

# Valor Optimal

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

Ec. bellman

$$V^*(s) = \underbrace{R(s)}_{\text{recomp. inmediatas}} + \max_a \underbrace{\sum_{s'} P_{sa}(s') V^*(s')}_{\text{recomp. esperadas futuras}}$$

Encontrar  $\pi^*$  (Política Optimal)  
Sabiendo  $V^*$  (fn. Valor optimal)

$$\pi^* = \underset{a}{\operatorname{argmax}} \left[ \sum_{s'} P_{sa}(s') V^*(s') \right]$$

♦ Nos da una acción óptima para cada estado.

notations

$$V^*(s) = V(s) \geq V^\pi(s)$$

↓  
fn. de Valor  
óptimo

↓  
fn. Valor para  
la pol. óptima

Resumen 1 → Encontrar  $V^*$

2 → Usar agente para  
encontrar  $\pi^*$

$V^* \rightarrow$  Algo:

Inicial  $V(s) = 0$

y para cada estado, se  
actualiza.

$$V^* \left\{ V(s) = R(s) + \max_a \sum_{s'} P_{sa}(s') V(s') \right\}$$

$$\begin{bmatrix} V(1,1) \\ \vdots \\ V(1,2) \end{bmatrix}$$



$$[V(4,3)]$$

$$R''$$

Actualizando  $\begin{cases} \text{Sincronizado} \\ \text{No-sincronizado} \end{cases}$

Estados observados (+1, -1)

$$\underbrace{V(s)}_{\text{new values}} = R(s) + \max_a \left( \underbrace{\sum_{s' \in \mathcal{S}_a} P(s') V(s')}_{\text{Old Values}} \right)$$

(Sincronizado)

Python  $\rightarrow$  Operador de Bellman backup

$$V := B(V)$$

$$V(s) \rightarrow V^*(s)$$

Converge

$K=9 \rightarrow \text{convergence}$

Exemplo = Temos

$$V^* =$$

0.86	0.9	0.95	+1
0.81	<del>0.8</del>	0.69	-1
0.78	0.75	0.71	0.49

$$H = \frac{y}{3}$$

	$\rightarrow$	$\rightarrow$	$\rightarrow$	+1
2	$\uparrow$	$\rightarrow$	$\uparrow$	-1
1	$\uparrow$	$\leftarrow$	$\leftarrow$	$\leftarrow$
	1	2	3	4

$x$

$(x, y)$   
 $(3, 1)$   
 $\leftarrow$

$$\begin{aligned}
 \sum_{s'} P(s') V^*(s') &= 0.8 \cdot 0.75 \\
 &+ 0.1 \cdot 0.69 \\
 &+ 0.1 \cdot 0.49 \\
 &= 0.71
 \end{aligned}$$

$V^*$

Espacio de Acciones

$\{h, s, e, o\}$



Estado 3,1 {n, e, o}

Condiciones

$$V^* \rightarrow \pi^*$$

$$\text{Des } V^\pi \Rightarrow V^*$$

$$\Rightarrow \pi^*$$

Sist. lcs.  
lineal  
11 ins  
11 eds.

Iteración de la  
Política

(Policy Iteration)

→ Iniciar con  $\pi$  aleatoria

Repeat

$$\rightarrow \text{Set } V = V^\pi \text{ (Eds. Solenar)}$$

$$\rightarrow \text{Set } \pi(s) = \underset{a}{\operatorname{argmax}} \sum_{s'} P_a(s') V(s')$$

if  $\pi$  is not stable then

It. II

→ + lento, + claro  
→ espacio de acciones  
→  $\pi^*$

It. V

→ Convergencia a  $V^* \rightarrow \pi^*$   
→ + usado  $V$ , + rápida

Qué pasa si no sé  $P_{sa}(s')$ ?

$P_{sa}(s') = ?$

$P_{sa}(s') \Rightarrow 0.1 \leftarrow 0.8 \rightarrow 0.1$

$P_{sa}(s')$

= # veces q' el robot tome  
la acción  $a$  en estado  $s$   
y llegue a  $s'$

# veces toma acción  
 $a$  en estado  $s$

Prob. de transición



Algo  $\Rightarrow$  Repeat {

$\epsilon$ -greedy  
 $\rightarrow 90\% \pi$   
 $10\% \text{ Random}$

pas (1) Toma acción a con respecto  $\pi$

- Actualize estimacion de  $P_{sa}(s')$

(det. estimacion de  $R$ )

- Resuelve Ecs de Bellman por la iteracion

- Actualizo

$$\pi(s) = \underset{a}{\operatorname{argmax}} \sum_{s'} P_{sa}(s') V(s')$$

Cuando no se fu. de Recomp

$\rightarrow$  Algo RL Acciones de bolsa

$\rightarrow$  Algo RL AA

# Dilema de la Explotación vs Explotación

