

APRENDIZAJE REFORZADO

CLASE 4

Julían Martínez

DOS PROBLEMAS:

- EVALUAR UNA POLÍTICA
- MEJORAR UNA POLÍTICA

DETALLE COMPUTACIONAL

$$\mu_k = \frac{\sum_{j=1}^k x_j}{k} = \mu_{k-1} \frac{(k-1)}{k} + \frac{1}{k} x_k$$

Media Incremental

$$= \left[\mu_{k-1} + \frac{1}{k} (x_k - \mu_{k-1}) \right]$$

A simple bandit algorithm

Initialize, for $a = 1$ to k :

$$Q(a) \leftarrow 0$$

$$N(a) \leftarrow 0$$

Loop forever:

$$A \leftarrow \begin{cases} \operatorname{argmax}_a Q(a) & \text{with probability } 1 - \varepsilon \\ \text{a random action} & \text{with probability } \varepsilon \end{cases} \quad (\text{breaking ties randomly})$$

$$R \leftarrow \text{bandit}(A)$$

$$N(A) \leftarrow N(A) + 1$$

$$Q(A) \leftarrow Q(A) + \frac{1}{N(A)} [R - Q(A)]$$

Evaluación de una política

$$v_{\pi}^{k+1}(s) = R(s) + \gamma \sum_{s'} v_{\pi}^k(s') p_{ss'}^{\pi}$$

Ecuaciones de Bellman

$$v_*(s) = \max_a \left[\sum_{s'} \left\{ r(s, a, s') + \gamma v_*(s') \right\} p_{s, s'}^a \right]$$

Mejora de una política

$$\pi_{k+1}(s) := \arg \max_a q_{\pi_k}(s, a)$$

Evaluación y mejora

$$\pi_k \longrightarrow v_{\pi_k} \longrightarrow q_{\pi_k} \longrightarrow \pi_{k+1}$$

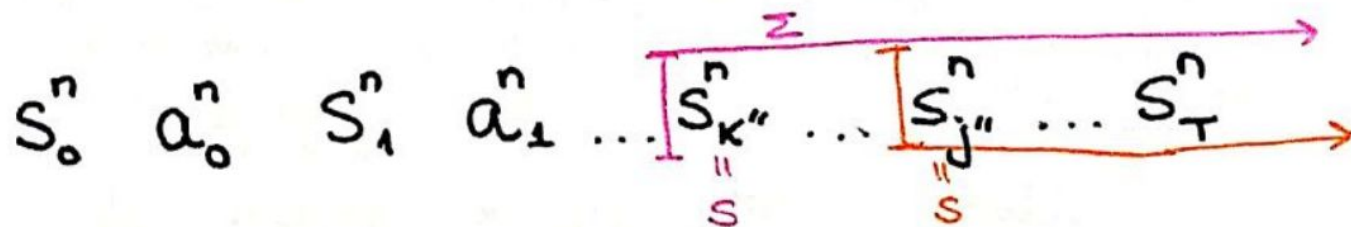
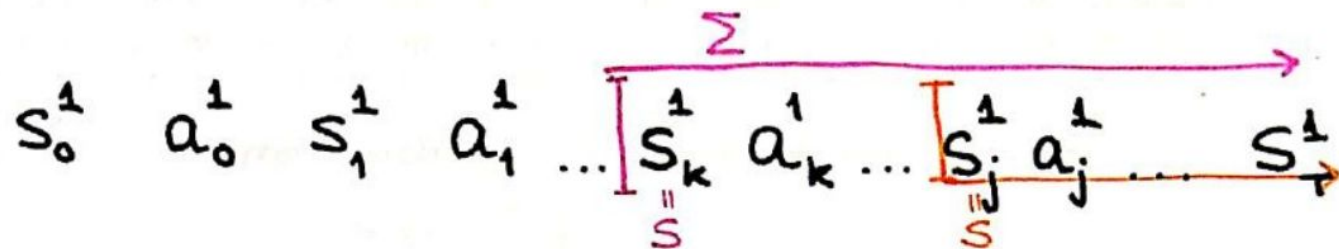
PREDICCIÓN LIBRE DE MODELO

Episode: $s_0 \ a_0 \ s_1 \ a_1 \ s_2 \ a_2 \ \dots \ s_j \ a_j \ \dots \ s_T$

¿Podemos hacer evaluación, aprendiendo tan sólo de la experiencia?

