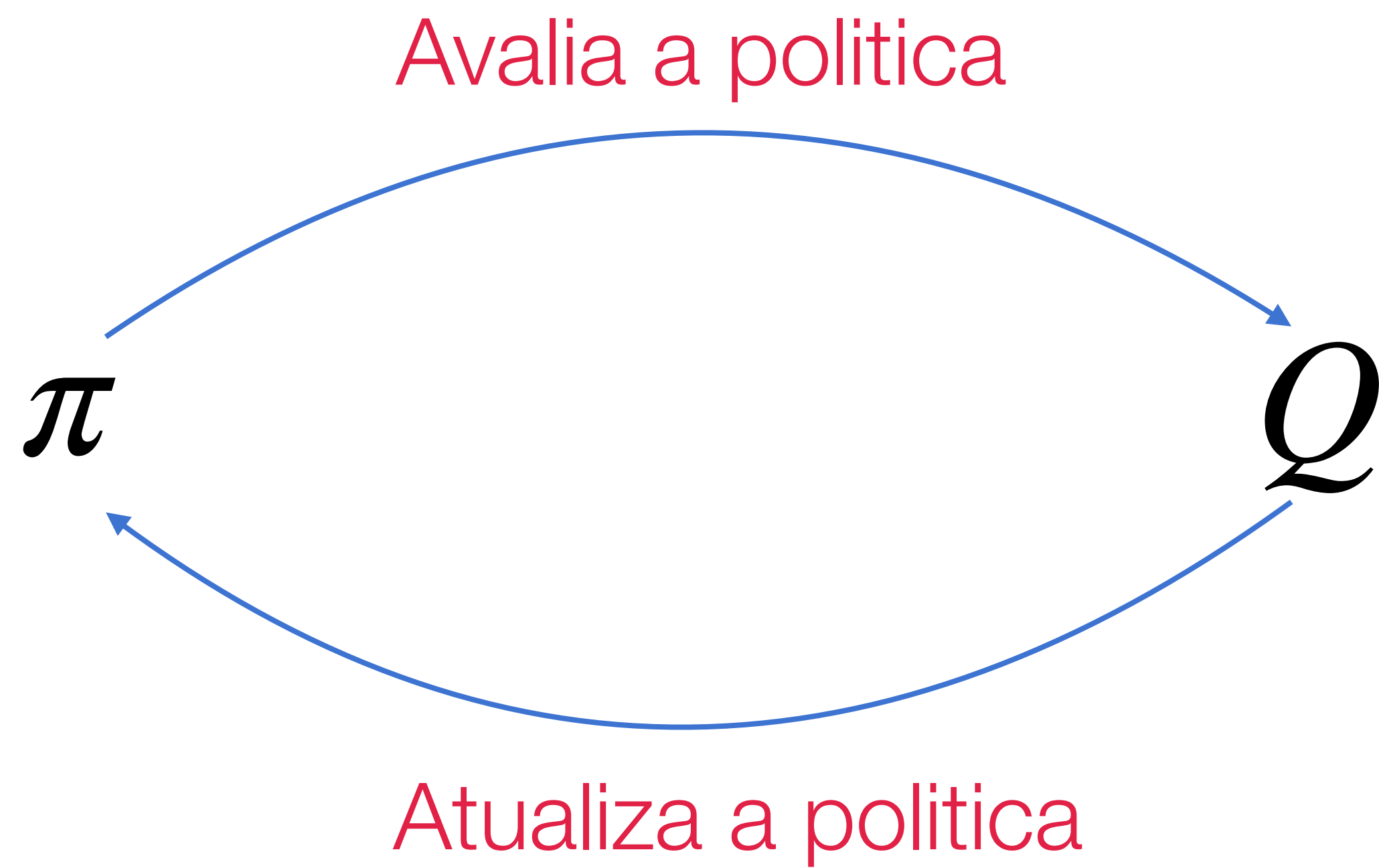


# Métodos de Diferença temporal (TD)

# O que são

- Base para a maioria dos métodos avançados de RL:
  - Tentativa e erro do MMC: aprende por experiência
  - Amostragem da PD: estimativas de valores  $Q(s, a)$
  - Memória

O que são



# O que são

## *Comparação*

- Métodos de Monte-Carlo esperam até o fim de um episódio para calcular  $G_t$  e atualizar  $Q(s, a)$

$$G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots + \gamma^{T-t-1} R_T$$

- Métodos de TD atualizam  $Q(s, a)$  toda vez que um agente realiza uma ação

# Alicerce de um Método de TD

- Como nos outros métodos, uma tabela de valores  $q(s, a)$  é construída para cada par  $S_t, A_t$

$$q_{\pi}(s, a) = \mathbb{E}[G_t | S_t = s, A_t = a] \quad q_{\pi}(s, a) = \sum_{s', r} p(s', r | s, a) [r + \gamma v_{\pi}(s')]$$

$$q_{\pi}(s, a) = \sum_{s', r} p(s', r | s, a) [r + \gamma \sum_{a'} \pi(a' | s') q_{\pi}(s', a')]$$

