# Generalizaciones del MDP

1. Recompensas Estado - Acción
2. Horizonte Finito $(T)$    $\dfrac{r(s,a)}{T}$
3. Sistemas Dinámicos Lineales
   - LQR
   - Basados Modelo

## (1) Recompensas Estado - Acción

$LVA$

$$S \longrightarrow R$$
$$S \times A \longrightarrow R$$

$$S_0 \xrightarrow{a_0} S_1 \xrightarrow{a_1} S_2 \ldots \ldots$$

$$R(S_0, a_0) + \gamma R(S_1, a_1) + \gamma^2 R(S_2, a_2)$$

## Ec. Bellman

$$V^*(s) = \max_a \left[ R(s,a) + \gamma \sum_{s'} P_{sa}(s') V^*(s') \right]$$

Las recompensas dependen tbn de la acción

$$\Pi^*(s) = \arg\max_a R(s,a) + \gamma \sum P_{sa}(s') V(s')$$

## 2$^{da}$ Generalización:

## Horizonte Finito

$$MDP \{ S, A, \{P_{sa}\}, \underline{\underline{T}}, R \}$$

$$(\sin \gamma)$$

Ocupando la Gen 1$^{ra}$

$$R(s_0, a_0) + R(s_1, a_1) + \ldots + R(s_T, a_T)$$

ganancia Total

$$E\left[ \phantom{xxxxxxxxxxxxx} \right]$$

↪ maximiza

E vs E
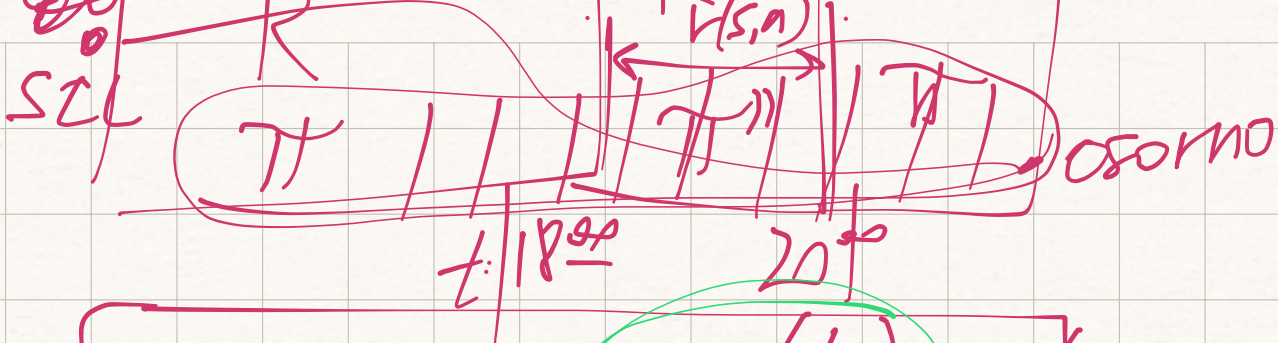


La ácion óptima depende del reloj del robot.

$$\pi_T^*(s) \rightarrow \text{No Estacionaria (cambia con el tiempo)}$$

$$\pi_1^*(s) \rightarrow \text{Estacionaa}$$

$$S_{t+1} \sim P_{S_t a_t}^{(t)}$$

$$R_{(s,a)}^{(t)}$$

## Ejemplo: Avión con combustible

$$V_T^*(s) = E\left[ R(S_t, a_t) + R(S_{t+1}, a_{t+1}) \right.$$
$$\left. + \dots + R(S_T, a_T) \,/\, \pi , S_0 = s \right]$$

