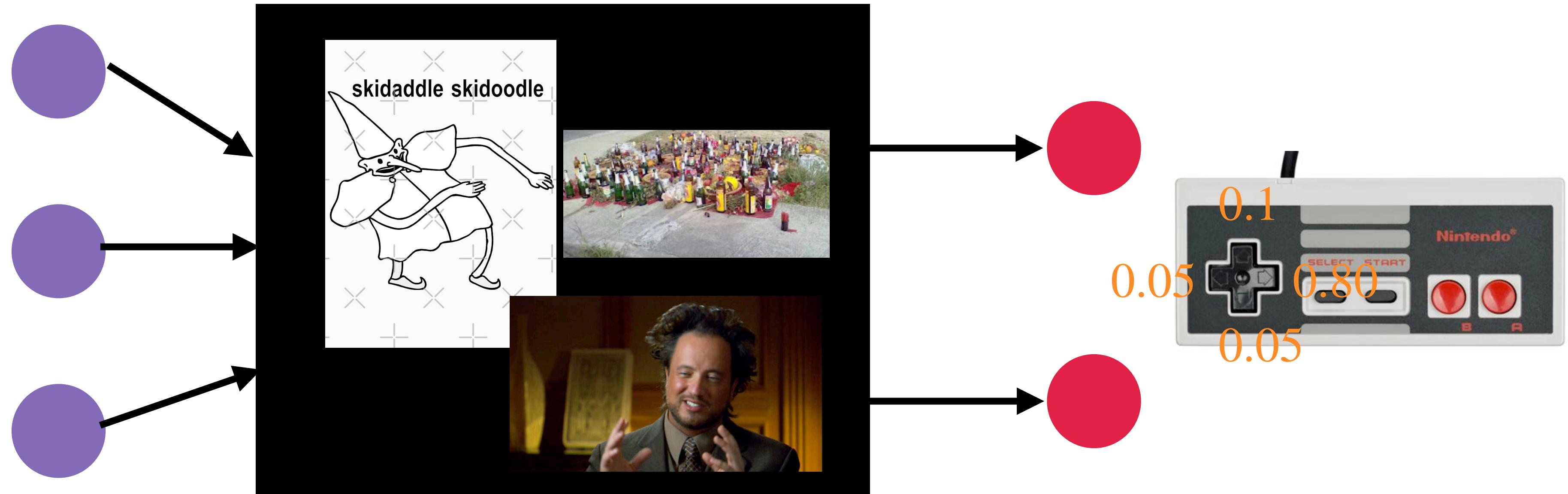


Probabilidade de
Ações

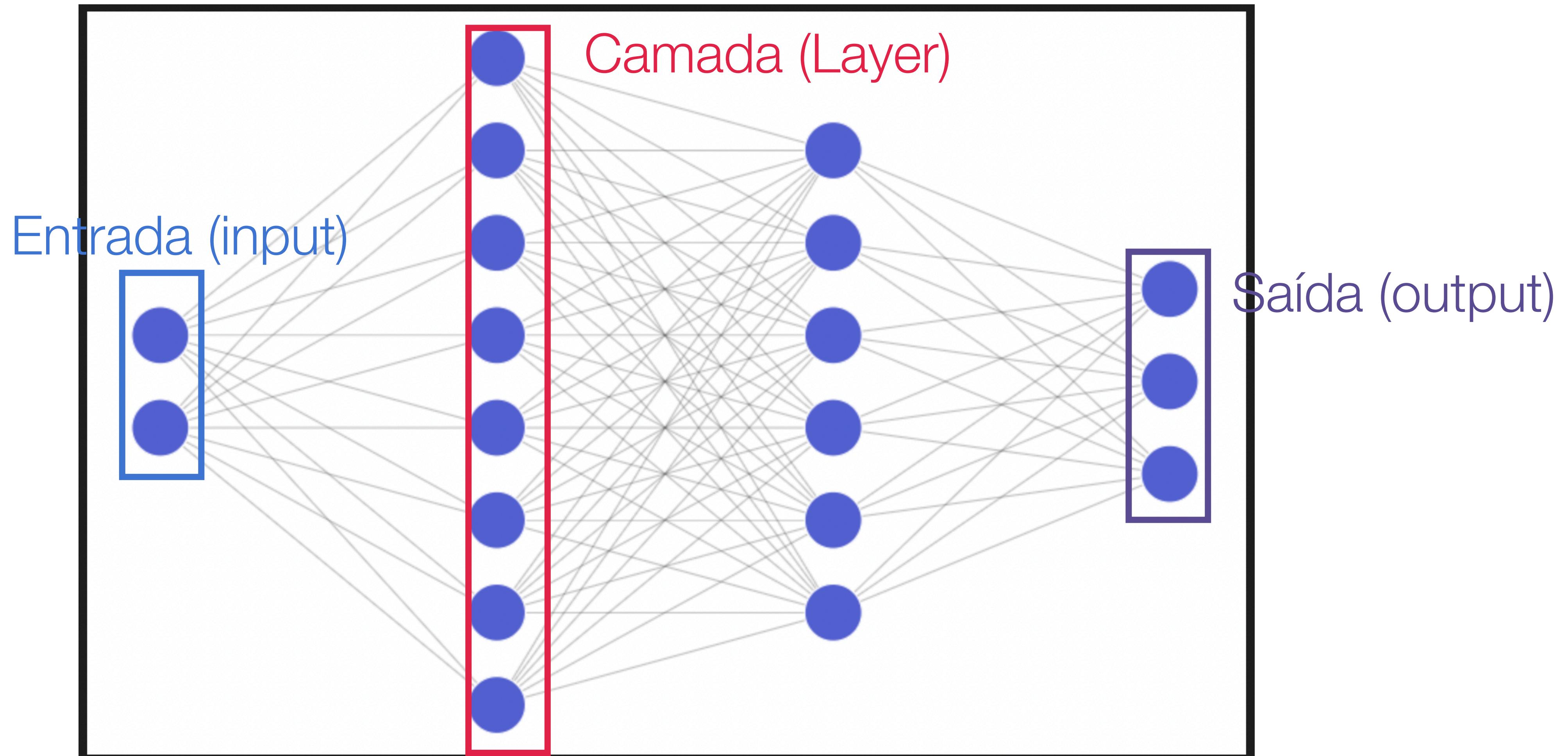
Introdução a redes Neurais



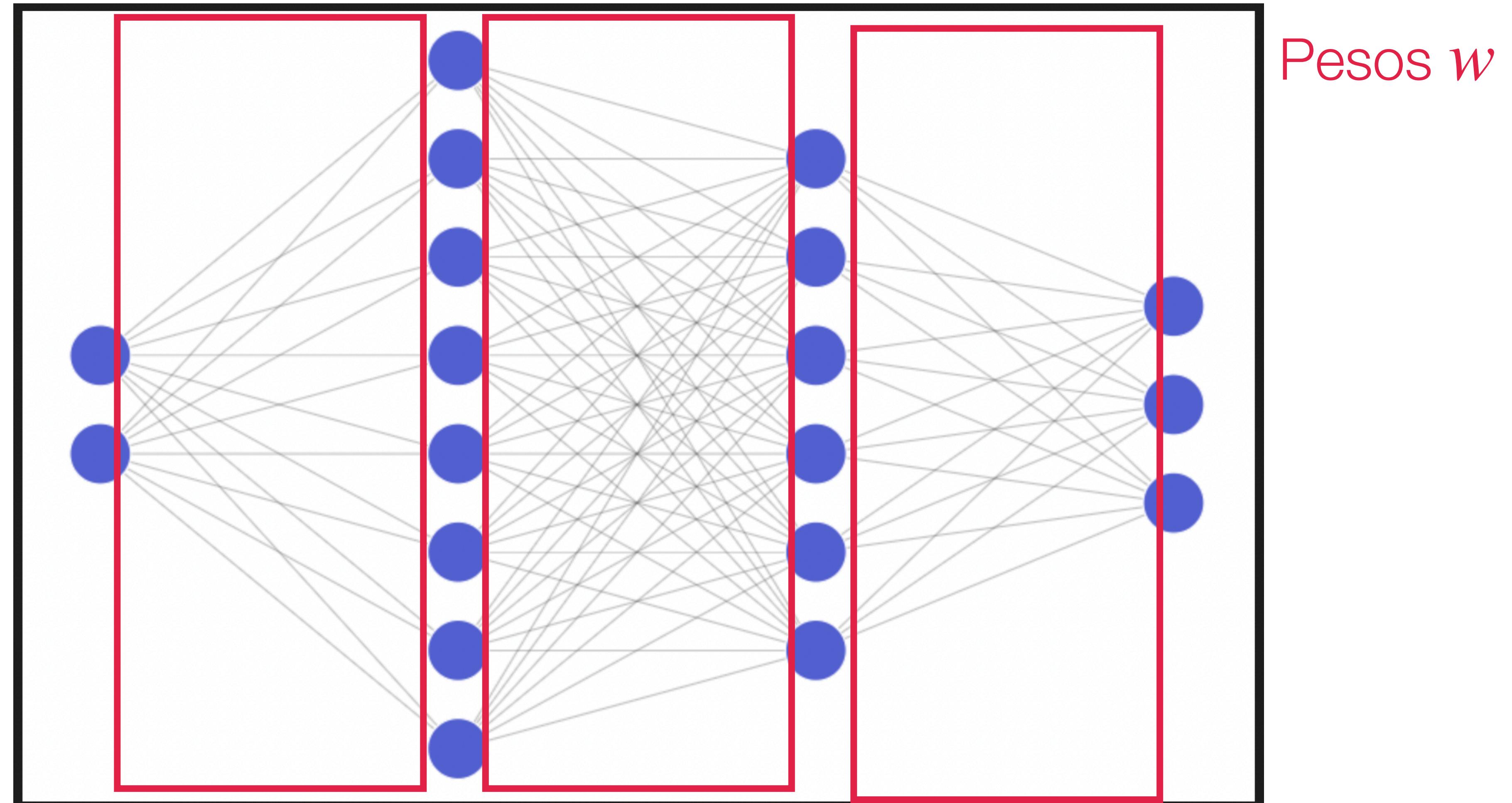
Introdução a redes Neurais



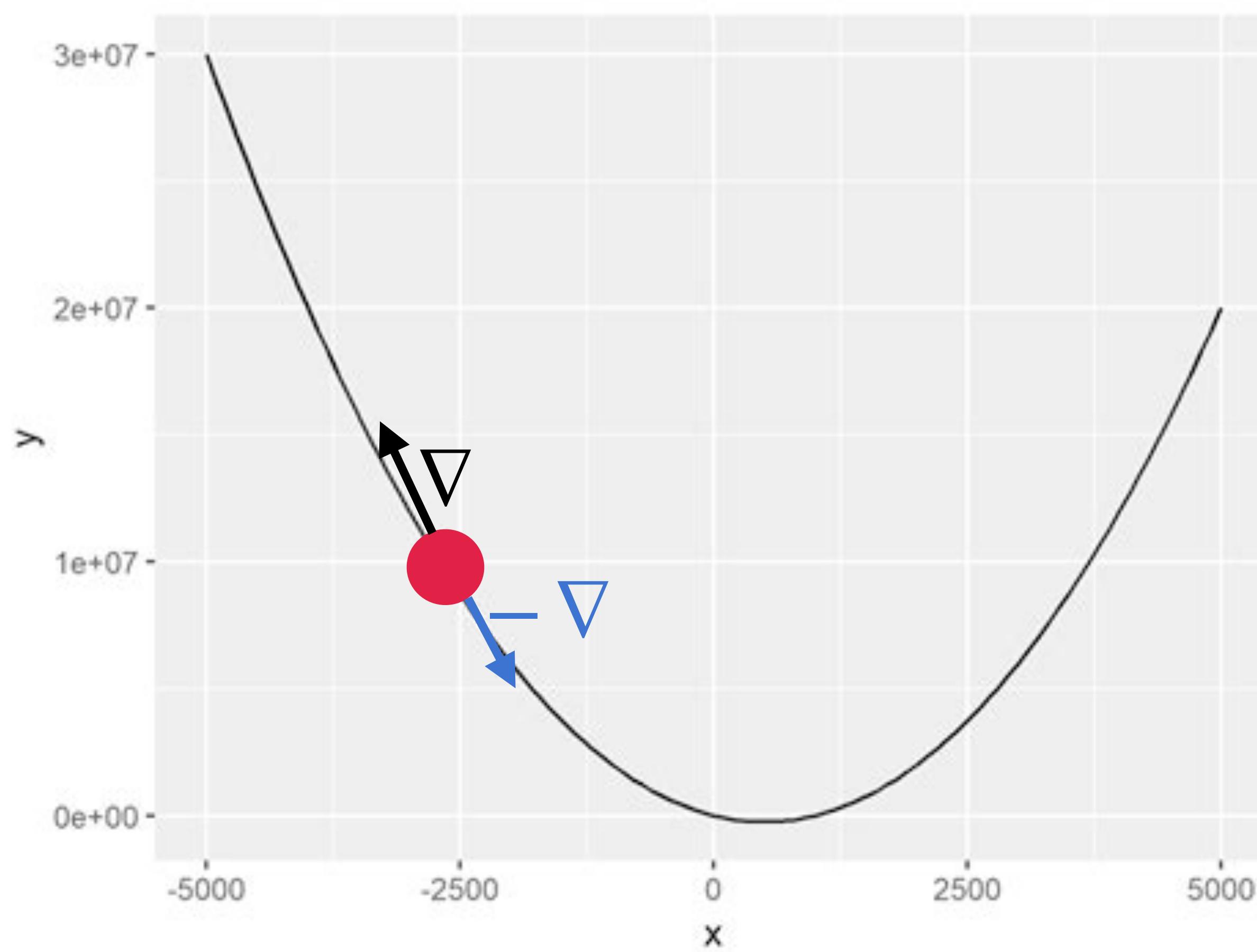
Introdução a redes Neurais



Introdução a redes Neurais



Cálculo do Gradiente



Deep SARSA

- SARSA + Redes neurais
- Como nos estados contínuos são infinito, só é necessário saber dos valores $q(s, a)$ deste momento.
- Aproxima os valores $\hat{q}(s, a)$ das ações disponíveis no estado atual s através da rede neural

Deep SARSA

Algorithm 1 Deep SARSA

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

