

APRENDIZAJE REFORZADO

CLASE 2

Julían Martínez

DADA UNA POLÍTICA, ¿CÓMO CALCULAR LA FUNCIÓN DE VALOR?

$$v(s) = E[G_t | S_t = s] = E[R_{t+1} + \gamma G_{t+1} | S_t = s]$$

$$= R(s) + \gamma E[G_{t+1} | S_t = s]$$

$$= R(s) + \gamma \sum_{s'} E[G_{t+1} | S_{t+1} = s'] P(S_{t+1} = s' | S_t = s)$$

$$= R(s) + \gamma \sum_{s'} v(s') P_{ss'}^\pi$$

ECUACIÓN DE
BELLMAN

SOBRE LA ECUACIÓN DE BELLMAN

$\mathcal{L} = \phi \frac{\partial}{\partial t}$
 $f(w) = \int_{-\infty}^{\infty} f(x) e^{-2\pi i x w} dx \quad \frac{d}{dt}$
 $\nabla \cdot E = 0 \quad \nabla \times E = -\frac{1}{c} \frac{\partial H}{\partial t}$
 $\nabla \cdot H = 0 \quad \nabla \times H = \frac{1}{c} \frac{\partial E}{\partial t}$
 $-i\hbar \frac{\partial}{\partial t} \Psi = H \Psi$
 $\rho \left(\frac{\partial v}{\partial t} + v \cdot \nabla v \right) = -\nabla p + \nabla \cdot T + f$
 $H = -\sum p(x) \log p(x)$
 $\frac{1}{2} G^2 S^2 \frac{\partial^2 V}{\partial S^2} + r S \frac{\partial V}{\partial S} + \frac{\partial V}{\partial t} - r \cdot V = 0$
 $(Q, q_i, m_i) = \sum_{i=1}^n \left[\frac{D_i}{m_i q_i} S_i + c_i v D_i + \frac{q_i H_i v}{2} \left(m_i \left(1 - \frac{D_i}{P_i} \right) - 1 + 2 \right) \right]$
 $\left[\frac{d \Delta p(s, \phi)}{d \phi} \right] = \begin{bmatrix} \gamma & -\mathcal{L} \\ -\beta & 0 \end{bmatrix} \begin{bmatrix} \Delta p(s, \phi) \\ \Delta M(s, \phi) \end{bmatrix}$
 $\int_0^{\frac{\pi}{2}} (\log \sin x)^2 dx = \int_0^{\frac{\pi}{2}} (\log \cos x)^2 dx = \frac{\pi}{2} \left\{ \frac{\pi^2}{12} + (\log 2)^2 \right\}$

- Sistema de ecuaciones lineales
- Mezcla de cadenas de markov

ALGORITMO ITERATIVO (EVALUATION) (JACOBI, GAUSS-SEIDEL)

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

O puedo ir reemplazando a medida que calculo

$$v^k(s) \leftrightarrow v^{k+1}(s)$$

¿POR QUÉ CONVERGE?

Bellman
Expectation
Backup Operator

$$T^{\pi}(v) := R^{\pi} + \gamma P^{\pi} v$$

$$\begin{aligned} \|T^{\pi}(u) - T^{\pi}(v)\|_{\infty} &= \|\gamma P^{\pi}(u - v)\|_{\infty} \\ &\leq \gamma \|u - v\|_{\infty} \end{aligned}$$

Contraction Mapping Theorem

POLÍTICA ÓPTIMA


$$V_*(s) = \max_{\pi} V_{\pi}(s) = V_{\pi^*}(s)$$

FUNCIÓN DE ESTADO-ACCIÓN

$$q_{\pi}(s, a) := E_{\pi}[G_t | S_t = s, A_t = a]$$

¿CÓMO MEJORAR UNA POLÍTICA?

$$\begin{aligned} V_{\pi}(s) &= E[G_t | S_t = s] \\ &= \sum E[G_t | S_t = s, A_t = a] \times P(A_t = a | S_t = s) \end{aligned}$$

