

CLASE 2: MARKOV DECISION PROCESS

C2
H1

MARKOV REWARD PROCESS

$$(\underbrace{S, p}_{\text{markov}}, \underbrace{R, \gamma}_{\text{discount factor}})$$

$$R_s = E[R_{t+1} | S_t = s]$$

REWARD
FUNCTION

discount factor
 $\gamma \in [0, 1]$

En general ; $R_{t+1} = \psi(S_{t+1})$

Recordar: $E[g(X)] = \sum_x g(x) p(x)$

PONER EJEMPLO EN SLIDES.

RETURN $G_t = R_{t+1} + \gamma R_{t+2} + \dots = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}$

Obs: $\gamma > 0$ evita ciclos en el proceso de Markov.

$\gamma \approx 0$ corto plazo
 $\gamma \approx 1$ largo plazo

VALUE
FUNCTION

$$v(s) := E[G_t | S_t = s]$$

(no depende de t
pues el markov es
homogéneo)

$$= E[R_{t+1} + \gamma(R_{t+2} + \gamma R_{t+3} + \dots) | S_t = s] \text{ RATA}$$

$$= E[R_{t+1} + \gamma G_{t+1} | S_t = s] = E[E[R_{t+1} + \gamma G_{t+1} | S_{t+1}] | S_t = s]$$

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

Agregar antes $E[\psi(X) \cdot Y | X] = \psi(X) E[Y | X]$

BELLMAN
EQUATION

$$v(s) = R_s + \gamma \sum_{s'} v(s') \underbrace{P(S_{t+1} = s' | S_t = s)}_{P_{s,s'}}$$

Obs: Dados $R, p \Rightarrow$ me da ecuaciones por v

$$\begin{bmatrix} v(s_1) \\ \vdots \\ v(s) \end{bmatrix} = \boxed{v = R + \gamma p v}$$

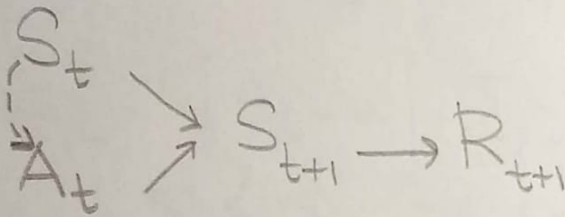
$\cdot O(|S|^3)$ complejidad
 \cdot Programación Dinámica, MC...

MARKOV DECISION PROCESS: $(\mathcal{S}, \mathcal{P}, \mathcal{R}, \gamma, \mathcal{A})$ C2
H2
 \mathcal{A} = conjunto de acciones

Obs: LA ACCIÓN DEL AGENTE MODIFICA EL AMBIENTE

$$P_{ss'}^a := P(S_{t+1}=s' | S_t=s, A_t=a)$$

$$R_s^a = E[R_{t+1} | S_t=s, A_t=a] \quad R_{t+1} = \gamma(S_{t+1})$$



POLÍTICA

$$\pi(a|s) := P(A_t=a | S_t=s)$$

NO DEPENDE DE t !!!

Decide la acción de menor reward en función del ambiente (s)

$$P(S_{t+1}=s' | S_t=s) = \sum_a P_{s_t=s}^{s_{t+1}=s'}(A_t=a)$$

PROB. TOTAL PARA LA CONDICIONAL

Menor reward condicional de la condicional:

$$P_B(A|C) = P(A|B, C)$$

$$= \sum_a \pi(a|s) P_{s,s'}^a =: P_{s,s'}^\pi$$

MEZCLA DE CADENAS DE MARKOV

$$R_s^\pi := E[R_{t+1} | S_t=s] = \sum_a \pi(a|s) R_s^a$$

STATE-VALUE

$$V_\pi(s) = E_\pi[G_t | S_t=s]$$

Explicar diferencia entre ambas funciones!