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

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

Deep SARSA

Algorithm 1 Deep SARSA

```
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:  $\pi \leftarrow \epsilon$ -greedy policy w.r.t  $\hat{q}(s, a|\theta)$   
5: Initialize replay buffer  $B$   
6: for episode  $\in 1..N$  do  
7:   Restart environment and observe the initial state  $S_0$   
8:   for  $t \in 0..T - 1$  do  
9:     Select action  $A_t \sim \pi(S_t)$   
10:    Execute action  $A_t$  and observe  $S_{t+1}, R_{t+1}$   
11:    Insert transition  $(S_t, A_t, R_{t+1}, S_{t+1})$  into the buffer  $B$   
12:     $K = (S, A, R, S') \sim B$   
13:    Select actions  $A' \sim \pi(S')$   
14:    Compute loss function over the batch of experiences:  
          
$$L(\theta) = \frac{1}{|K|} \sum_{i=1}^{|K|} [R_i + \gamma \hat{q}(S'_i, A'_i | \theta_{targ}) - \hat{q}(S_i, A_i | \theta)]^2 \quad (1)$$
  
15:   end for  
16:   Every  $k$  episodes synchronize  $\theta_{targ} \leftarrow \theta$   
17: end for  
18: Output: Near optimal policy  $\pi$  and q-value approximations  $\hat{q}(s, a|\theta)$ 
```

Agente interage
com o ambiente

Deep SARSA

Algorithm 1 Deep SARSA

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

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

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

Atualização dos
Parâmetros da rede
Baseado nas experiencias

Processo de otimização

- Erro médio quadrado (MSE):

$$L(\theta) = \frac{1}{N} \sum_{i=1}^N [y_i - \hat{y}_i]^2$$

Minimização dos erros das estimativas

Processo de otimização

- No nosso caso:

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

Batch

Objetivo

Estimativa

$$y_i = R_i + \gamma \hat{q}(S'_i, A'_i | \theta_{targ})$$

Valor que queremos levar
As estimativas

$$y_i = \hat{q}(S_i, A_i | \theta)$$

Estimativa do valor $q(s, a)$

Processo de otimização

- Calculo do gradiente da função L em respeito aos pesos w :

