

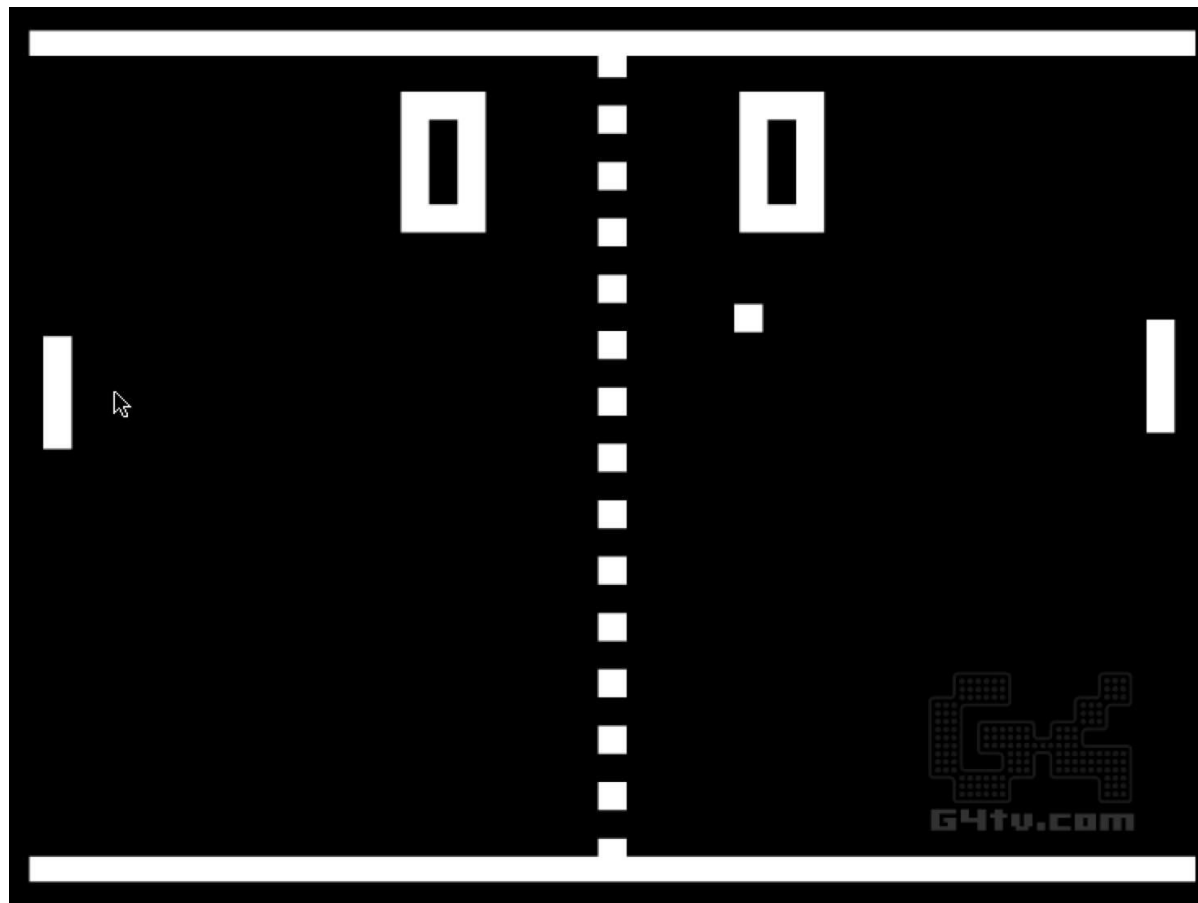


# PONG

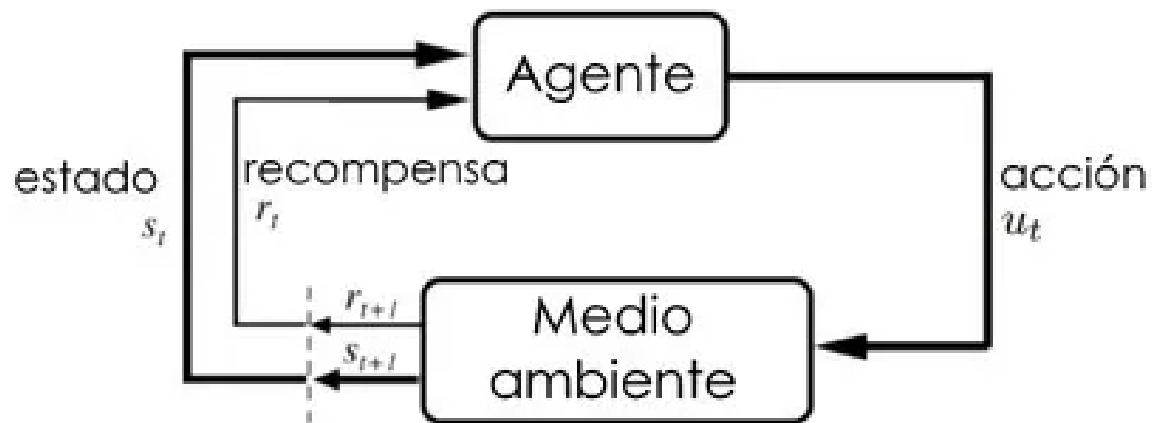
APLICACIÓN DE Q-LEARNING EN EL JUEGO PONG

Ribot, santiago  
ribot@frba.utn.edu.ar

# PONG



# APRENDIZAJE POR REFUERZO



[Figure source: Sutton & Barto, 1998]

Repetir

Inicializar Q-Table

Elegir una acción

Realizar la acción

Obtener una recompensa

Actualizar la Q-table

# APRENDIZAJE POR REFUERZO

$$\hat{Q}(s,a) = Q(s,a) + \alpha \left[ R + \left( \lambda \max_{s'} Q(s',a) \right) - Q(s,a) \right]$$

