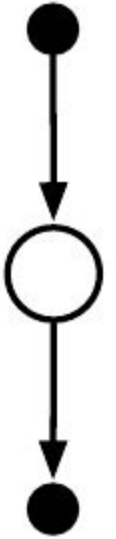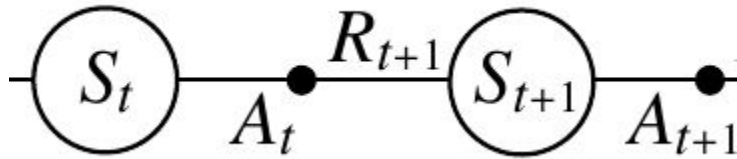# Aprendizaje Reforzado
# Clase 6

Julián Martínez

# Sarsa: TD con on-policy control



Sarsa

$$\hat{q}^{t+1}(S_t, A_t) = \hat{q}^t(S_t, A_t) + \alpha \left[ R_{t+1} + \gamma \hat{q}^t(S_{t+1}, A_{t+1}) - \hat{q}^t(S_t, A_t) \right]$$

$$A_{t+1} \sim \Pi_{\hat{q}^t} - \varepsilon \; \text{greedy}$$

# Pseudo-código del sutton

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

Algorithm parameters: step size $\alpha \in (0, 1]$, small $\varepsilon > 0$
Initialize $Q(s, a)$, for all $s \in \mathcal{S}^+, a \in \mathcal{A}(s)$, arbitrarily except that $Q(terminal, \cdot) = 0$