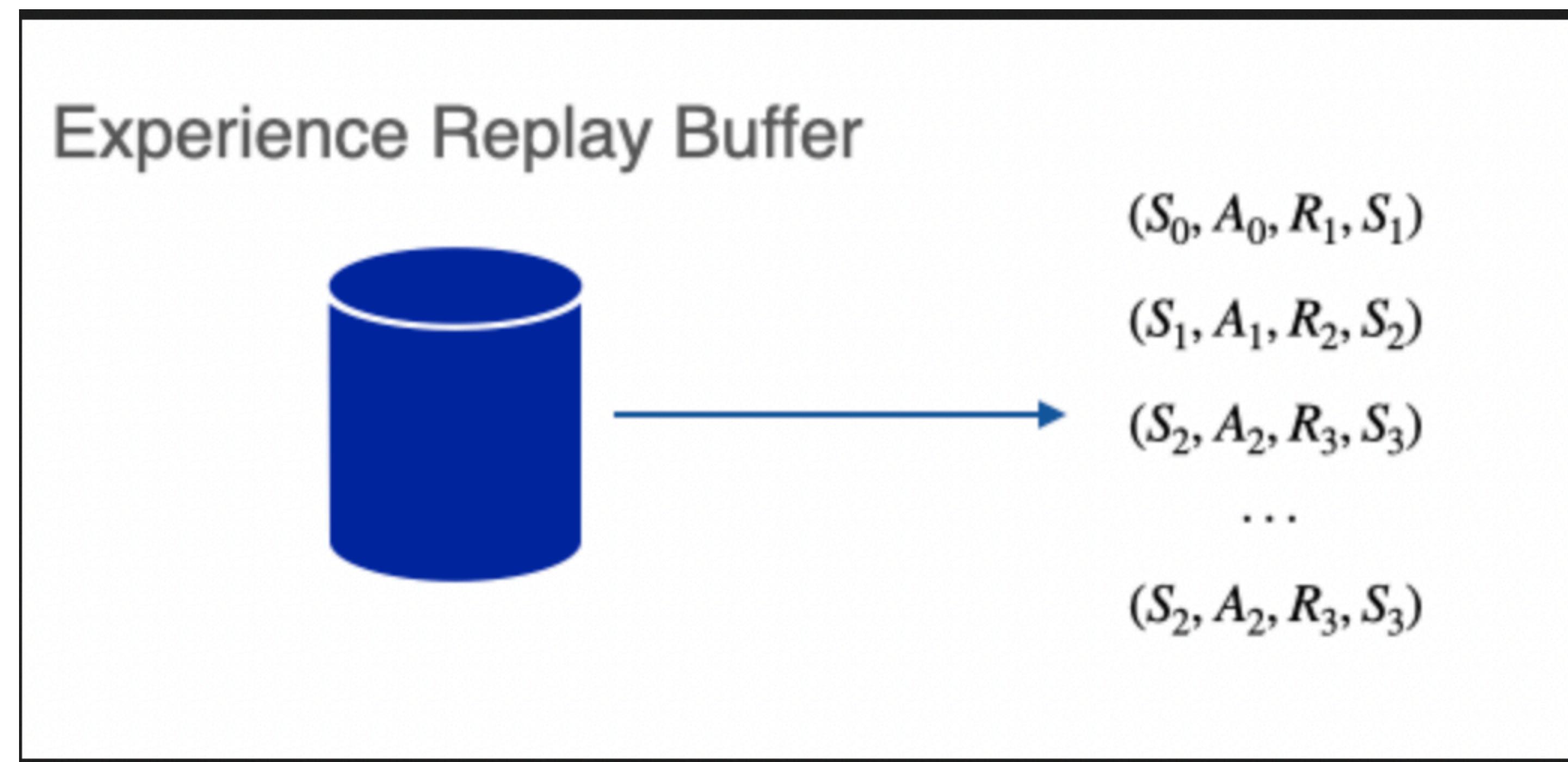
$$\nabla L(\theta) = \left[\frac{\partial L}{\partial w_1}, \frac{\partial L}{\partial w_2}, \dots, \frac{\partial L}{\partial w_n} \right]$$

