

1. Recap (MDP,  $V^*$ ,

Iteración de Valor)

2. Iteración de Política

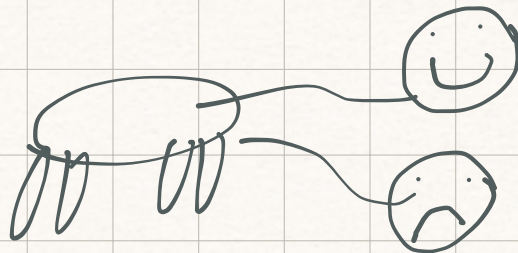
3. MDP Estados Continuos

\* Tarea 1 (PDF)  $\Leftarrow$

Recap

RL  $\rightarrow$  AA, helicóptero, carrito  
con pole

Dinámica  $\rightarrow$  Control  
RL



• AlphaGo  
• Ajedrez  
• AlphaStar

12L  
bames

# función de Recompensas

$$R(s) = \begin{cases} +1 & \text{(para ganar)} \\ -1 & \text{(para perder)} \\ 0 & \text{(el resto)} \end{cases}$$

recomp. (rewards)      estado

$$R(3,2) = 0$$

$$R(4,3) = +1$$

## MDP

$$\text{MDP} \{S, A, P_{sa}, \gamma, R\}$$

$S \rightarrow$  set de estados

$A \rightarrow$  set de acciones

$P_{sa} \rightarrow$  Prob. de transición

$\gamma \rightarrow$  factor de descuento  $\gamma \in [0, 1)$

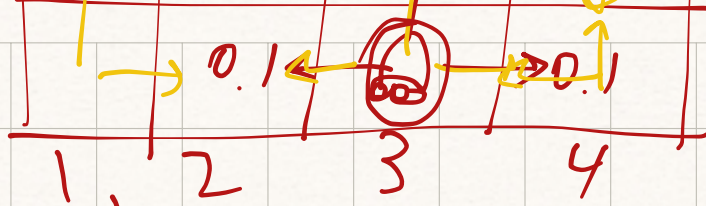
$R \rightarrow$  función de recompensas  $\gamma = 0.9$   
 $0.99$   
 $0.8$

$$P(s')_{sa}$$



recomp. inmediata vs recomp.





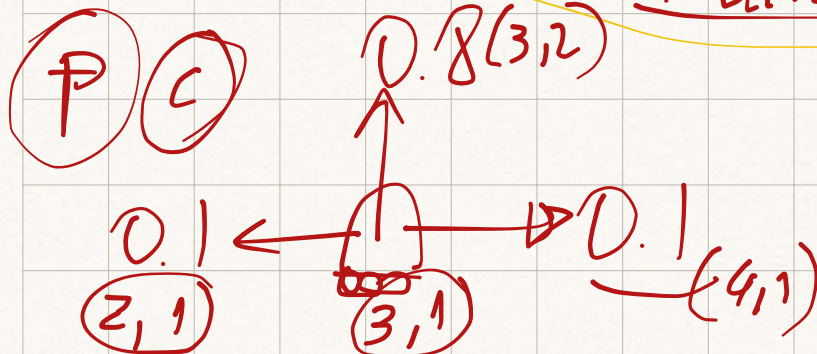
future

$$R_0 + \gamma R_1 + \gamma^2 R_2 + \dots$$

5) 11 estados (pos.)

A) Acciones  $\{S, N, E, O\}$   
4-dir

$S = \{-, -, -, -\}$   
 $N = \uparrow$   
 $E = \rightarrow$   
 $O = \leftarrow$



$$P_{(3,1)N}((3,2)) = 0.8$$

$$P_{(3,1)E}((4,1)) = 0.1$$

$$P_{(3,1)O}((2,1)) = 0.1$$

$$P_{(3,1)}((3,3)) = 0$$

d)  $R(\underline{(4,3)}) = +1$

$R(\underline{(4,2)}) = +1$

$$R(9, 1) = -1$$

$$R(s) = 0$$

$\forall s$

Algorithm

$S_0$  ('despertar')



$a_0$



$S_1 \sim P_{S_0 a_0}$



$a_1$



$S_2 \sim P_{S_1 a_1}$

¿Cuál es el pago total?  
(ganancia)

$$\Rightarrow \underline{R(S_0) + \gamma R(S_1) + \gamma^2 R(S_2) + \dots + \gamma^n R(S_n)}$$

$$\gamma \cong 0.9$$



# Objetivo de un algo de RL

Arrojar la acción a través del tiempo para maximizar la sumatoria de recompensas.

$$V = \max_a \left( E \left[ R(s_0) + \gamma^1 R(s_1) + \gamma^2 R(s_2) + \dots \right] \right)_{\pi}$$

Política ( $\pi$ ) mapea acción para cada estado

$$\pi : S \rightarrow A$$

(controlador)

Política dada  $\pi$  ( $\pi^*$ )

→	→	→	+1
↑	<del>↘</del>	↑	-1
↑	←	←	←

$$\pi(3,1) = \text{"←"}$$

Estando en el estado  $s$ ,  
tome acción  $\pi(s)$

Obj.  $\rightarrow$  Maximizar el total  
de recompensa esperada.

### Conclusiones

MDP  $\rightarrow$  formula el problema

Algo. de RL  $\rightarrow$  Encontrar la  
política  $\pi$  que  
maximice la  
ganancia esperada