

Inteligencia Artificial

Procesos de Decisión de Markov

Edgar Andrade, Ph.D.

Matemáticas Aplicadas y Ciencias de la computación

Última revisión: Noviembre de 2023



MACC
Matemáticas Aplicadas y
Ciencias de la Computación

Contenido

Evaluación de políticas

Políticas óptimas

Mejoramiento

Policy iteration

Value iteration



Contenido

Evaluación de políticas

Políticas óptimas

Mejoramiento

Policy iteration

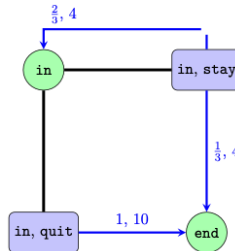
Value iteration



Procesos de decisión markovianos

Para cada ronda $k = 1, 2, \dots$

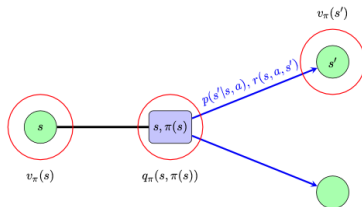
- ▶ Estados: in o end
- ▶ Acciones: stay o quit
- ▶ Si acción == quit:
recompensa = \$10
utilidad += recompensa y
end
- ▶ Si acción == stay:
recompensa = \$4
utilidad += recompensa y
lanza el dado:
 - ▶ Si dado en $\{1, 2\}$: end
 - ▶ Si no: siguiente ronda



Ecuación de Bellman

$$v_{\pi}(s) = \begin{cases} 0, & \text{si } s \text{ es terminal} \\ q_{\pi}(s, \pi(s)), & \text{en otro caso} \end{cases}$$

$$q_{\pi}(s, \pi(s)) = \sum_{s'} (p(s'|s, \pi(s)) [r + \gamma v_{\pi}(s')])$$



$$v_{\pi}(s) = \begin{cases} 0, & \text{si } s \text{ es terminal} \\ \sum_{s'} (p(s'|s, \pi(s)) [r + \gamma v_{\pi}(s')]) , & \text{en otro caso} \end{cases}$$



Policy evaluation

👉 Usar la ecuación de Bellman como una regla iterativa:

$$v_{k+1}(s) = \sum_{s'} \left(p(s'|s, \pi(s)) \left[r + \gamma v_k(s') \right] \right)$$



Policy evaluation

Algorithm 1: Iterative policy evaluation

Data: una política π

Result: valor $V_\pi(s)$ para cada s

$V(s) \leftarrow 0$ para cada estado s ;

repeat

$\Delta \leftarrow 0$;

for *cada estado s* **do**

$v \leftarrow V(s)$;

$V(s) \leftarrow \sum_{s'} (p(s'|s, \pi(s)) [r + \gamma V(s')])$;

$\Delta \leftarrow \max(\Delta, |v - V(s)|)$;

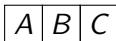
end

until $\Delta < \theta$;



Ejercicio

Suponga el siguiente entorno con $\gamma = 0,8$:



- ▶ El agente parte de la casilla A y puede moverse a izquierda o a derecha.
- ▶ Si el agente llega a C, el juego termina.
- ▶ Moverse a una casilla distinta a C tiene como recompensa -1. La recompensa de moverse a C es 10.
- ▶ Moverse a la izquierda resulta en la casilla de la izquierda con probabilidad 1. Moverse a la izquierda en la casilla A deja al agente en A.
- ▶ Moverse a la derecha resulta en la casilla de la derecha con probabilidad 0.9 y en la misma casilla con probabilidad 0.1.



Ejercicio

Sea $\pi(s)$ la política de moverse a la derecha para todo s .

Realice dos iteraciones del algoritmo iterativo de evaluación de política para la política π .



