



# **Aprendizaje Reforzado**

## **Maestría en Ciencia de Datos, DC - UBA**

Julián Martínez  
Javier Kreiner



# Plan de la clase

## Teoría:

- Cadenas de Markov
- Procesos de Decision de Markov
- Procesos de Recompensa de Markov
- Función de Valor
- Funcion de Acción-Valor
- Ecuación de Esperanza de Bellman
- Ecuación de Optimalidad de Bellman
- Evaluación de una política
- Iteración de política
- Iteración de (función de) valor

# Markov Reward Processes

  
*Función de recompensa*

$$\mathcal{R}_s := E[R_{t+1} | S_t = s]$$

*Retorno*

$$G_t := R_{t+1} + \gamma R_{t+2} + \dots$$

*Función de Valor*

$$v(s) = E[G_t | S_t = s]$$

## Ecuación de Bellman



$$v(s) = \mathcal{R}_s + \gamma \sum_{s'} p_{s,s'} v(s')$$

Un sistema de K ecuaciones!

$$v = (v(s_1), \dots, v(s_K))$$

$$v = \mathcal{R} + \gamma p v$$

# Proceso Markovianos de Decisión

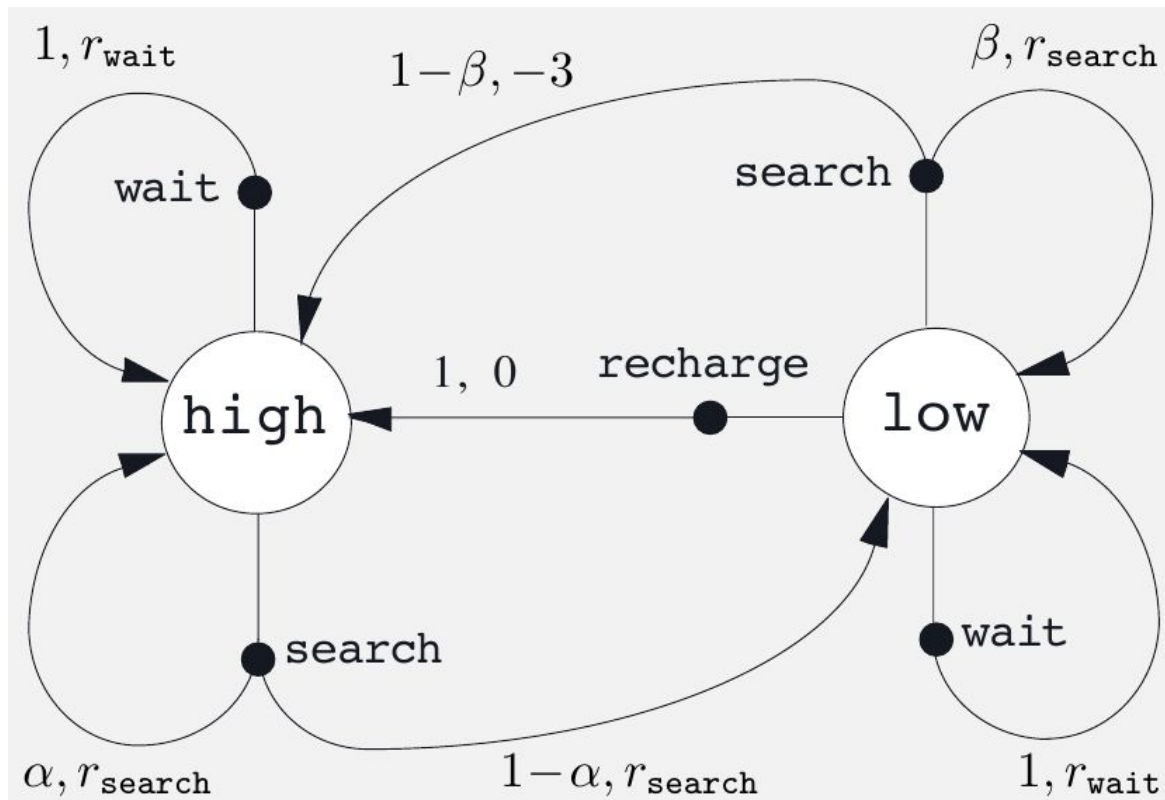


Se agrega un *acción* la cual modifica el ambiente.

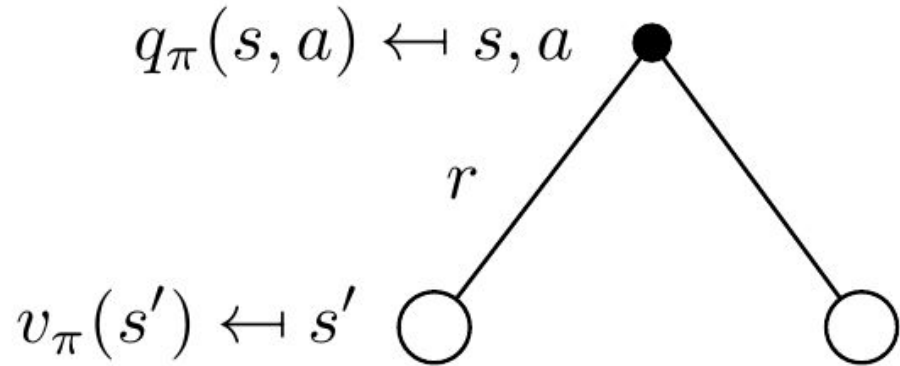
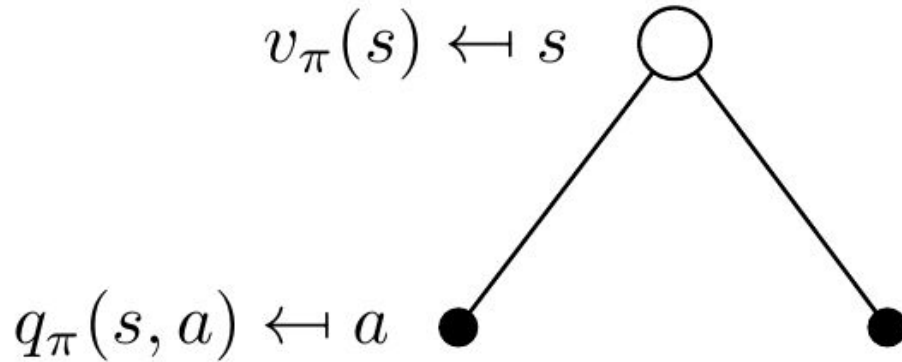
$$p_{s,s'}^a := P(S_{t+1} = s' | S_t = s, A_t = a)$$