Loop for each episode:
    Initialize $S$
    Choose $A$ from $S$ using policy derived from $Q$ (e.g., $\varepsilon$-greedy)
    Loop for each step of episode:
        Take action $A$, observe $R, S'$
        Choose $A'$ from $S'$ using policy derived from $Q$ (e.g., $\varepsilon$-greedy)
        $Q(S, A) \leftarrow Q(S, A) + \alpha\big[R + \gamma Q(S', A') - Q(S, A)\big]$
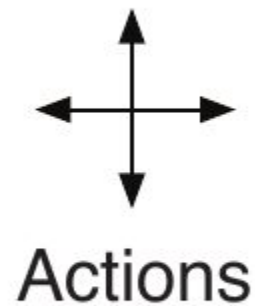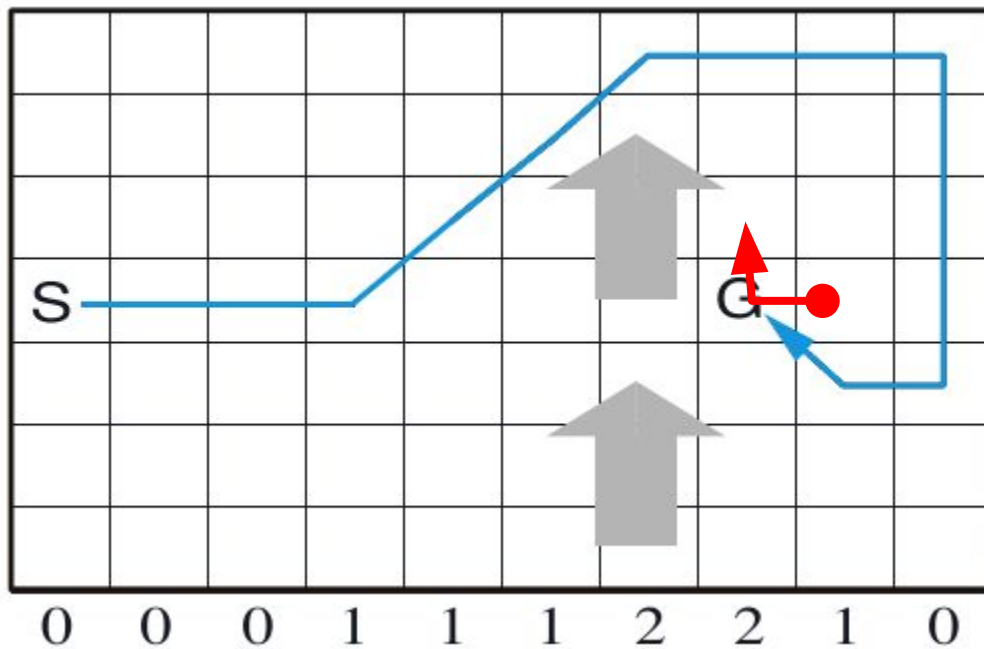        $S \leftarrow S'; A \leftarrow A';$
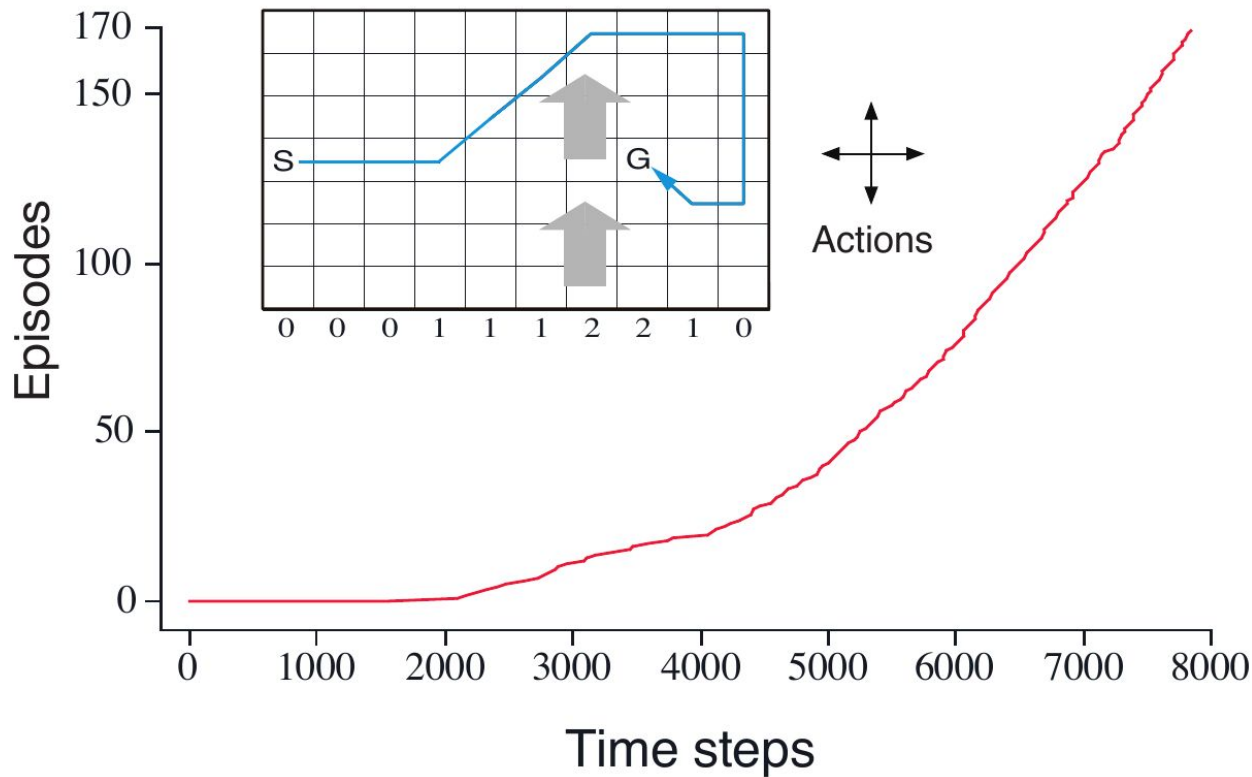    until $S$ is terminal