Com o SGD atualiza-se por um pequeno passo

$$\theta \leftarrow \theta - \alpha \nabla L(\theta)$$

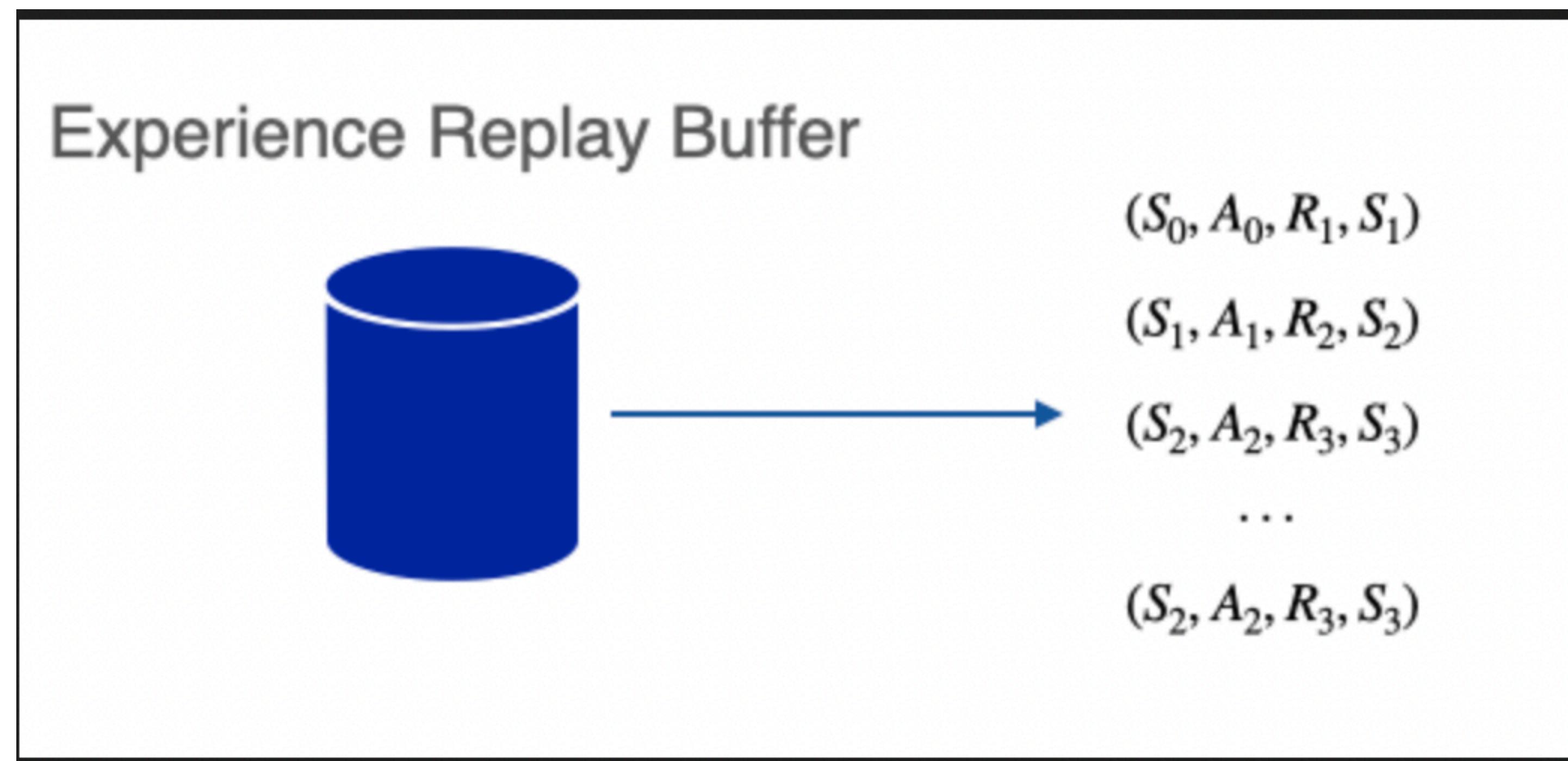
Memória de replay

- Memória que guarda as transições de estados que o agente observa ao tomar as ações