$$V_t^*(s) = \max_a R^{(t)}(s,a) + \sum_{s'} P_{sa}^{(t)}(s') V_{t+1}^*(s')$$

$$\pi^*(s) = \arg\max_a \left[ \phantom{xxxx} \right]$$

$$V_t^*(s) = \max_a R(s, a)$$

## Programación Dinámica

### 3) Sist. Dinámicos Lineales

$$\Rightarrow MDP \{S, A, P_{sa}, T, R\}$$

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

$$\Rightarrow A = \mathbb{R}^d$$

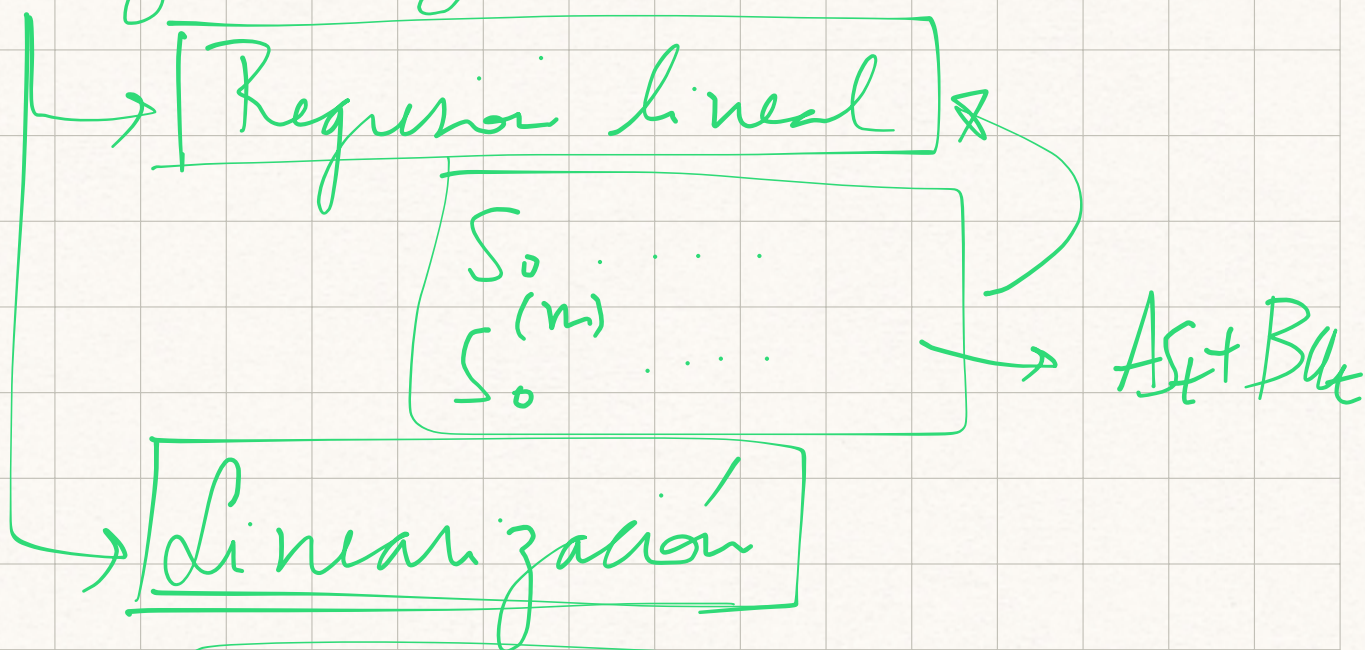$$\Rightarrow P_{sa} \rightarrow \boxed{S_{t+1} = A^{(t)} S_t + B^{(t)} a_t + \boxed{W_t}}$$

ruido

$$\Rightarrow W_t \sim N(0, \Sigma_w)$$

- - - - - Tarea

$$\Rightarrow \text{Recompensas Cuadráticas}$$

$$R(s, a) = \left(- \left(S^T U_t S_t + a_t^T W_t a_t\right.\right.$$

$$R(s,a) = -(\|s\|^2 + \|a\|^2)$$

¿Cómo se encuentran estas
matriz $A$ y $B$?

→ [ Regresión lineal ] ↗

$$\begin{matrix} S_0 \cdots \cdots \\ \qquad (m) \\ S_0 \cdots \cdots \end{matrix}$$

→ $As_t + Ba_t$

→ [ Linearización ]

$$\boxed{S_{t+1} = f(S_t, a_t)}$$

Cart Pole $\begin{pmatrix} X_{t+1} \\ \dot{X}_{t+1} \\ \Theta_{t+1} \\ \dot{\Theta}_{t+1} \end{pmatrix} = f\left( \begin{pmatrix} X_t \\ \dot{X}_t \\ \Theta_t \\ \dot{\Theta}_t \end{pmatrix}, a_t \right)$

Series de Taylor

$S_{t+1}$ ↑ , $f$

$$S_{t+1} = f(S_t)$$

$$\boxed{f(S_t) \approx f(\bar{S}_t) + f'(\bar{S}_t)(S_t - \bar{S}_t)}$$

caso
gral

$$S_{t+1} = f(S_t, a_t)$$

$$S_{t+1} \approx \left[ f(\bar{S}_t, \bar{a}_t) \right.$$
$$+ \left( \nabla_S f(\bar{S}_t, \bar{a}_t)(S_t - \bar{S}_t) \right.$$
$$+ \left. \left( \nabla_a f( \qquad ) (a_t - \bar{a}_t) \right] \right.$$