

$$\begin{aligned}
 v_*(s) &= \max_a E[R_{t+1} + \gamma v_*(S_{t+1}) | S_t = s, A_t = a] \\
 &= \max_a \sum_{s', r} [r + \gamma v_*(s')] p(s', r | s, a)
 \end{aligned}$$

(Handwritten note: v_ and q_* are related)*

$$\begin{aligned}
 q_*(s, a) &= E[R_{t+1} + \gamma \max_{a'} q_*(S_{t+1}, a') | S_t = s, A_t = a] \\
 &= \sum_{s', r} p(s', r | s, a) [r + \gamma \max_{a'} q_*(s', a')]
 \end{aligned}$$

~~Problema~~ ¿Cómo obtener q_* , v_* , π_* ?

POLICY EVALUATION (PREDICTION) v_π ?

Menú 1: Vimos ecuaciones lineales (Δ si $|S| \uparrow$)

Método iterativo: $v_{k+1}(s) := E_\pi [R_{t+1} + \gamma v_k(S_{t+1}) | S_t = s]$

Esto quizás se podría poner en slides recordando la ecuación lineal!

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

Obs: $v_k \equiv v_\pi$ es el punto fijo! ($\gamma < 1$ o desde todos los estados, bajo π , se llega a un estado terminal)

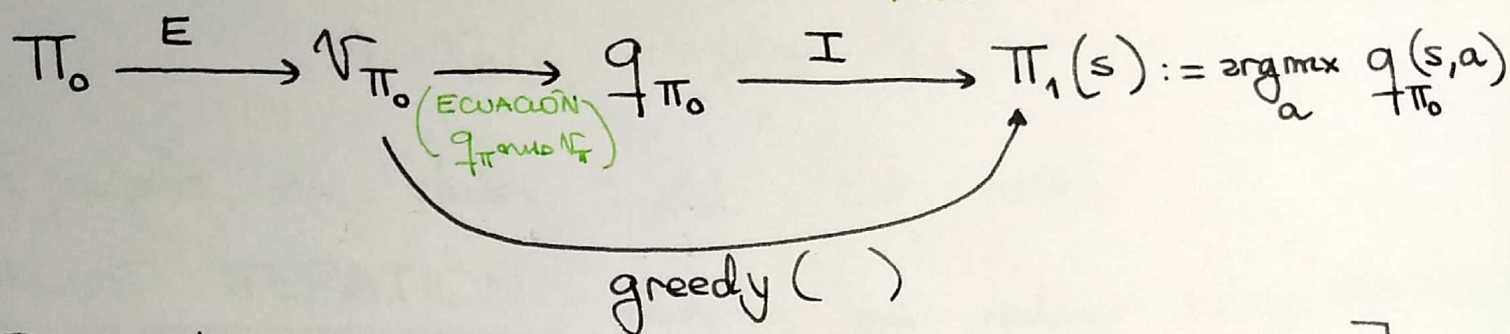
Obs 2: Podría ir actualizando $v_k(s')$ por $v_{k+1}(s')$ a medida que voy calculando!

HACER SLIDE con pseudocódigo
y ¿Ejemplo 4.1 Gridworld?

Mejorando una política π (o hasta encerr y pulir :)

C3
H2

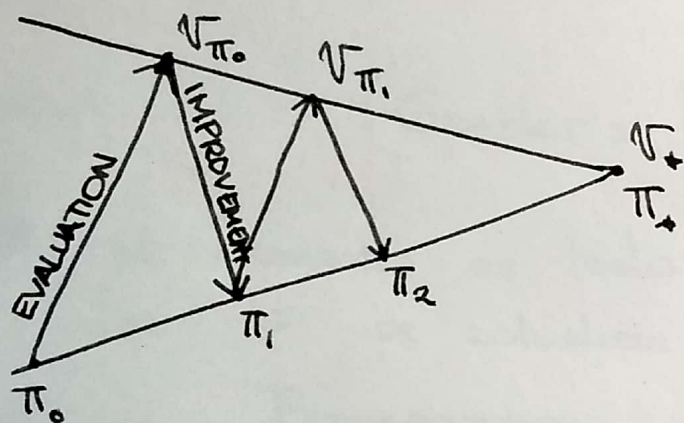
HACER UN SLIDE !!!