# Windy GridWorld

El reward es -1
hasta llegar a G

Intensidad del viento



Actions

0 0 0 1 1 1 2 2 1 0

# $\varepsilon$ – GREEDY SARSA, $\varepsilon = 0.1$, $\alpha = 0.5$

# Q-learning: TD con off-policy control

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha\left[R_{t+1} + \gamma \max_a Q(S_{t+1}, a) - Q(S_t, A_t)\right].$$

# Q-LEARNING: TD CON OFF-POLICY CONTROL

$$\hat{q}^{t+1}(S_t, A_t) = \hat{q}^t(S_t, A_t) + \alpha \left[ R_{t+1} + \gamma \, \hat{q}^t(S_{t+1}, A_{t+1}^*) - \hat{q}^t(S_t, A_t) \right]$$

$$A_{t+1}^* = \Pi_{\hat{q}^t}^*(S_{t+1}) \quad , \quad A_{t+1} \sim \Pi_{\hat{q}^t}^{*, \varepsilon}(\cdot \mid S_{t+1})$$

# Pseudo-código del sutton

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

Algorithm parameters: step size $\alpha \in (0, 1]$, small $\varepsilon > 0$
Initialize $Q(s, a)$, for all $s \in \mathcal{S}^+, a \in \mathcal{A}(s)$, arbitrarily except that $Q(terminal, \cdot) = 0$

