

$$MDP = \langle S, A, a, T, r, S_F, \gamma \rangle$$

S := estados

A := Acciones

a := Acciones legales

T := Función de transición $S \times A \times S \rightarrow \mathbb{R}$ $T(s, a, s') = P_i[S_{t+1} = s' | S_t = s, A_t = a]$

S := $S \rightarrow \mathbb{R}$ (s, a, s') Recompensa de ir a s' desde s haciendo a

$0 \leq \gamma \leq 1$ Factor de descuento

Retorno

$$R_t = r_t + \gamma R_{t+1} + \gamma^2 R_{t+2} \dots$$

$$R_t = r_t + \gamma R_{t+1}$$

Encontrar una política

$\pi: S \rightarrow A$ política determinista

$$\pi(s) = a$$

Necesito comparar políticas

$$\pi_1 < \pi_2 \Leftrightarrow V^{\pi_1}(s) \leq V^{\pi_2}(s) \quad \forall s \in S$$

$$V^{\pi}(s) = E[R_t | S_t = s]$$

$$V^{\pi}(s) = \sum_{s' \in S} T(s, \pi(s), s') [r(s, \pi(s), s') + \gamma V^{\pi}(s')]$$

(cálculo $V^{\pi}(s)$)

def valor_politica(π, MDP, ϵ):

$V(s) :=$ random if s not in S_T else 0 for s in S

while true:

$\Delta = 0$
while true:

for s in S/S_T

$$V := V^{\pi}(s)$$

$$V^{\pi}(s) = \sum T(s, \pi(s), s') [r(s, \pi(s), s') + \gamma V^{\pi}(s')]$$

$$\Delta := \max(\Delta, |V - V^{\pi}(s)|)$$

if $\Delta < \epsilon$ break;

return V^{π}

$$V^{\pi}(s) \leq \underline{V^{\pi^*}(s)} \quad \forall s \in S \quad \forall \pi \in \Pi$$

$V^*(s) \leftarrow$ Valor óptimo

$$V^*(s) = \max_{a \in A(s)} \sum_{s' \in S} T(s, a, s') [r(s, a, s') + \gamma V^*(s')]$$

Algoritmos de programación dinámica

$$\pi_0 \rightarrow V^{\pi_0} \rightarrow \pi_1 \rightarrow V^{\pi_1} \rightarrow \pi_2 \rightarrow V^{\pi_2} \rightarrow \dots \pi^*$$

$$\pi^*(s) = \arg \max_{a \in A(s)} \sum_{s' \in S} T(s, a, s') [r(s, a, s') + \gamma V^*(s')]$$

def iteracion_politicas (MDP, ϵ):

$\pi(s) = \text{random}(A(s)) \quad \forall s \in S/S_T$

Optima = False

While True:

$V^{\pi} = \text{Valor_politica}(\pi, \text{MDP}, \epsilon)$

Optima = True

Para cada s en S/S_T :

anterior = $\pi(s)$

$$\pi(s) = \arg \max_{a \in A} \sum_{s' \in S} T(s, a, s') + [r(s, a, s') + \gamma V^{\pi}(s')]$$

Si $\pi(s) \neq \text{anterior}$

Optima = False

Si Optima break

return π

Iteración

def Iteracion_valor (MDP, ϵ)

$V^*(s) = \text{random if } s \notin S_T \text{ else } 0$

While True:

$\Delta = 0$

for S in S/S_T :

$V = V^*(s)$

$V^*(s) = \max_{a \in A(s)} \sum_{s' \in S} (T(s, a, s') + [r(s, a, s') + \gamma V^*(s)])$

$\Delta = \max(\Delta, |V^*(s) - V|)$

Si $\Delta < \epsilon$ break

for S in S/S_T :

$\pi^*(s) = \arg \max_{a \in A(s)} \sum_{s' \in S} (T(s, a, s') + [r(s, a, s') + \gamma V^*(s)])$

return π

$O(\text{max_epoch} * |S|^2 * |A|) \leftarrow \text{polynomial}$