

10/3/25

Continuación de procesos de Markov

Distribución de probabilidad para todos los valores
 X_t v. a. en el tiempo t

$$\text{Valores}(X_t) = \{s_1, \dots, s_n\}$$

$$P_t(X_t) = \begin{bmatrix} P_t(X_t = s_1) \\ \vdots \\ P_t(X_t = s_n) \end{bmatrix} \rightarrow \sum_{i=1}^n P_t(X_t = s_i) = 1$$

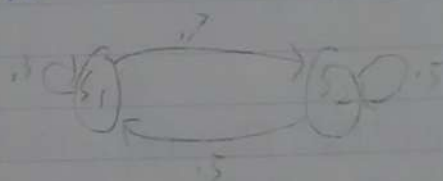
 X_t es un proceso de Markov de primer orden

$$P_t(X_{t+n} | X_t, X_{t-1}, X_{t-2}, \dots, X_0) = P_t(X_{t+n} | X_t)$$

$$P_t(X_{t+1} = s_i | X_t = s_j) = P_{ij} \rightarrow \sum_{i=1}^n P_{ij} = 1$$

Ejemplo:

$$P = \begin{bmatrix} 0.3 & 0.5 \\ 0.7 & 0.5 \end{bmatrix}$$



$$P(X_t) = \begin{bmatrix} .4 \\ .6 \end{bmatrix}$$

$$P(X_{t+1} = s_i) = \sum_{j=1}^n P_{ij} * P(X_t = s_j)$$

$$P(X_{t+1} = s_1) = .3 * .4 + .5 * .6$$

$$P(X_{t+1} = s_2) = .7 * .4 + .5 * .6$$

$$P(X_1) = P P(X_0)$$

$$P(X_2) = P P(X_1) = P P P(X_0)$$

$$\therefore P(X_t) = P^t P(X_0)$$

10/3/25

Si el proceso es estacionario, entonces:
 $P(X_{t+1}) = P(X_t)$ para una t grande
 $P(X_{t+1}) = P(X_t) \rightarrow$

Si es estacionaria entonces hay una t tal que:

$$P(X_t = s_1) = .3 P(X_t = s_1) + .5 P(X_t = s_2)$$

$$P(X_t = s_2) = .7 P(X_t = s_1) + .5 P(X_t = s_2)$$

$$-.7 P(X_t = s_1) + .5 P(X_t = s_2) = 0$$

$$-.5 P(X_t = s_1) - .5 P(X_t = s_2) = -.5$$

$$-1.2 P(X_t = s_1) = -.5 \rightarrow P(X_t = s_1) = \frac{5}{12}$$

$$P(X_t = s_2) = 7/12$$

Concepto Ejemplo:

S conjunto estados

A conjunto acciones

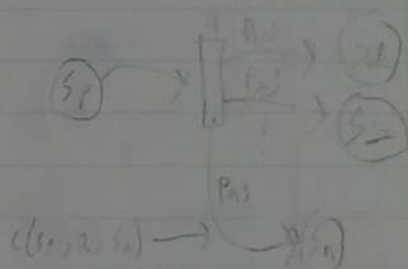
$s_0 \in S$ estado inicial

$S_f \subseteq S$ conjunto de estados finales

Acciones legales $S \rightarrow P(A)$

transición: $S \times A \rightarrow S$

costo-recompensa: $S \times A \rightarrow \mathbb{R} \cup \{-\infty, \infty\}$



Modelo de Decisión
de Markov
(MDP)

11/3/28

Математика де Марков:

MDP: Modelo matemático para toma de decisiones bajo incertidumbre. Fue introducido en 1950's-1960's. El término "Markov" se refiere a Andrey Markov, pues MDP ~~son~~ extensiones de las cadenas de Markov.

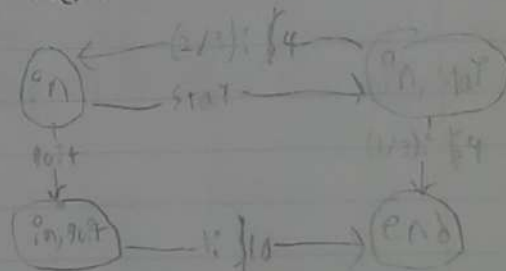
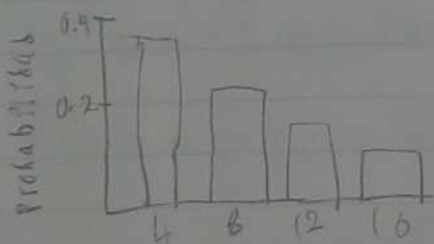
Este se usa en áreas de Robótica, Almacenamiento y Agricultura, como ejemplo.