Solución

$$V(A) = 0$$

$$V(B) = 0$$

$$V(C) = 0$$



Solución

$$V(A) = 0$$

$$V(B) = 0$$

$$V(C) = 0$$

$$V(A) = -1$$

$$V(B) = 8,9$$

$$V(C) = 0$$



Solución

$$V(A) = 0$$

$$V(B) = 0$$

$$V(C) = 0$$

$$V(A) = -1$$

$$V(B) = 8,9$$

$$V(C) = 0$$

$$V(A) = 5,328$$

$$V(B) = 9,612$$

$$V(C) = 0$$



Solución

$$V(A) = 0$$

$$V(B) = 0$$

$$V(C) = 0$$

$$V(A) = -1$$

$$V(B) = 8,9$$

$$V(C) = 0$$

$$V(A) = 5,328$$

$$V(B) = 9,612$$

$$V(C) = 0$$

$$V(A) = 6,347$$

$$V(B) = 9,669$$

$$V(C) = 0$$



Solución

$$V(A) = 0$$

$$V(B) = 0$$

$$V(C) = 0$$

$$V(A) = -1$$

$$V(B) = 8,9$$

$$V(C) = 0$$

$$V(A) = 5,328$$

$$V(B) = 9,612$$

$$V(C) = 0$$

$$V(A) = 6,347$$

$$V(B) = 9,669$$

$$V(C) = 0$$

$$V(A) = 6,469$$

$$V(B) = 9,674$$

$$V(C) = 0$$



Contenido

Evaluación de políticas

Políticas óptimas

Mejoramiento

Policy iteration

Value iteration



Política óptima

$$\pi' \geq \pi \quad \text{sii} \quad v_{\pi'}(s) \geq v_{\pi}(s) \text{ para todo } s$$



Política óptima

$$\pi' \geq \pi \quad \text{sii} \quad v_{\pi'}(s) \geq v_{\pi}(s) \text{ para todo } s$$

$$\pi^* \text{ es óptima} \quad \text{sii} \quad \pi^* \geq \pi \text{ para toda } \pi$$



Política óptima

$$\pi' \geq \pi \quad \text{sii} \quad v_{\pi'}(s) \geq v_{\pi}(s) \text{ para todo } s$$

$$\pi^* \text{ es óptima} \quad \text{sii} \quad \pi^* \geq \pi \text{ para toda } \pi$$

$$v_*(s) = \max_{\pi} v_{\pi}(s)$$



Política óptima

$$\pi' \geq \pi \quad \text{sii} \quad v_{\pi'}(s) \geq v_{\pi}(s) \text{ para todo } s$$

$$\pi^* \text{ es óptima} \quad \text{sii} \quad \pi^* \geq \pi \text{ para toda } \pi$$

$$v_*(s) = \max_{\pi} v_{\pi}(s)$$

$$q_*(s, a) = \max_{\pi} q_{\pi}(s, a)$$



Backup diagrams

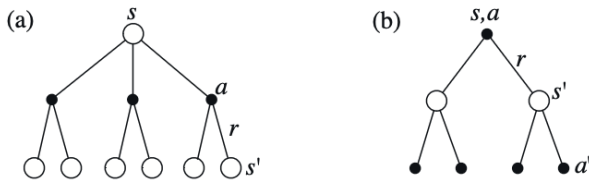


Figure 3.4: Backup diagrams for (a) v_π and (b) q_π .



Backup diagrams

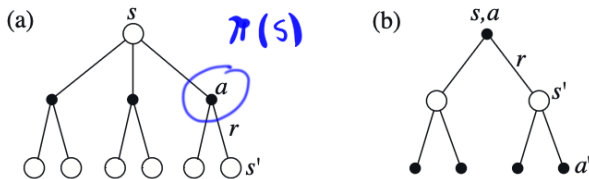


Figure 3.4: Backup diagrams for (a) v_π and (b) q_π .



Backup diagrams

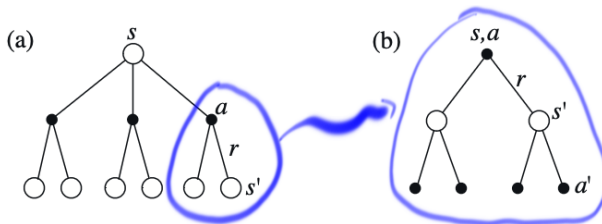
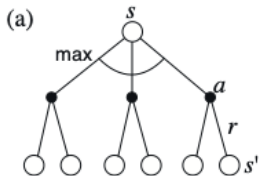


Figure 3.4: Backup diagrams for (a) v_π and (b) q_π .



Ecuación de Bellman (2/2) v_*

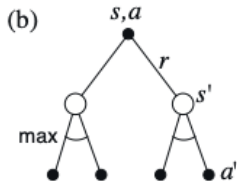


Backup diagram for (a) v_*

$$v_*(s) = \max_a q_*(s, a)$$



Ecuación de Bellman (2/2) q_*



Backup diagram for (b) q_*

$$q_*(s, a) = \sum_{s'} \left(p(s'|s, a) \left[r + \gamma v_*(s') \right] \right)$$



Contenido

Evaluación de políticas

Políticas óptimas

Mejoramiento

Policy iteration

Value iteration



Mejoramiento de una política

- ▶ Mejorar π cambiando la acción para un solo estado s .

$$\pi'(s') = \begin{cases} a, & \text{si } s' = s \\ \pi(s'), & \text{en otro caso} \end{cases}$$



Mejoramiento de una política

- ▶ Mejorar π cambiando la acción para un solo estado s .

$$\pi'(s') = \begin{cases} a, & \text{si } s' = s \\ \pi(s'), & \text{en otro caso} \end{cases}$$

- ▶ Si $q_{\pi'}(s, a) \geq q_{\pi}(s, a)$, entonces $v_{\pi'}(s) \geq v_{\pi}(s)$ para todo s .