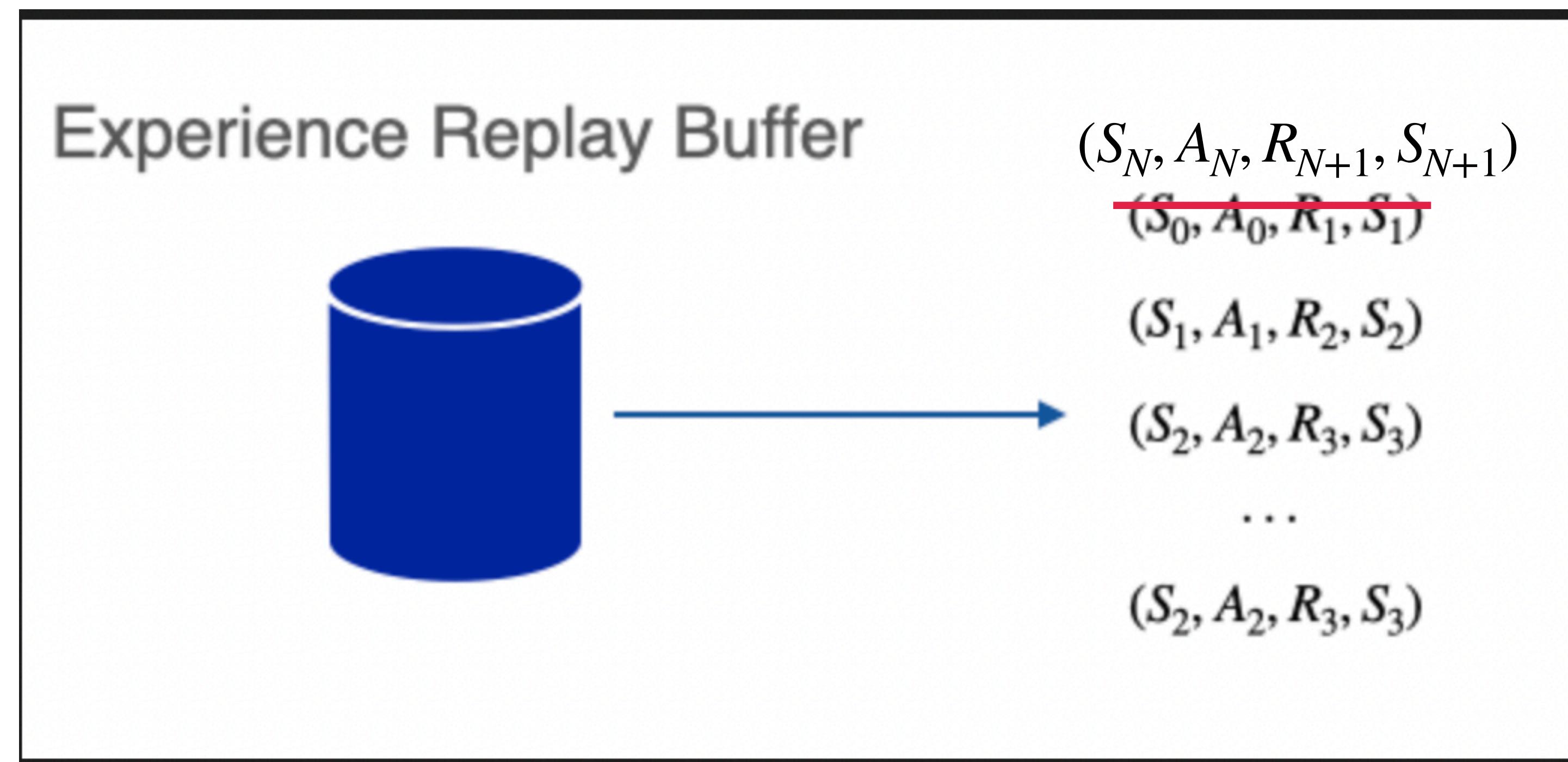
Memória de replay

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



Memória de replay

- A memória possui um espaço limitado e quando enche, os valores mais antigos são trocados pelos mais novos



Rede Neural Alvo

- Processo de aprendizado instável
- Quando um valor $q(s, a)$ muda, todos os valores serão afetados
- Quando se muda $\hat{q}(S_i, A_i | \theta)$, também modifica-se o valor alvo de $\hat{q}(S'_i, A'_i | \theta_{targ})$
- O alvo (Target) também precisa ser **estável!**

Rede Neural Alvo

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

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

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

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

Deep SARSA

Algorithm 1 Deep SARSA

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

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

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

Deep SARSA

**Abra o Notebook
Cap_9_Deep_SARSA.ipynb**