$$\mathcal{R}_s^a := E[R_{t+1} | S_t = s, A_t = a]$$

# Un ejemplo

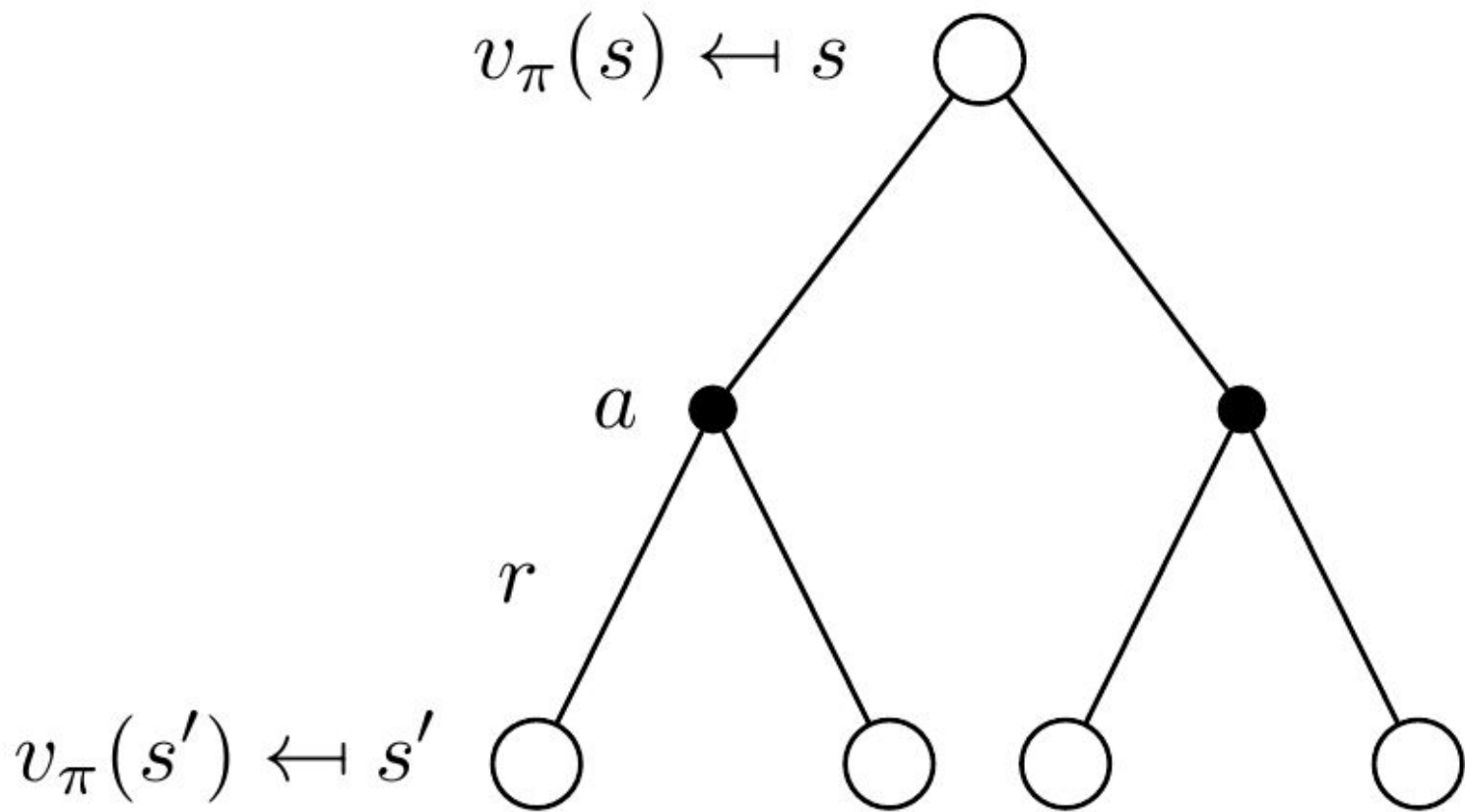


## Ecuación de Bellman (bis)



$$v_\pi(s) = \sum_{a \in \mathcal{A}} \pi(a|s) q_\pi(s, a)$$

$$q_\pi(s, a) = \mathcal{R}_s^a + \gamma \sum_{s' \in \mathcal{S}} \mathcal{P}_{ss'}^a v_\pi(s')$$







# Introducción a OpenAI Gym



## Lectura recomendada:

- Desafíos de Reinforcement Learning: <https://www.alexirpan.com/2018/02/14/rl-hard.html>



# Competición MIT

- Competición en MIT: <https://selfdrivingcars.mit.edu/deeptraffic/>