Loop for each episode:
    Initialize $S$
    Loop for each step of episode:
        Choose $A$ from $S$ using policy derived from $Q$ (e.g., $\varepsilon$-greedy)
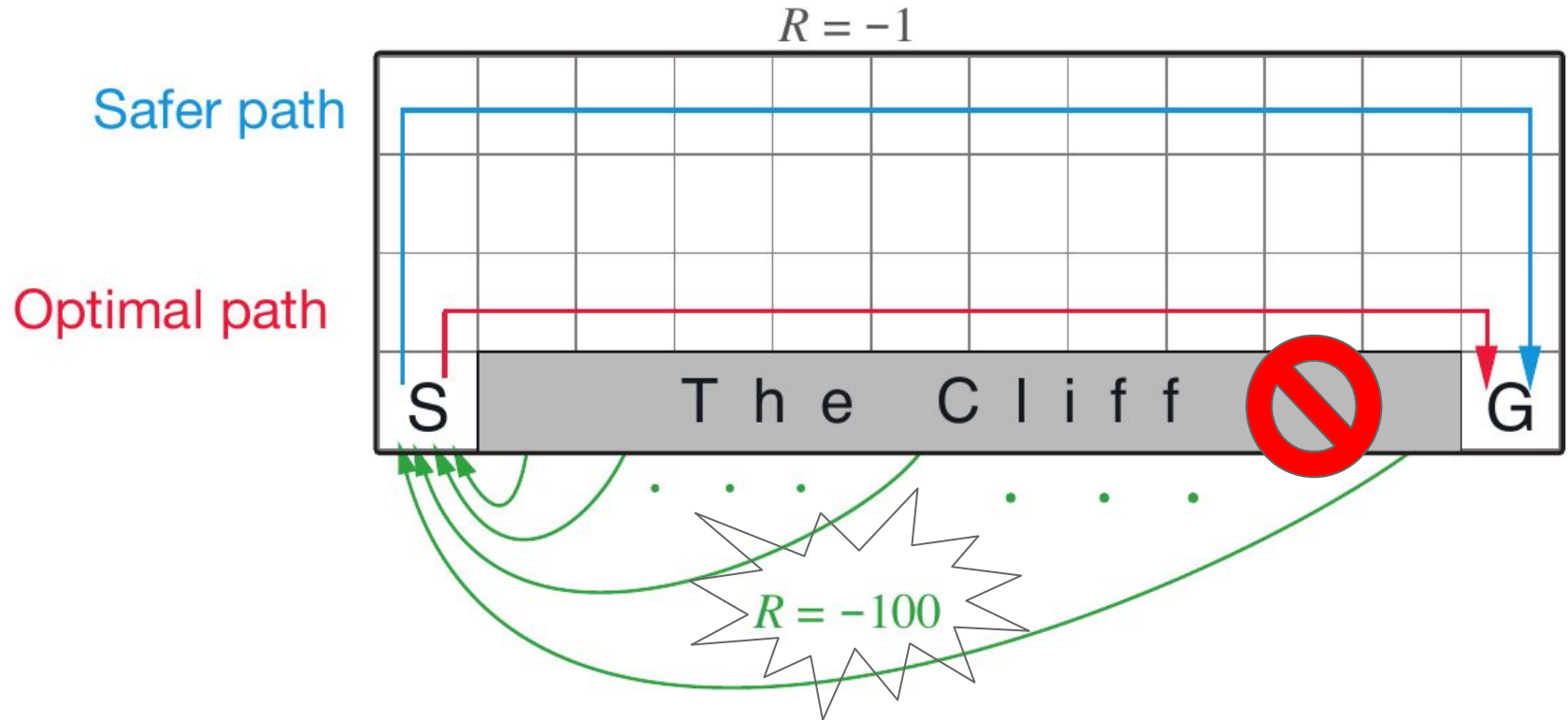        Take action $A$, observe $R$, $S'$
        $Q(S, A) \leftarrow Q(S, A) + \alpha \big[ R + \gamma \max_a Q(S', a) - Q(S, A) \big]$
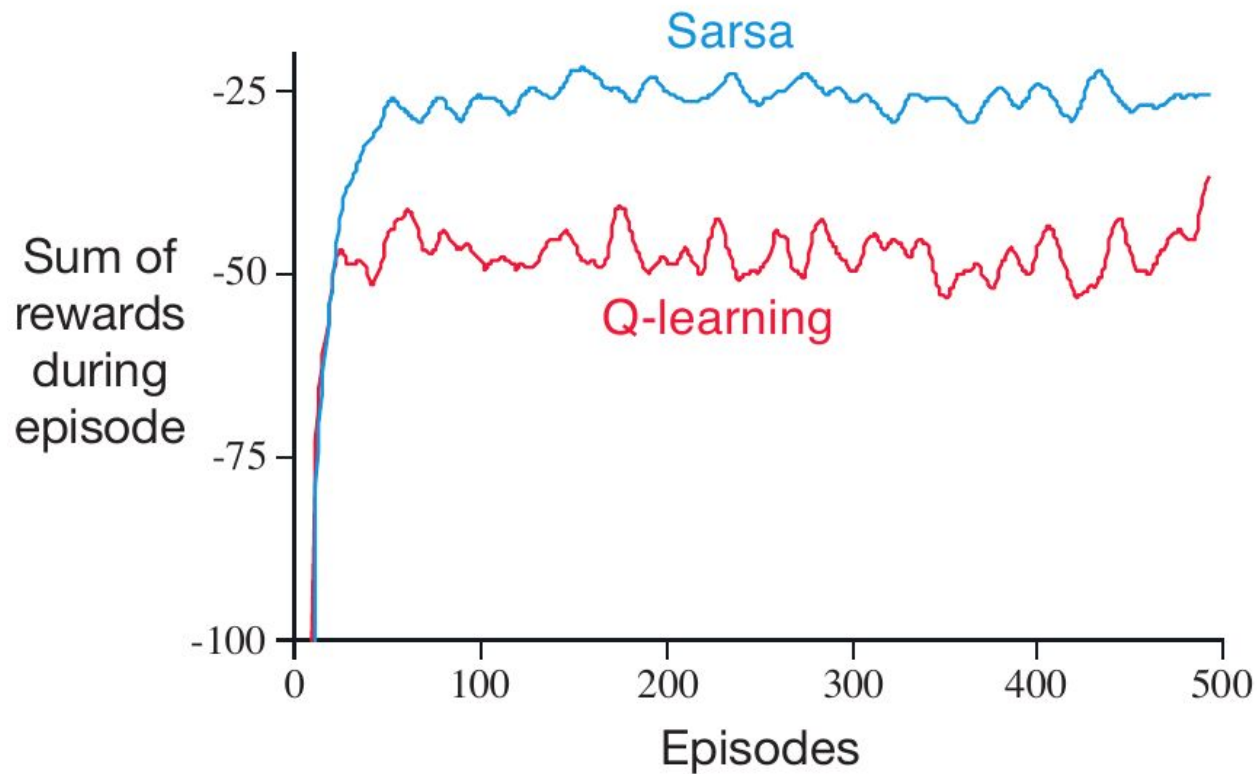        $S \leftarrow S'$
    until $S$ is terminal