MONTE CARLO

$$v_{\pi}(s) = E_{\pi}[G_t | S_t = s]$$



LEY DE LOS GRANDES NÚMEROS (APROXIMANDO AL ESPERANZA)

$$\frac{\sum_{i=1}^n X_i}{n} \xrightarrow{n \rightarrow \infty} \mathbb{E}[X] \quad X_i \text{ iid}$$

FIRST/EVERY VISIT MC

$$G_{\tau_s^l}^l = \sum_{t=\tau_s^l+1}^T R_t^l \cdot \gamma^{t-\tau_s^l-1}$$

τ_s^l = tiempo de la PRIMERA visita a s

$$v_{\pi}(s) \approx \sum_{l=1}^n G_{\tau_s^l}^l / N(s)$$

$N(s)$ = # episodios donde s
fue visitado

SUTTON

First-visit MC prediction, for estimating $V \approx v_\pi$

Input: a policy π to be evaluated

Initialize:

$V(s) \in \mathbb{R}$, arbitrarily, for all $s \in \mathcal{S}$

$Returns(s) \leftarrow$ an empty list, for all $s \in \mathcal{S}$

Loop forever (for each episode):

Generate an episode following π : $S_0, A_0, R_1, S_1, A_1, R_2, \dots, S_{T-1}, A_{T-1}, R_T$

$G \leftarrow 0$

Loop for each step of episode, $t = T-1, T-2, \dots, 0$:

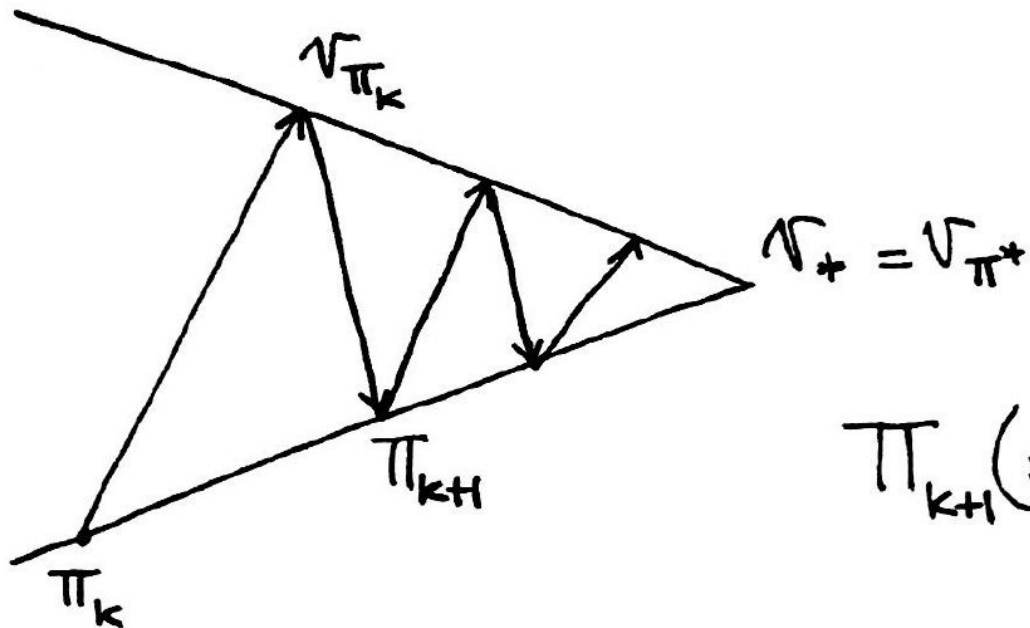
$G \leftarrow \gamma G + R_{t+1}$

Unless S_t appears in S_0, S_1, \dots, S_{t-1} :