Recordemos

$$v_{\pi}(s) = \sum_a \left[R_s^a + \sum_{s'} v_{\pi}(s') p_{ss'}^a \right] \pi(a|s)$$

DOS "INGREDIENTES": v_{π} y π . Busco $v^* \equiv \max_{\pi} v_{\pi} \equiv v_{\pi^*}$



TEOREMA: π', π deterministic / $q_{\pi}(s, \pi'(s)) \geq v_{\pi}(s)$
 $\Rightarrow v_{\pi'}(s) \geq v_{\pi}(s)$

Obs: $\pi'(s) := \operatorname{argmax}_a q_{\pi}(s, a)$ ($q_{\pi}(s, a) = R_s^a + \gamma \sum_{s'} v_{\pi}(s') p_{ss'}^a$)
 Como π determinista $v_{\pi}(s) = q_{\pi}(s, \pi(s))$

$$\begin{aligned} v_{\pi}(s) &\leq q_{\pi}(s, \pi'(s)) = E[R_{t+1} + \gamma v_{\pi}(S_{t+1}) | S_t = s, A_t = \pi'(s)] \\ &\stackrel{\text{Como } \pi' \text{ determinista}}{\Rightarrow A_t = \pi'(S_t)} \\ &\equiv E_{\pi'}[R_{t+1} + \gamma v_{\pi}(S_{t+1}) | S_t = s] \\ &\leq E_{\pi'}[R_{t+1} + \gamma q_{\pi}(S_{t+1}, \pi'(S_{t+1})) | S_t = s] \\ &\leq E_{\pi'}[R_{t+1} + \gamma R_{t+2} + \gamma^2 v_{\pi}(S_{t+2}) | S_t = s] \leq \dots \leq v_{\pi'}(s) \end{aligned}$$

Obs: Si $q_{\pi}(s, \pi'(s)) = v_{\pi}(s)$

C3
H3

$$\Rightarrow v_{\pi}(s) = \underset{a}{\operatorname{argmax}} q_{\pi}(s, a)$$

Ejemplo en slides del cr rental.

VALUE ITERATION: Variante para reducir policy EVALUATION

Se utiliza sólo 1 paso de update cuando hacemos la evaluación de la policy. Combinar una de ~~evaluación~~ y una de improvement

Recordar: $v_{*}^{k+1}(s) = \max_a \sum_{s', r} [r + \gamma v_{*}^{k+1}(s')] p(s', r | s, a)$ (4.10)

Ejemplo 4.3 (Gambler's Problem)
 En slides!
 c/iteración puede NO corresponder a un VALUE FUNCTION!

Hasta el momento es todo SINCRÓNICO: Todos los s se actualizan en SIMULTÁNEO.

x ejemplo: Blackjackman $\approx 10^{20} = |\mathcal{S}|$

ASINCRÓNICO (capítulo 8)

- Usar 4.10 actualizando 1 sólo s en c/iteración.
- Actualizar algunos ~~estados~~ ^{estados} s en evaluación y otros en improvement.
- No actualizar estados poco probables!

(Meter el slide donde hace el summary de los algoritmos)

¿POR QUÉ CONVERGEN ESTAS ITERACIONES?

C3
H4

BELLMAN EXPECTATION
BACKUP OPERATOR

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

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

COMO $\gamma < 1$ es un contracción! \Rightarrow CONTRACTION MAPPING THM.

(T^{π} tiene ! punto fijo)

Idem con $T^*(v) := \max_a R^a + \gamma P^a v$