# Ejemplo 6.6 - Sutton (Acantilado Walking)

# Sarsa vs Q-learning

# ¿Por qué es off-Policy?

$$q_*(s,a) = \mathbb{E}_{\pi_*^!}\left[G_t \mid S_{t=s}, A_{t=a}\right]$$

$$= \mathbb{E}_{\pi_*}\left[R_{t+1} + \gamma\, q_{\pi_*}(S_{t+1}, A_{t+1}) \mid S_{t=s}, A_{t=a}\right]$$

$$= \mathbb{E}_{\pi_*^\varepsilon}\left[\{R_{t+1} + \gamma\, q_{\pi_*}(S_{t+1}, A_{t+1})\}\, \omega_{t+1:T} \mid S_{t=s}, A_{t=a}\right]$$

$$\omega(s', a') = \frac{1\{\pi^*(s') = a'\}}{(1-\varepsilon)\, 1\{\pi^*(s') = a'\} + \varepsilon \cdot \frac{1}{|A|}}$$

# Resumiendo

$$\hat{q}^{t+1}(S_t, A_t) = \hat{q}^t(S_t, A_t) + \alpha \left[ R_{t+1} + \gamma \, \hat{q}^t(S_{t+1}, A^*_{t+1}) - \hat{q}^t(S_t, A_t) \right]$$

$$A^*_{t+1} = \Pi^*_{\hat{q}^t}(S_{t+1}) \qquad ; \qquad A_{t+1} \sim \Pi^{*,\varepsilon}_{\hat{q}^t}(\cdot \mid S_{t+1})$$

$$= \mathbb{E}_{\Pi^\varepsilon_*} \left[ \left\{ R_{t+1} + \gamma \, q_{\pi_*}(S_{t+1}, A_{t+1}) \right\} \omega(S_{t+1}, A_{t+1}) \mid S_{t=s}, A_{t=a} \right]$$

$$\omega(s', a') = \frac{1_{\{\pi^*(s') = a'\}}}{(1-\varepsilon) \, 1_{\{\pi^*(s') = a'\}} + \varepsilon \cdot \frac{1}{|A|}}$$

# Stochastic Approximation

$$E\left[f(\theta, W)\right]\big|_{\theta=\theta^*} = 0$$

Robbins–Monro, 1951

$$\theta(n+1) = \theta(n) + \alpha_n \, f\big(\theta(n), W(n+1)\big)$$

$$\sum \alpha_n = \infty \qquad \sum \alpha_n^2 < \infty \qquad \alpha_n = \frac{1}{n}$$

# ODES / Método de euler

$$\bar{F}(\theta) = E[f(\theta, w)]$$

$$\frac{d\,x(t)}{dt} = \bar{F}(x(t))$$

$$\frac{\theta(n+1) - \theta(n)}{1/n} = \bar{F}(\theta(n))$$

$$\theta(n+1) = \theta(n) + \alpha_n \left[ \overline{f}(\theta(n)) + \Delta(n+1) \right]$$

$$\Delta(n+1) = f(\theta(n), W(n+1)) - \overline{f}(\theta(n))$$

__THM__   If the ODE has a unique asymptotically stable equilibrium $x^*$ $\Rightarrow$ $x_n \longrightarrow x^*$ with probability one

