

# Reinforcement Learning (RL)

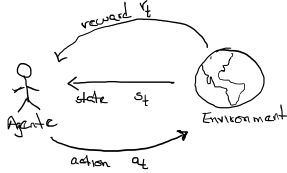
## Intro Deep Q-Learning

RL: Enfoque matemático para aprender toma de decisiones mediante ensayo y error

No pensamos en términos  $f: x \rightarrow y$   $\times$

Si pensamos en términos  $f: \max_{\pi} \pi(R|S \times A)$ , donde  $\pi$ : policy

$R$ : recompensa a largo plazo  
 $r_t$ : recompensa inmediata en el tiempo  $t$



## Terminología:

- Markov decision process (MDP)
- Q learning
- Policy gradient
- Deep Q Network (DQN)
- Bellman Equation
- Exploration vs exploitation
- Ambiente
- Agente  $a_t$
- Estado  $s_t$
- Acción
- Recompensa
- Episodio

$$s_t \in S$$

$$a_t \in A$$

$$r_t \in R$$

$$\pi(s_t, a_t) \rightarrow r_t$$

$$\text{Evaluar } \pi^* = \arg \max_{\pi} \pi(R|S \times A)$$

$$\text{donde } R_t = \sum_{k=0}^{\infty} \gamma^k r_{t+k+1}, \quad 0 \leq \gamma \leq 1$$

## Q-Learning

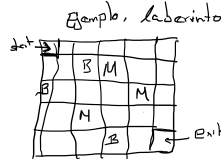
Uno de varios tipos de RL

- Construimos Q-table
- Q: quality

Q-table

states	0	1	2	3	4
1	0	0	0	0	0
2	0	0	0	0	0
3	0	0	0	0	0
4	0	0	0	0	0
5	0	0	0	0	0

matriz de transición  
(probabilidad de realizar alguna acción dado un estado)



$$S = \begin{matrix} s_0 = (1,1) \\ s_1 = (1,2) \\ s_2 = (1,3) \\ s_3 = (1,4) \\ \vdots \\ s_{29} = (5,5) \end{matrix}$$

$$A = \begin{matrix} a_1 = \rightarrow 1 \\ a_2 = \leftarrow 2 \\ a_3 = \uparrow 3 \\ a_4 = \downarrow 4 \end{matrix}$$

## función de calidad

$$q_{\pi}(s, a) = E[R_t | s_t, a_t]$$

$$= E\left[\sum_{k=0}^{\infty} \gamma^k r_{t+k+1} | s_t, a_t\right]$$

$$q^*(s, a) = \max_{\pi} q_{\pi}(s, a)$$

## Bellman Ecuación

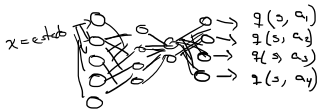
$$q^*(s, a) = E[r_t + \gamma \max_{a'} q^*(s', a')]$$

$s'$  = estado en  $t+1$ , después de tomar acción  $a$  sugerida por  $q^*$

$a'$  = acción en  $t+1$  después de  $q^*$

De manera iterativa actualizar tabla Q emulando estrategias "Programación lineal"

## Reemplazar Tabla Q con una Red Neuronal

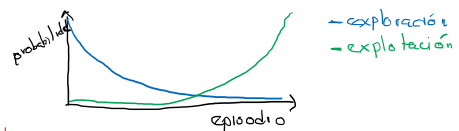


- Backprop usa recompensas para ajustar pesos
- Debemos encontrar la manera de considerar recompensas futuras (Bellman)

## Trade-off exploración y explotación

- Exploración: tomar acciones al random y calcular recompensas inmediatas
- Explotación: tomar acciones usando conocimiento adquirido hasta ahora

$$\alpha_t = f(\alpha_{\text{explor}} + (1 - \alpha_{\text{explor}}) \alpha_{\text{explot}})$$



## Deep Q Learning

$$\text{loss} = q^*(s, a) - q(s, a), \quad \leftarrow \text{usado por backprop}$$

donde:  $q(\cdot)$ : valor de calidad (recompensa) inmediata

$q^*(\cdot)$ : valor de calidad de la mejor acción después del siguiente estado.

## - Replay memory (experience replay)

• Memoria (liste) con experiencias (tuples):  $e_t = (s_t, a_t, r_t, s_{t+1})$