Append G to $Returns(S_t)$

$V(S_t) \leftarrow \text{average}(Returns(S_t))$

EVALUACIÓN Y MEJORA



$$\pi_{k+1}(s) := \arg \max_a q_{\pi_k}(s, a)$$

$$\pi_k \longrightarrow \tilde{V}_{\pi_k} \longrightarrow q_{\pi_k} \longrightarrow \pi_{k+1}$$



EXPLORATION



preciso reconstruir

$$q_{\pi} \Rightarrow \{P_{s,s'}^a\}_{s,s'} \quad \forall a$$

\Rightarrow preciso visitar

los $|A| \times |\mathcal{S}|$ estados!

- *Exploration starts*: Cada corrida comenzarla recorriendo todos los posibles pares (s,a)

Monte Carlo ES - ONLINE

Monte Carlo ES (Exploring Starts), for estimating $\pi \approx \pi_*$

Initialize:

$\pi(s) \in \mathcal{A}(s)$ (arbitrarily), for all $s \in \mathcal{S}$

$Q(s, a) \in \mathbb{R}$ (arbitrarily), for all $s \in \mathcal{S}, a \in \mathcal{A}(s)$

$Returns(s, a) \leftarrow$ empty list, for all $s \in \mathcal{S}, a \in \mathcal{A}(s)$

Loop forever (for each episode):

Choose $S_0 \in \mathcal{S}, A_0 \in \mathcal{A}(S_0)$ randomly such that all pairs have probability > 0

Generate an episode from S_0, A_0 , following π : $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T$

$G \leftarrow 0$

Loop for each step of episode, $t = T-1, T-2, \dots, 0$:

$G \leftarrow \gamma G + R_{t+1}$

Unless the pair S_t, A_t appears in $S_0, A_0, S_1, A_1, \dots, S_{t-1}, A_{t-1}$:

Append G to $Returns(S_t, A_t)$

$Q(S_t, A_t) \leftarrow \text{average}(Returns(S_t, A_t))$

$\pi(S_t) \leftarrow \operatorname{argmax}_a Q(S_t, a)$

OTRA MANERA, ON POLICY

- On policy MC **control**: En cada corrida, calcular la política mejorada, y con su versión ϵ -greedy samplear la siguiente corrida

BLACKJACK - DAVID SILVER

- States (200 of them):
 - Current sum (12-21)
 - Dealer's showing card (ace-10)
 - Do I have a "useable" ace? (yes-no)
- Action **stick**: Stop receiving cards (and terminate)
- Action **twist**: Take another card (no replacement)
- Reward for **stick**:
 - +1 if sum of cards $>$ sum of dealer cards
 - 0 if sum of cards = sum of dealer cards
 - -1 if sum of cards $<$ sum of dealer cards
- Reward for **twist**:
 - -1 if sum of cards $>$ 21 (and terminate)
 - 0 otherwise
- Transitions: automatically **twist** if sum of cards $<$ 12



COMENTARIOS SOBRE MC

- La aproximación para un estado no depende de la aproximación para los otros estados. Por lo tanto puedo aproximar la función de valor para sólo algunos estados.
- En algunos casos, inclusive conociendo la distribución del ambiente, puede ser inclusive más conveniente que usar ecuaciones de Bellman.

OFF POLICY PREDICTION VIA IMPORTANCE SAMPLING

- Busco una política **óptima** (*evaluation*), pero para eso preciso **aproximar bien** la dinámica del ambiente (*exploration*).

Tengo dos políticas: la que quiero obtener, π (*target*) y la que voy a usar para explorar b (*behavior*).

Por ahora, sólo quiero evaluar la política π
¿Cómo estimar π usando muestras generadas con b ?

Haz lo que yo digo, pero no lo que yo hago.