$$P\left( \lim_{n \to \infty} x_n = x^* \right) = 1$$

# Ejemplo 1 - Monte Carlo

$$E\left[c(X)\right] \stackrel{?}{=}$$

$$f(\theta, x) = c(x) - \theta$$

$$\theta_{n+1} = \theta_n + \alpha_n \left[c(X_{n+1}) - \theta_n\right]$$

# Ejemplo 2 – Ajuste

$$F(\theta) := \frac{1}{2} E\left[(Y - f_\theta(X))^2\right]$$

$$\nabla_\theta F(\theta) = E\left[(Y - f_\theta(X)) \cdot \nabla_\theta f_\theta(X)\right]$$
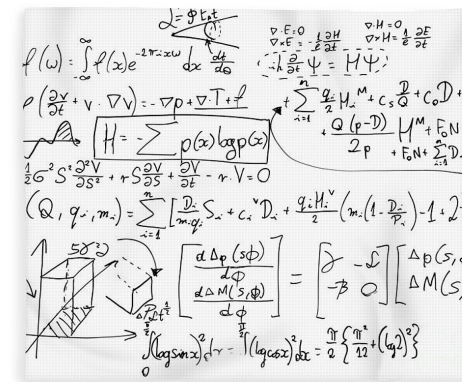
$$\theta_{n+1} = \theta_n + \alpha_n \left[(Y_n - f_{\theta_n}(X_n)) \nabla_\theta f_{\theta_n}(X_n)\right]$$

# Ejemplo 3 - Temporal Difference

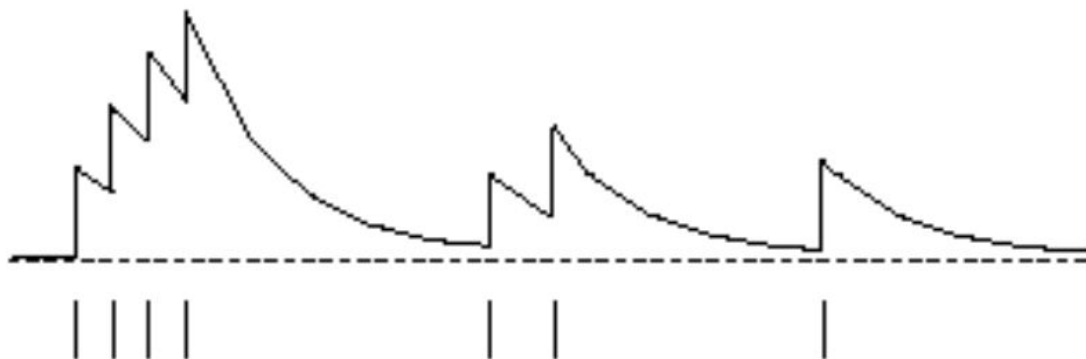$$E\left[R_s^+ + \gamma \sum_{s'} v(s') \mathbb{1}_{\{S_s^+ = s'\}} - v(s)\right] = 0 \quad \forall s$$

$$V_{n+1} = V_n + \alpha \left[R_{n+1} + \gamma V_n(S_{n+1}) - V_n(S_n)\right]$$

# Eligibility traces

$$E_0(s) = 0$$
$$E_t(s) = \gamma\lambda E_{t-1}(s) + \mathbf{1}(S_t = s)$$

accumulating eligibility trace

times of visits to a state

# El lado oscuro...



- Artículo de divulgación sobre Stochastic Approximation
  https://www.ias.ac.in/article/fulltext/reso/018/12/1086-1094
- THE O.D.E. METHOD FOR CONVERGENCE OF STOCHASTICAPPROXIMATION AND REINFORCEMENT LEARNING
  http://repository.ias.ac.in/5333/1/351.pdf
- Reinforcement Learning: Hidden Theory and New Super-Fast Algorithms (Charla de Meyn sobre Stochastic Approximation en el Simons Institute) https://www.youtube.com/watch?v=dhEF5pfYmvc

TENGO HASTA AHÍ...