

# Introducción al uso de MDP pero para Estados Continuos

Hasta Ahora: MDP  $(S, A, P_{sa}, \gamma, R)$

$$V^\pi(s) = E[R(s_0) + \gamma R(s_1) \dots]$$
$$V^*(s) = \max_a V^\pi(s)$$
$$\pi^*(s)$$

$$\pi^* = \arg \max_a \sum_{s'} P_{sa}(s') V^*(s')$$

$$V(s) = R(s) + \max_a \gamma \sum_{s'} P_{sa}(s') V(s')$$

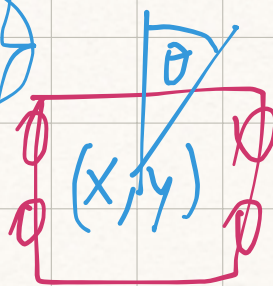
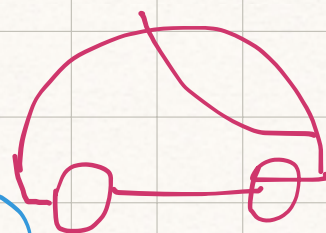
$$V(s) = V(s') \Rightarrow V^*(s)$$

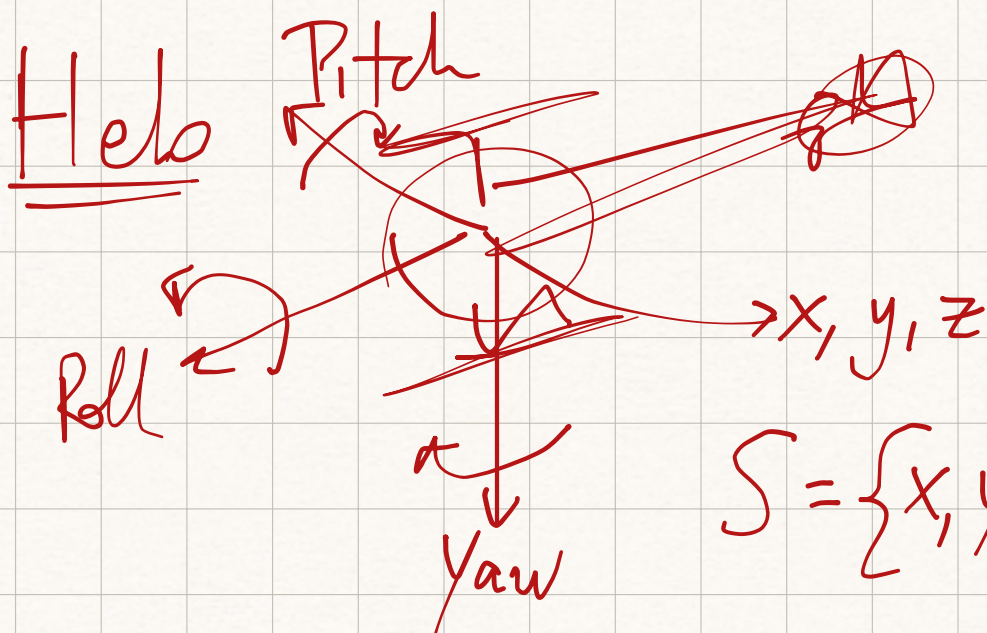
Ejemplos: AA

$$S = \{x, y, \theta, x, y, w\}$$

6-dim

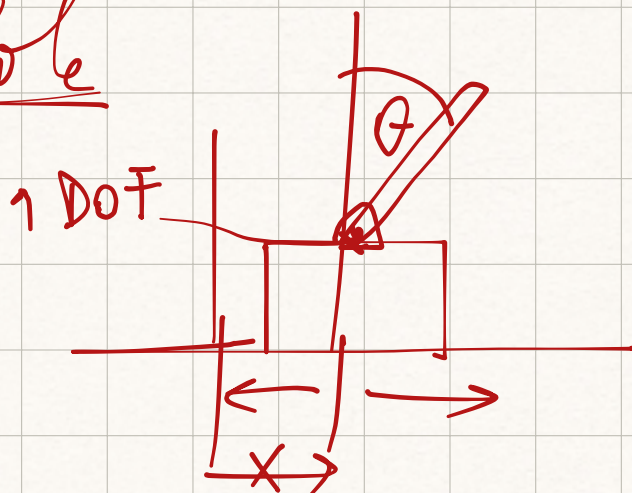
$$S = (3, 1)$$





$$S = \{x, y, z, \phi, \theta, \psi\}$$

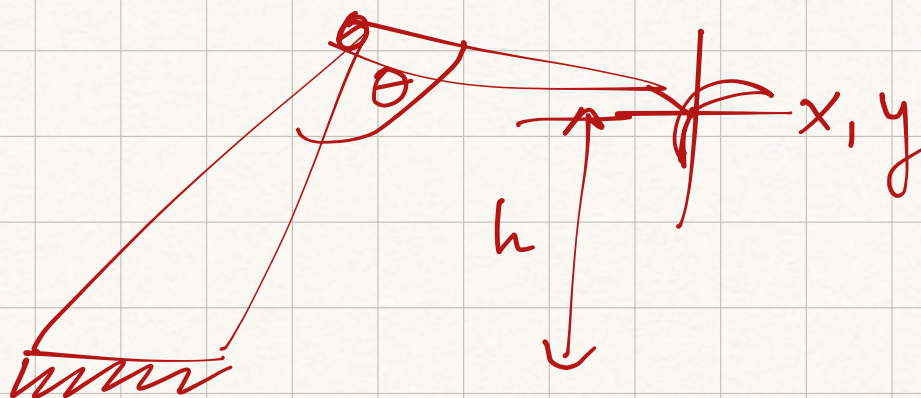
Cart-Pole



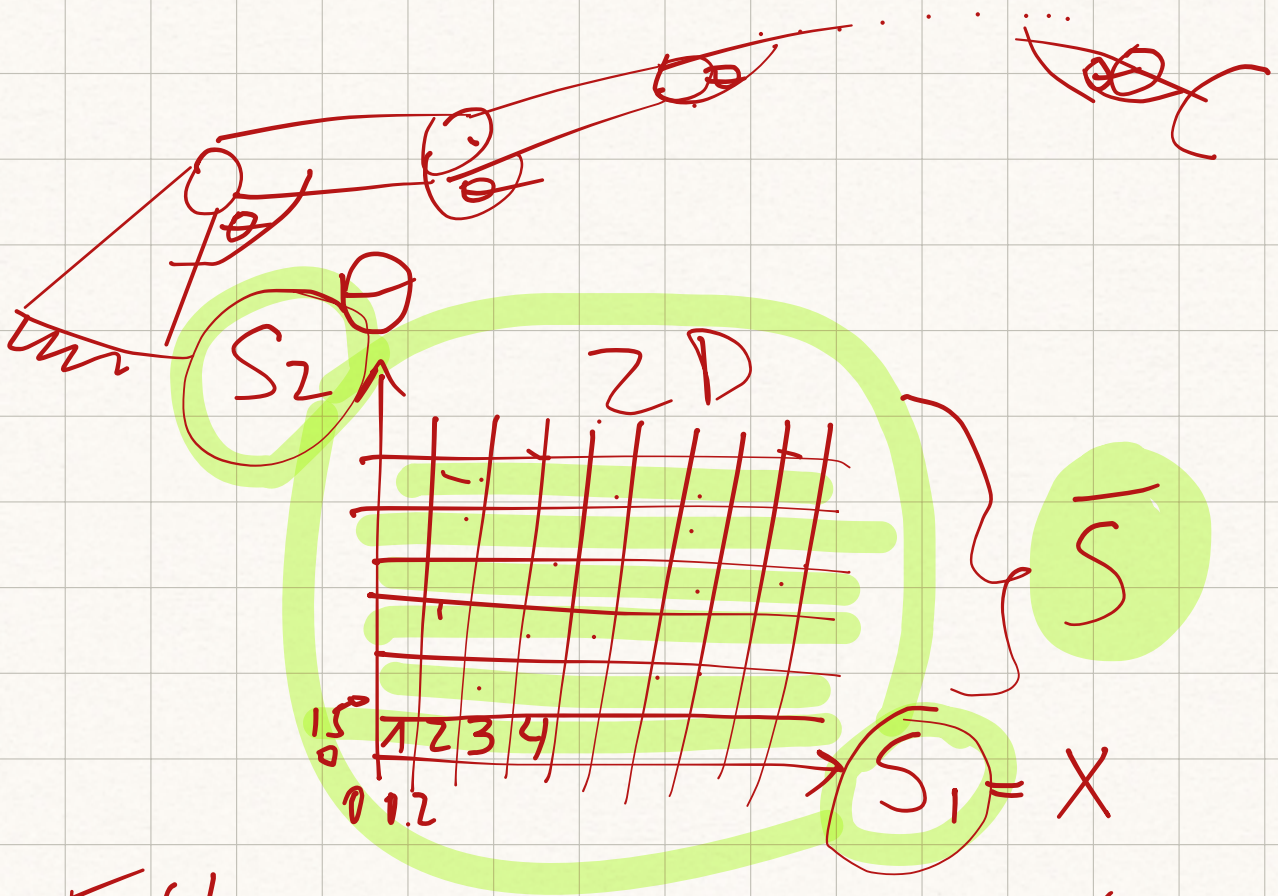
$$S = \{x, \theta, \dot{x}, \dot{\theta}\}$$