Mejoramiento de una política

- ▶ Mejorar π cambiando la acción para un solo estado s .

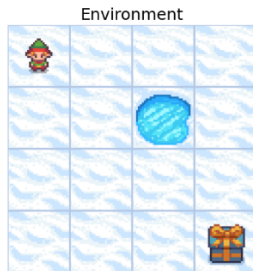
$$\pi'(s') = \begin{cases} a, & \text{si } s' = s \\ \pi(s'), & \text{en otro caso} \end{cases}$$

- ▶ Si $q_{\pi'}(s, a) \geq q_{\pi}(s, a)$, entonces $v_{\pi'}(s) \geq v_{\pi}(s)$ para todo s .
- ▶ Mejoramiento avaro dada π :

$$\begin{aligned} \pi'(s) &= \arg \max_a q_{\pi}(s, a) \\ &= \arg \max_a \sum_{s'} \left(p(s'|s, a) \left[r + \gamma v_{\pi}(s') \right] \right) \end{aligned}$$



Ejemplo ilustrativo (2/2)



Lago congelado

Policy

↓	→	↑	→
↓	→	→	→
↓	↓	←	↓
↓	→	→	→

Values

0.17	0.19	0.21	0.32
0.23	0.24	0.0	0.44
0.31	0.4	0.31	0.71
0.35	0.47	0.71	0.0

Adelante: $\frac{1}{3}$, Izquierda: $\frac{1}{3}$, Derecha: $\frac{1}{3}$,

Recompensa: 1 el regalo.

Descuento: 0.9



El valor aumenta

$$v_{k+1}(s) = q_{\pi_{k+1}}(s, \pi_{k+1}(s))$$

👉 Por definición de $v_{k+1}(s)$



El valor aumenta

$$\begin{aligned}v_{k+1}(s) &= q_{\pi_{k+1}}(s, \pi_{k+1}(s)) \\ &= \max_a q_{\pi_k}(s, a)\end{aligned}$$

☞ Por definición de $\pi_{k+1}(s) = \arg \max_a q_{\pi_k}(s, a)$



El valor aumenta

$$\begin{aligned}v_{k+1}(s) &= q_{\pi_{k+1}}(s, \pi_{k+1}(s)) \\&= \max_a q_{\pi_k}(s, a) \\&\geq q_{\pi_k}(s, \pi_k(s))\end{aligned}$$

👉 Por definición de máx



El valor aumenta

$$\begin{aligned}v_{k+1}(s) &= q_{\pi_{k+1}}(s, \pi_{k+1}(s)) \\&= \max_a q_{\pi_k}(s, a) \\&\geq q_{\pi_k}(s, \pi_k(s)) \\&= v_{\pi_k}(s)\end{aligned}$$

☞ Por definición de $v_k(s)$



Contenido

Evaluación de políticas

Políticas óptimas

Mejoramiento

Policy iteration

Value iteration



Policy iteration

$$\pi_0 \xrightarrow{E} v_{\pi_0} \xrightarrow{I} \pi_1 \xrightarrow{E} v_{\pi_1} \xrightarrow{I} \dots \pi^* \xrightarrow{E} v_* \xrightarrow{I} \pi^*$$

$$\pi_k \xrightarrow{E} v_{\pi_k}$$

Encontrar $v_{\pi_k}(s)$ para todo s



Policy iteration

$$\pi_0 \xrightarrow{E} v_{\pi_0} \xrightarrow{I} \pi_1 \xrightarrow{E} v_{\pi_1} \xrightarrow{I} \dots \pi^* \xrightarrow{E} v_* \xrightarrow{I} \pi^*$$

$$\pi_k \xrightarrow{E} v_{\pi_k}$$

Encontrar $v_{\pi_k}(s)$ para todo s

$$v_{\pi_k} \xrightarrow{I} \pi_k$$

Mejorar π_k con base en $v_{\pi_k}(s)$
para todo s



Policy iteration

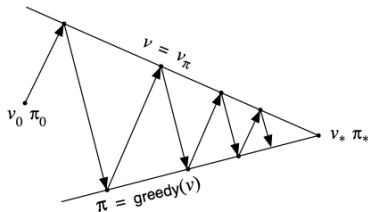
$$\pi_0 \xrightarrow{E} v_{\pi_0} \xrightarrow{I} \pi_1 \xrightarrow{E} v_{\pi_1} \xrightarrow{I} \dots \pi^* \xrightarrow{E} v_* \xrightarrow{I} \pi^*$$