$$V_{\pi}(s) = \sum_a q_{\pi}(s, a) \pi(a|s)$$


ALGUNAS OBSERVACIONES

$$V_{\pi}(s) = \sum_a q_{\pi}(s, a) \pi(a|s)$$

$$\pi^+(s) := \arg \max_a q_{\pi}(s, a)$$

$$V_{\pi}(s) \leq q_{\pi}(s, \pi^+(s))$$

EFFECTIVAMENTE, MEJORA LA FUNCIÓN DE VALOR

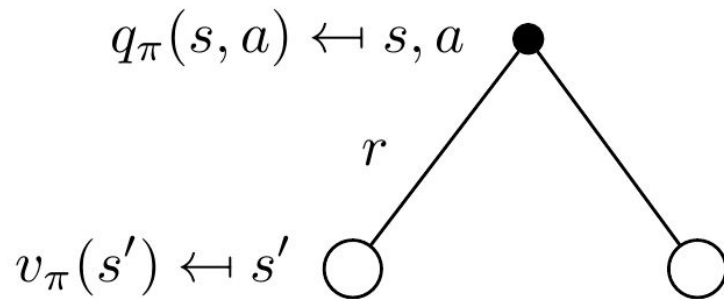
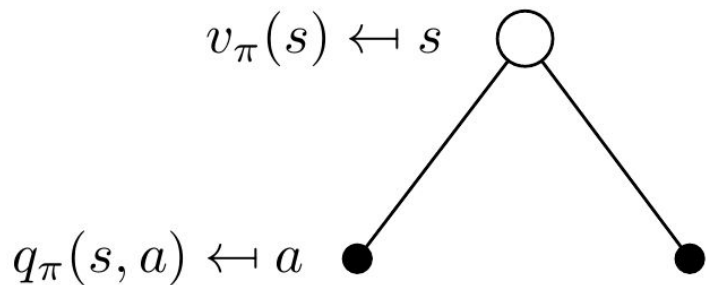
$$\begin{aligned} v_{\pi}(s) &\leq q_{\pi}(s, \pi'(s)) = E[R_{t+1} + \gamma v_{\pi}(S_{t+1}) | S_t = s, A_t = \pi'(s)] \\ &\stackrel{\textcircled{=}}{=} E_{\pi'}[R_{t+1} + \gamma v_{\pi}(S_{t+1}) | S_t = s] \quad \begin{array}{l} \text{Como } \pi' \text{ determinística} \\ \Rightarrow A_{t+1} = \pi'(S_{t+1}) \end{array} \\ &\leq E_{\pi'}[R_{t+1} + \gamma q_{\pi}(S_{t+1}, \pi'(S_{t+1})) | S_t = s] \end{aligned}$$

Condiciono en A_t

$$(V-EA) \quad v_{\pi}(s) = \sum_a q_{\pi}(s, a) \pi(a|s)$$

Condiciono en $S_{\{t+1\}}$

$$(EA-V) \quad q_{\pi}(s, a) = \sum_{s'} \left\{ r(s, a, s') + \gamma v_{\pi}(s') \right\} \cdot p_{s, s'}^a$$

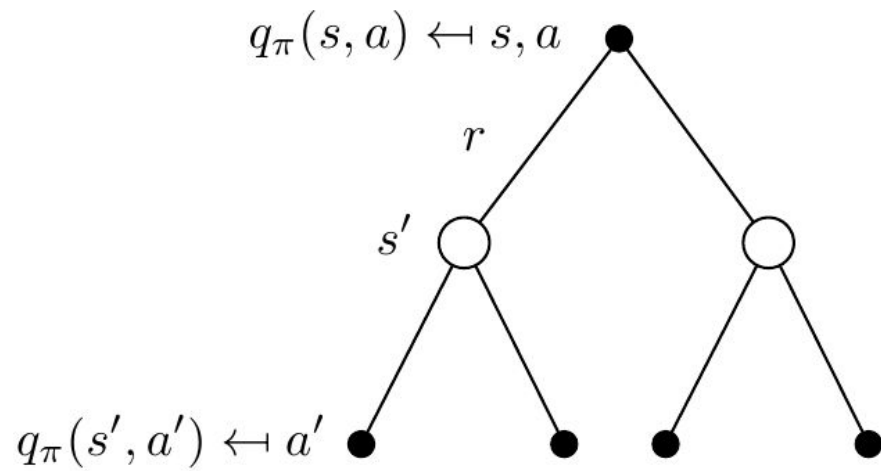
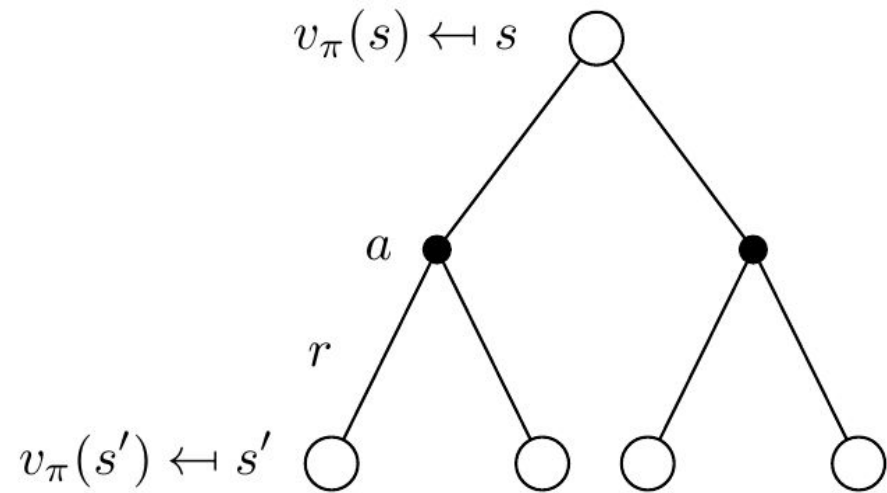