# Alicerce de um Método de TD

$$q_{\pi}(s, a) = \sum_{s', r} \underbrace{p(s', r | s, a)}_{\substack{\text{red arrow to } R_{t+1} \\ \text{red arrow to } S_{t+1}}} [r + \gamma \underbrace{\sum_{a'} \pi(a' | s') \underbrace{q_{\pi}(s', a')}_{\substack{\text{blue arrow to } Q}}]}_{\substack{\text{red arrow to } A_{t+1}}}]$$

# Alicerce de um Método de TD

- $Q_{\pi}(S_t, A_t)$  podem ser estimados por:

$$R_{t+1} + \gamma Q(S_{t+1}, A_{t+1})$$

$$R_{t+1} + \gamma Q(S_{t+1}, A_{t+1}) - Q(S_t, A_t)$$



Erro de Diferença temporal

# Alicerce de um Método de TD

- Diferentes métodos possuem diferentes estimativas de  $Q(S_t, A_t)$

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha [R_{t+1} + \gamma Q(S_{t+1}, A_{t+1}) - Q(S_t, A_t)]$$

Lembre-se do Método de Monte-Carlo com  $\alpha$  constante



# Alicerce de um Método de TD

- Reorganizando:

$$Q(S_t, A_t) \leftarrow (1 - \alpha)Q(S_t, A_t) + \alpha[R_{t+1} + \gamma Q(S_{t+1}, A_{t+1}) - Q(S_t, A_t)]$$

# SARSA

*Primeiro algoritmo de TD*

- SARSA:  $S_t, A_t, R_t, S_{t+1}, A_{t+1}$

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha[R_{t+1} + \gamma Q(S_{t+1}, A_{t+1}) - Q(S_t, A_t)]$$

- O SARSA é *on-policy*, só tem uma política
- Como no MMC, usa a política  $\epsilon$ -gulosa para exploração

# SARSA

*Primeiro algoritmo de TD*

---

**Algorithm 1** SARSA

---

```
1: Input:  $\alpha$  learning rate,  $\epsilon$  random action probability,  $\gamma$  discount factor
2:  $\pi \leftarrow \epsilon$ -greedy policy w.r.t  $Q(s, a)$ 
3: Initialize  $Q(s, a)$  arbitrarily, with  $Q(\text{terminal}, \cdot) = 0$ 
4: for episode  $\in 1..N$  do
5:   Reset the environment and observe  $S_0$ 
6:    $A_0 \sim \pi(S_0)$ 
7:   for  $t \in 0..T - 1$  do
8:     Execute  $A_t$  in the environment and observe  $S_{t+1}, R_{t+1}$ 
9:      $A_{t+1} \sim \pi(S_{t+1})$ 
10:     $Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha [R_{t+1} + \gamma Q(S_{t+1}, A_{t+1}) - Q(S_t, A_t)]$ 
11:  end for
12: end for
13: Output: Near optimal policy  $\pi$  and action values  $Q(s, a)$ 
```

---

# SARSA

*Primeiro algoritmo de TD*

**Abra o Notebook SARSA.ipynb**

# Q-Learning

*Segundo Algoritmo de TD*

- O Q-Learning é *off-policy*, possui uma politica exploratória  $\pi_b$
- Funciona exatamente como o SARSA, mas usa  $\pi_b$  para gerar uma trajetória

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha[R_{t+1} + \gamma Q(S_{t+1}, A_{t+1}) - Q(S_t, A_t)]$$

$$A_{t+1} = \pi(S_{t+1}) = \operatorname{argmax}_a Q(S_{t+1}, a)$$



# Q-Learning

*Segundo Algoritmo de TD*

---

**Algorithm 2** Q-Learning

---

```
1: Input:  $\alpha$  learning rate,  $\gamma$  discount factor
2:  $\pi \leftarrow$  greedy policy w.r.t  $Q(s, a)$ 
3:  $b \leftarrow$  exploratory policy with coverage of  $\pi$ 
4: Initialize  $Q(s, a)$  arbitrarily, with  $Q(\text{terminal}, \cdot) = 0$ 
5: for episode  $\in 1..N$  do
6:   Reset the environment and observe  $S_0$ 
7:   for  $t \in 0..T - 1$  do
8:      $A_t \sim b(S_t)$ 
9:     Execute  $A_t$  in the environment and observe  $S_{t+1}, R_{t+1}$ 
10:     $Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha [R_{t+1} + \gamma \pi(S_{t+1}) - Q(S_t, A_t)]$ 
11:   end for
12: end for
13: Output: Approximately optimal policy  $\pi$  and action values  $Q(s, a)$ 
```

---

# Q-Learning

*Segundo algoritmo de TD*

**Abra o Notebook QLearning.ipynb**