IMPORTANCE SAMPLING

$$\mathbb{E}_p[f(x)] = \sum_x f(x) p(x)$$

puede que $p(x) \gg 0$ PERO $|f(x)| \simeq 0$
 \Rightarrow MC no funciona bien!

- p puede ser difícil de samplear o tener varianza grande.

APROXIMAR LA ESPERANZA CON RESPECTO A P MEDIANTE
MUESTRAS OBTENIDAS CON Q

$$E_p[f(X)] = \sum_x f(x) \cdot q(x) \cdot \frac{p(x)}{q(x)}$$

$$= E_q \left[f(X) \cdot \frac{p(X)}{q(X)} \right]$$

$$\approx \frac{\sum_{i=1}^n f(X_i) \frac{p(X_i)}{q(X_i)}}{n}$$

$$; X_i \sim q$$

$$w_s = \frac{p(X_s)}{q(X_s)}$$

IMPORTANCE
WEIGHTS

$$V(\bar{X}) = \frac{1}{n} \cdot V(X)$$

Possible criteno p2r determinur q :

$$\min_q V_q(f(X)w(X))$$

$$q^*(x) = \frac{|f(x)| p(x)}{\sum_x |f(x)| p(x)}$$

$f(w) = \int_{-\infty}^{\infty} f(x) e^{-w \cdot x} dx$
 $\frac{\partial f(w)}{\partial w} = \int_{-\infty}^{\infty} f(x) (-x) e^{-w \cdot x} dx = -\int_{-\infty}^{\infty} x f(x) e^{-w \cdot x} dx$
 $H = -\sum p(x) \log p(x)$
 $\frac{\partial}{\partial S} \left(\frac{1}{2} G^2 S^2 + r S \frac{\partial V}{\partial S} + \frac{\partial V}{\partial k} - r \cdot V \right) = 0$
 $(Q, q, m) = \sum_{i=1}^n \left[\frac{D_i}{m q_i} S_i + c_i D_i + \frac{q_i H_i}{2} \left(m_i \left(1 - \frac{D_i}{p_i} \right) - 1 \right) \right]$
 $\left[\frac{\frac{\partial \Delta p(s, \phi)}{\partial \phi}}{\frac{\partial \Delta M(s, \phi)}{\partial \phi}} \right] = \begin{bmatrix} \gamma & -\omega \\ -\beta & 0 \end{bmatrix} \begin{bmatrix} \Delta p(s, \phi) \\ \Delta M(s, \phi) \end{bmatrix}$
 $\int_0^{\pi} (\log \sin x)^2 dx = -\frac{1}{2} \int_0^{\pi} (\log \cos x)^2 dx = -\frac{\pi}{2} \left\{ \frac{\pi}{12} + (\log 2)^2 \right\}$

VOLVIENDO AL PROBLEMA ORIGINAL

$$\begin{aligned} & P_{\pi}(A_t = a, S_{t+1} = s_{t+1}, A_{t+1} = a_{t+1}, \dots, S_T = s_T | S_t = s) \\ &= P_{\pi}(A_t = a, S_{t+1:T} = s_{t+1:T}, A_{t+1:T} = a_{t+1:T} | S_t = s) \\ &= \pi(a|s) P_{s, s_{t+1}}^a \pi(a_{t+1}|s_{t+1}) P_{s_{t+1}, s_{t+2}}^{a_{t+1}} \dots P_{s_{T-1}, s_T}^{a_{T-1}} \\ \Rightarrow w_{t:T-1} &= \pi(a|s) \cdot \prod_{k=t+1}^{T-1} \frac{\pi(a_k | s_k)}{b(a_k | s_k)} \end{aligned}$$

VOLVIENDO AL PROBLEMA ORIGINAL

$$\begin{aligned} E_{\pi}[G_t | S_t = s] &= E_{\pi}[\psi(A_t, S_{t+1}, A_{t+1}, \dots, S_T) | S_t = s] \\ &= E_b[\psi(A_t, S_{t+1:T}, A_{t+1:T-1}) \cdot w(A_{t:T-1})] \end{aligned}$$

- Tiene mayor varianza que hacer MC directamente y converge más lento.
- Me permite usar experiencia generada con otra política.