ACTION-VALUE FUNCTION

$$Q_\pi(s,a) = E_\pi[G_t | S_t=s, A_t=a]$$

DESARROLLAR PARA C/DEF EL EJEMPLO DE SLIDES

Bellman Expectation Equation

C2
H3

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

$$q_{\pi}(s, a) = E_{\pi} [R_{t+1} + \gamma q_{\pi}(S_{t+1}, A_{t+1}) | S_t = s, A_t = a]$$

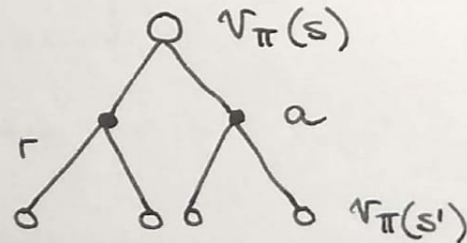
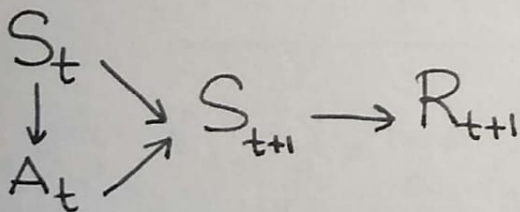
$$\boxed{v_{\pi}(s)} = E_{\pi} [E[G_t | A_t] | S_t = s]$$

$$\begin{aligned} E[G_t | A_t = a] &= \sum_w w P_{S_t=s}(w | A_t = a) = \sum_w w P(w | S_t = s, A_t = a) \\ &= \sum_a \left\{ \sum_w w P(w | S_t = s, A_t = a) \right\} P(A_t = a | S_t = s) \\ &= \sum_a E_{\pi}[G_t | S_t = s, A_t = a] \pi(a | s) = \boxed{\sum_a q_{\pi}(s, a) \pi(a | s)} \end{aligned}$$

$$\boxed{q_{\pi}(s, a)} = R_s^a + \gamma \sum_{s'} v_{\pi}(s') p_{ss'}^a \quad \text{(by conditioning on } S_{t+1})$$

Mezclas

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



$$q_{\pi}(s, a) = R_s^a + \gamma \sum_{s'} p_{ss'}^a \sum_{a'} \pi(a' | s') q_{\pi}(s', a')$$

$$\boxed{v_{\pi} = R^{\pi} + \gamma P^{\pi} v_{\pi}}$$

¿Por qué fórmulas recursivas?

C2
H4

FUNCION DE VALOR OPTIMA:

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

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

DEFINICION: $\pi \geq \pi'$ si $V_{\pi}(s) \geq V_{\pi'}(s) \quad \forall s$

TEOREMA: Para cualquier MDP: (Con Dynamic Programming diriz que se entiende porque ~~funcion~~)

- $\exists \pi_*$ optimo / $\pi_* \geq \pi \quad \forall \pi$
- $V_*(s) = V_{\pi_*}(s) ; \quad Q_*(s,a) = Q_{\pi_*}(s,a) \quad \forall s,a$

Obs: Si tenemos $Q_* \Rightarrow \pi_*(a|s) := \begin{cases} 1 & \text{si } a = \arg\max_a Q_*(s,a) \\ 0 & \text{cc} \end{cases}$

π_* es optimo y deterministico.

ECUACIONES DE OPTIMALIDAD DE BELLMAN

1) $V_*(s) = \max_{a \in A(s)} Q_{\pi_*}(s,a)$ Recordar: $V_{\pi_*}(s) = \sum_a Q_{\pi_*}(s,a) \pi_*(a|s)$

2) $V_*(s) = \max_a \sum_{s',r} p(s',r|s,a) [r + \gamma V_*(s')]$

3) $Q_*(s,a) = \sum_{s',r} p(s',r|s,a) [r + \gamma \max_{a'} Q_*(s',a')]$

NO LINEALES °° no hay fórmula cerrada → Métodos iterativos!!!