Ejemplo de Markov: Juego de Dados

En cada ronda:

- Eliges quedarte o salir
- Si te quedas, ganas \$4
- Si te sales, ganas 10

El juego termina si te sales o la cae 1 o 2.



Recompensa (utilidad)

$$E[Y | \pi_{\text{salir}}] = 12$$

1/3/25

Definición de MDP (Proceso de Decisiones de Markov):

class MDP:

S: conjunto de estados

A: conjunto de acciones

def acciones(self, s):

def transición(self, s, a, s'): $\Gamma(s, a, s')$

representa $P(s_{t+1} = s' | s_t = s \wedge A_t = a)$

def ganancia(self, s, a, s'): $R(s, a, s')$

representa un número

def terminal(self, s):

true si s es final, false

$0 \leq \gamma \leq 1$

$$R_t = \gamma^0 r_t + \gamma^1 r_{t+1} + \gamma^2 r_{t+2} + \dots = \sum_{i=0}^{\infty} \gamma^i r_{t+i}$$

12/3/25

Definición alternativa de MDP:
 $MDP = \langle S, A, T, R, \gamma \rangle$

$$S = \{s_1, s_2, \dots, s_n\}$$

$$A = \{a_1, \dots, a_m\}$$

$\pi: S \rightarrow P(A)$ $A(s)$ acciones legales en $s \in S$

$$T: S \times A \times S \rightarrow R \quad T(s, a, s') = P_r(s_{t+1} = s' | s_t = s, a_t = a)$$

$$\sum_{s' \in S} T(s, a, s') = 1$$

$R: S \times A \times S \rightarrow R$ $r(s, a, s')$ es la recompensa de ir de s a s' con la acción a .

$S_T \subseteq S$ conjunto de estados terminales

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

$$R_t = \sum_{i=0}^{\infty} \gamma^i r_{t+i} \quad \text{Retorno en el instante } t$$

Objetivo: Encontrar la mejor política

$\pi: S \rightarrow A$ tal que $\pi(s) = a$

π_1 y π_2 dos políticas diferentes

$$V^\pi(s) = E^\pi[R_t | s_t = s]$$

$$V^\pi(s) = \sum_{s' \in S} T(s, \pi(s), s') [r(s, \pi(s), s') + \gamma E^\pi[R_{t+1} | s_{t+1} = s']]$$

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

12/3/28

Ejemplo: Juego de dados

$$S = \{s_1, s_2\}$$

$$A = \{a, b\}$$

$$T(s_1, a, s_1) = 0$$

$$P(s_1, a, s_1) = 0$$

$$T(s_1, a, s_2) = 1$$

$$P(s_1, a, s_2) = 1/6$$

$$T(s_1, b, s_1) = 2/3$$

$$P(s_1, b, s_1) = 4$$

$$T(s_1, b, s_2) = 1/3$$

$$P(s_1, b, s_2) = 4$$

$$V^{\pi_1}(s_1) = T(s_1, b, s_1) [P(s_1, b, s_1) + \gamma V^{\pi_1}(s_1)] + T(s_1, b, s_2) [P(s_1, b, s_2) + \gamma V^{\pi_1}(s_2)]$$

$$V^{\pi_1}(s_1) = \frac{2}{3} [4 + V^{\pi_1}(s_1)] + \frac{1}{3} [4 + 0]$$

$$\frac{1}{3} V^{\pi_1}(s_1) = \frac{8}{3} + \frac{4}{3} = 12$$

$$\pi_1 \leq \pi_2 \text{ si y solo si}$$

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

y π^* es una política óptima si:

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

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

Ecuación de optimidad de estado de Bellman

14/3/25

Ejemplo amon magro (0000)

1 2 3 4 5 6

$$T(1, 1, 2) = 1$$

$$T(1, 1, x) = 0 \quad \forall x \neq 2$$

$$T(5, x, 5) = 1 \quad \forall x \in A, \quad \forall x \neq 5$$

$$T(5, x, 1) = 0$$

$$T(6, x, 6) = 1 \quad \forall x$$

$$T(6, x, 1) = 0$$

$$T(1, 6, 2) = 1$$

$$T(1, 6, 1) = 1 - 1 = 0$$

$$T(1, 6, x) = 0 \quad \forall x \neq 1, 2$$