HACIA TEMPORAL DIFFERENCE (EVALUATION)

$$v_{\pi}(s) = E[R_{t+1} + \gamma E[G_{t+2} | S_{t+1}] | S_t = s]$$

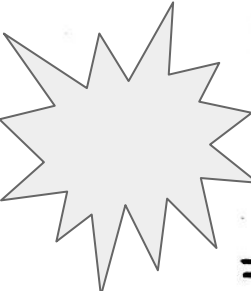
$$= E[R_{t+1} + \gamma v_{\pi}(S_{t+1}) | S_t = s]$$

$$= v_{\pi}(s) + E[\{R_{t+1} + \gamma v_{\pi}(S_{t+1})\} - v_{\pi}(s) | S_t = s]$$

LA MÁQUINA DE HACER CHORIZOS

$$v_{\pi}(s) = E[R_{t+1} + \gamma E[G_{t+2} | S_{t+1}] | S_t = s]$$

$$= E[R_{t+1} + \gamma v_{\pi}(S_{t+1}) | S_t = s]$$


$$= v_{\pi}(s) + E[\{R_{t+1} + \gamma v_{\pi}(S_{t+1})\} - v_{\pi}(s) | S_t = s]$$

MONTE-CARLO



TEMPORAL DIFFERENCE

$$v_{\pi}(s) + E \left[\left\{ R_{t+1} + \gamma v_{\pi}(S_{t+1}) \right\} - v_{\pi}(s) \middle| S_t = s \right]$$

MONTE-CARLO

$$v^{k+1}(s) = v^k(s) + \frac{(G^{k+1} - v^k(s))}{N^k(s) + 1}$$

TEMPORAL DIFFERENCE

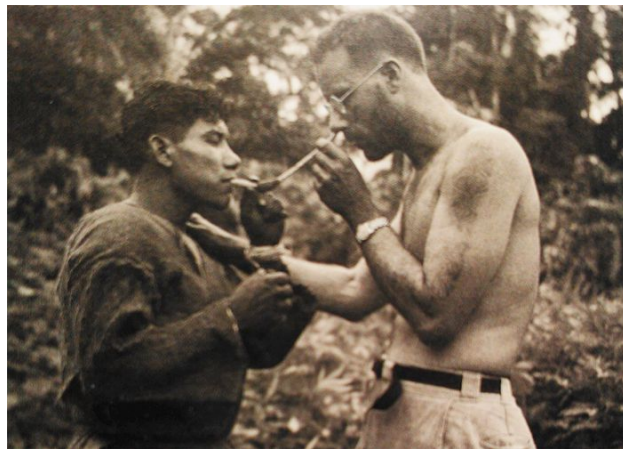
$$v^{t+1}(S_t) = v^t(S_t) + \alpha \left[\overbrace{R_{t+1} + \gamma v^t(S_{t+1})}^{\text{Estimated Return}} - v^t(S_t) \right]$$

TD Error

MC VS TD (VARIANZA VS SESGO)

- Monte Carlo usa todo el episodio para tener una nueva aproximación.
- Temporal difference en cada paso, genera una nueva aproximación

EL LADO OSCURO...



If both TD and Monte Carlo methods converge asymptotically to the correct predictions, then a natural next question is “Which gets there first?” In other words, which method learns **faster**? Which makes the more efficient use of limited data? At the current time this is an open question in the sense that no one has been able to prove mathematically that one method converges faster than the other. In fact, it is not even clear what is the most appropriate formal way to phrase this question! In practice, however, TD methods

Learning to predict by the methods of **temporal differences**, Sutton, 1988:

<https://link.springer.com/content/pdf/10.1007/BF00115009.pdf>