$$\pi_k \xrightarrow{E} v_{\pi_k}$$

Encontrar $v_{\pi_k}(s)$ para todo s

$$v_{\pi_k} \xrightarrow{I} \pi_k$$

Mejorar π_k con base en $v_{\pi_k}(s)$
para todo s



Policy iteration — pseudocódigo —

Algorithm 2: Policy iteration

Data: una política π
Result: una política óptima π^*
// 1. Inicialización:
 $V(s) \leftarrow 0$ para cada estado s ;
 $policy_stable \leftarrow \text{False}$;
repeat
 // 2. Evaluación de política:
 $V \leftarrow policy_evaluation(\pi)$;
 // 3. Mejoramiento de política:
 $policy_stable \leftarrow \text{True}$;
 for cada estado s **do**
 $a \leftarrow \pi(s)$;
 $\pi(s) \leftarrow \arg \max_a \sum_{s'} (p(s'|s, a) [r + \gamma V(s')])$;
 if $a \neq \pi(s)$ **then**
 $policy_stable \leftarrow \text{False}$;
 end
 end
until $policy_stable$;



Contenido

Evaluación de políticas

Políticas óptimas

Mejoramiento

Policy iteration

Value iteration



Value iteration (1/3)

- ▶ Truncar la evaluación de la política después de una iteración para cada estado.



Value iteration (1/3)

- ▶ Truncar la evaluación de la política después de una iteración para cada estado.
- ▶ Se combina el mejoramiento de la política con la evaluación truncada de la política:

$$v_{k+1}(s) = \max_a \sum_{s'} \left(p(s'|s, a) \left[r + \gamma v_k(s') \right] \right)$$



Value iteration (1/3)

- ▶ Truncar la evaluación de la política después de una iteración para cada estado.
- ▶ Se combina el mejoramiento de la política con la evaluación truncada de la política:

$$v_{k+1}(s) = \max_a \sum_{s'} \left(p(s'|s, a) \left[r + \gamma v_k(s') \right] \right)$$

- ▶ Se puede demostrar que la sucesión $\{v_k\}$ converge a v_* .



Value iteration (1/3)

- ▶ Truncar la evaluación de la política después de una iteración para cada estado.
- ▶ Se combina el mejoramiento de la política con la evaluación truncada de la política:

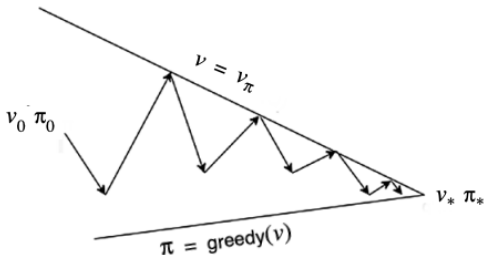
$$v_{k+1}(s) = \max_a \sum_{s'} \left(p(s'|s, a) \left[r + \gamma v_k(s') \right] \right)$$

- ▶ Se puede demostrar que la sucesión $\{v_k\}$ converge a v_* .
- ▶ Finalmente, se obtiene π_* :

$$\pi_*(s) = \arg \max_a \sum_{s'} \left(p(s'|s, a) \left[r + \gamma v_*(s') \right] \right)$$



Value iteration (2/3)



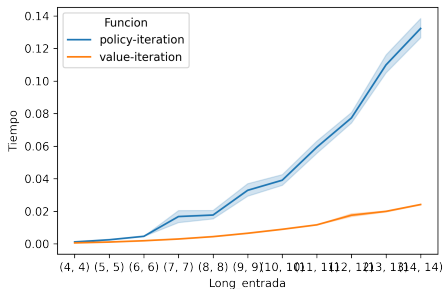
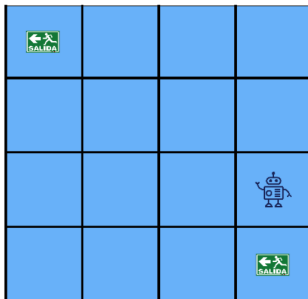
Value iteration (3/3)

Algorithm 3: Value iteration

Result: una política óptima π^* $V(s) \leftarrow 0$ para cada estado s ;**repeat** $\Delta \leftarrow 0$; **for** cada estado s **do** $v \leftarrow V(s)$; $V(s) \leftarrow \max_a \sum_{s'} (p(s'|s, a) [r + \gamma V(s')])$; $\Delta \leftarrow \max(\Delta, |v - V(s)|)$; **end****until** $\Delta < \theta$;**for** cada estado s **do** $\pi^*(s) \leftarrow \arg \max_a \sum_{s'} (p(s'|s, a) [r + \gamma V(s')])$;**end**



Comparación de tiempos



Take away

En esta sesión usted aprendió:

- ▶ Usar programación dinámica para evaluar una política.
- ▶ Definir el concepto de política óptima.
- ▶ Definir un proceso de mejora de una política.
- ▶ Usar los algoritmos de Policy iteration y Value iteration para buscar una política óptima.