EXISTENCIA DE POLÍTICA ÓPTIMA

Theorem

For any Markov Decision Process

- *There exists an optimal policy π_* that is better than or equal to all other policies, $\pi_* \geq \pi, \forall \pi$*
- *All optimal policies achieve the optimal value function, $v_{\pi_*}(s) = v_*(s)$*
- *All optimal policies achieve the optimal action-value function, $q_{\pi_*}(s, a) = q_*(s, a)$*



ECUACIONES DE OPTIMALIDAD DE BELLMAN

$$v_*(s) = \max_{a \in A(s)} q_{\pi_*}(s, a)$$

$$v_{\pi}(s) = \sum_a q_{\pi}(s, a) \pi(a|s)$$



$$q_{\pi}(s, a) = \sum_{s'} \left\{ r(s, a, s') + \gamma v_{\pi}(s') \right\} \cdot p_{s, s'}^a$$

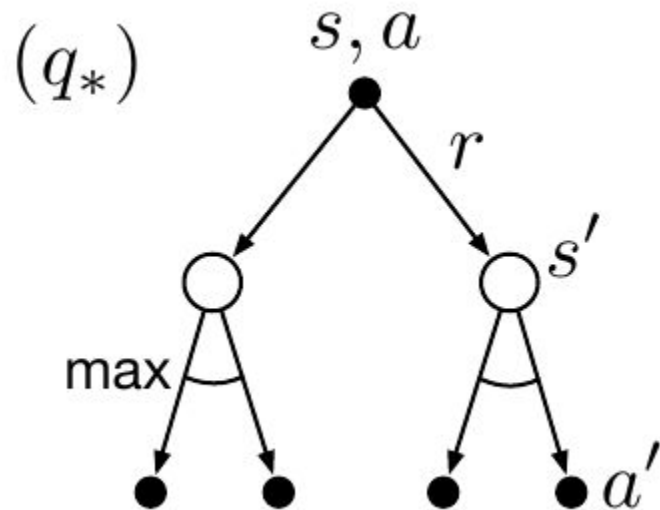
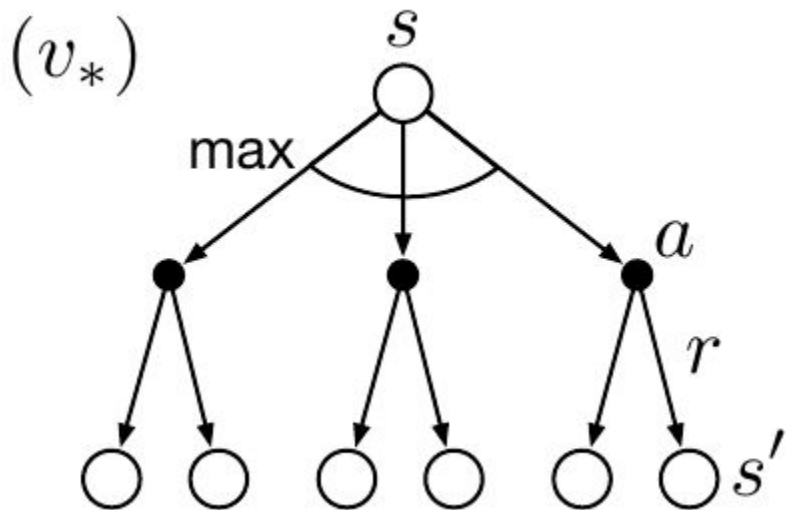
$$v_*(s) = \max_a \left[\sum_{s'} \left\{ r(s, a, s') + \gamma v_*(s') \right\} p_{s, s'}^a \right]$$

- Algoritmos iterativos

- Ecuaciones NO lineales!

$$q_*(s, a) = \sum_{s'} \left\{ r(s, a, s') + \gamma \max_{a'} q_*(s', a') \right\} p_{s, s'}^a$$