$$S = \{-, -, -, -, \dots, n\}$$

$$S = \mathbb{R}^n$$







Fábrica 100 robots  
Cada robot ejecuta 2 acciones  
Conjunto de Estados de la fábrica

Aproximación de la función de Valor

$$V^* \approx \Theta^T \phi(s)$$

Modelos

$S$

Modelo

Regresiones  
lineales

$S' \sim P_{S,a}$

Simulación física

$a$

$$S_{t+1} = AS_t + BA_t$$

Aprender por Data

Pilotos

$$S_0 \xrightarrow{a_0} S_1 \xrightarrow{a_1} S_2 \xrightarrow{a_2} \dots S_T$$

$$S_0^{(1)} \xrightarrow{a_0^{(1)}} S_1^{(1)} \xrightarrow{a_1^{(1)}} S_2^{(1)} \xrightarrow{a_3^{(1)}} \dots S_T^{(1)}$$

$$S_0^{(m)} \xrightarrow{a_0^{(m)}} S_1^{(m)} \xrightarrow{a_1^{(m)}} S_2^{(m)} \xrightarrow{a_3^{(m)}} \dots S_T^{(n)}$$

Carla

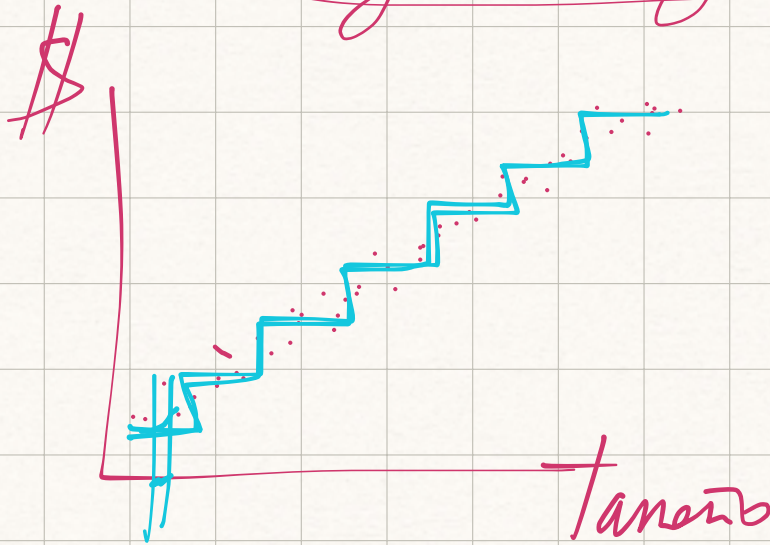
Olivera Modelos

$$(S_{t+1} = AS_t + BA_t)$$



$$(S_{t+1}, A_t, R_t)$$

Algoritmo de RL basado en  
Modelo  $\rightarrow$  ↓ generalización



Model-Free