BACKUP DIAGRAMS



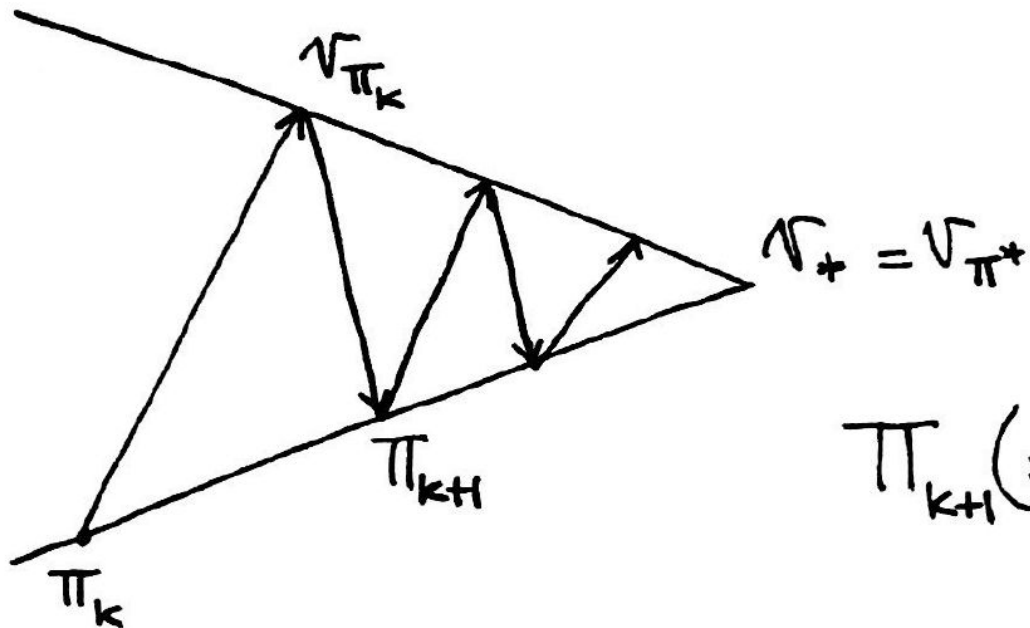
OPTIMALIDAD DE BELLMAN

T_{ij} = Costo de viajar de i a j

O_{ij} = Costo del viaje ÓPTIMO de i a j

$$O_{ij} = \min_k [T_{ik} + O_{kj}]$$

EVALUACIÓN Y MEJORA

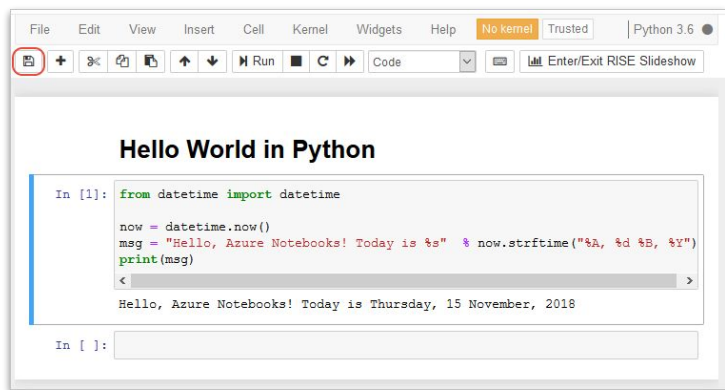


$$\pi_{k+1}(s) := \arg \max_a q_{\pi_k}(s, a)$$

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



INTRODUCCIÓN A OPENAI GYM



- Introducción General
- Para que “jueguen”: Mountain
- Ejemplo del Robot (Batería)

INTRODUCCIÓN A OPENAI GYM

- ¿Cómo instalarlo? ubuntu 18.04, Python, jupyter, open ai gym
 - Linux:
 - `sudo apt-get update`
 - `sudo apt-get install python3 python3-pip ipython3 python3-fontconfig`
 - `sudo apt-get install libglu1-mesa-dev freeglut3-dev mesa-common-dev python-opengl`
 - `pip3 install numpy pandas matplotlib jupyter gym`
 - Windows:
 - Virtual Box:
 - instalar ubuntu 16.04 y usar el instructivo de la parte de Linux
 - WSL:(inspirado en <https://github.com/openai/gym/issues/11#issuecomment-242950165>)
 - instalar Windows Subsystem for Linux (WSL): <https://docs.microsoft.com/en-us/windows/wsl/install-win10>
 - instalar ubuntu 16.04 LTS para WSL yendo a Microsoft Store (barra de búsqueda de Windows) y buscando Ubuntu 16.04
 - correr una consola WSL (buscar Ubuntu en la barra de búsqueda de Windows)
 - realizar los mismos pasos que en el instructivo de linux en esa consola
 - instalar vcXsrv/xming;
 - correr vcXsrv (elegir one large window); tipear en la consola de comandos de WSL: `export DISPLAY=:0`
 - Correr jupyter: `jupyter notebook --no-browser`
 - Google colab: ir a google colab, <https://colab.research.google.com/notebooks/welcome.ipynb#recent=true>, elegir la solapa Github y buscar en <https://github.com/javkrei/aprendizaje-reforzado-austral>