Diagram illustrating the components of the Bellman equation for Q-learning:

- $\hat{Q}(s,a)$ : valor actual
- $Q(s,a)$ : valor actual
- $\alpha$ : ratio aprendizaje
- $R$ : recompensa
- $\lambda$ : tasa descuento
- $\max_{s'} Q(s',a)$ : valor óptimo esperado
- $Q(s,a)$ : valor actual

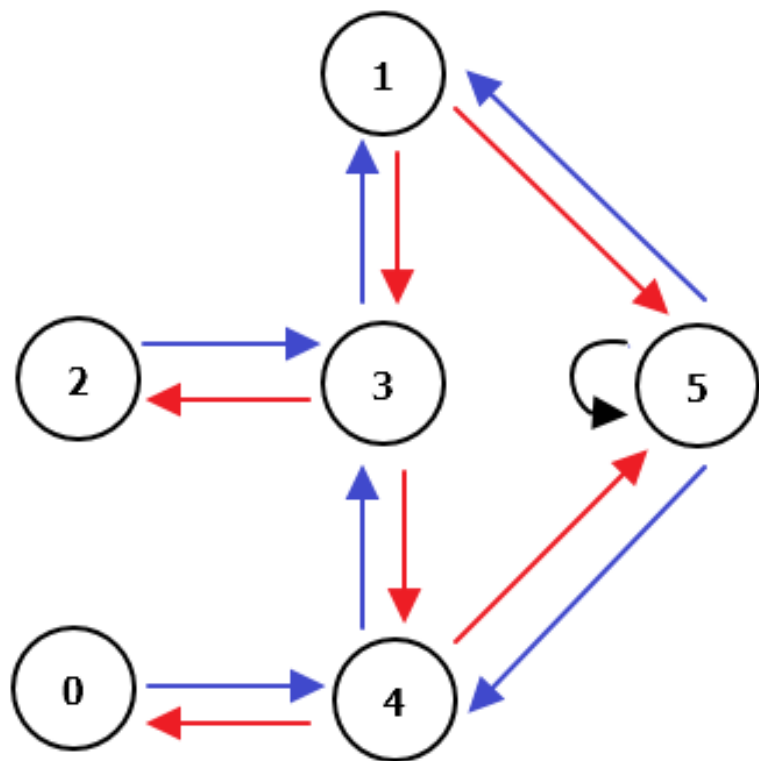
## Exploración vs Explotación:

- Ratio de exploración: (0, 1)

## Ecuación de Bellman:

- Factor de descuento: (0, 1)
- Ratio de aprendizaje: (0, 1)
- Recompensa: cualquier valor

# APRENDIZAJE POR REFUERZO - EJEMPLO

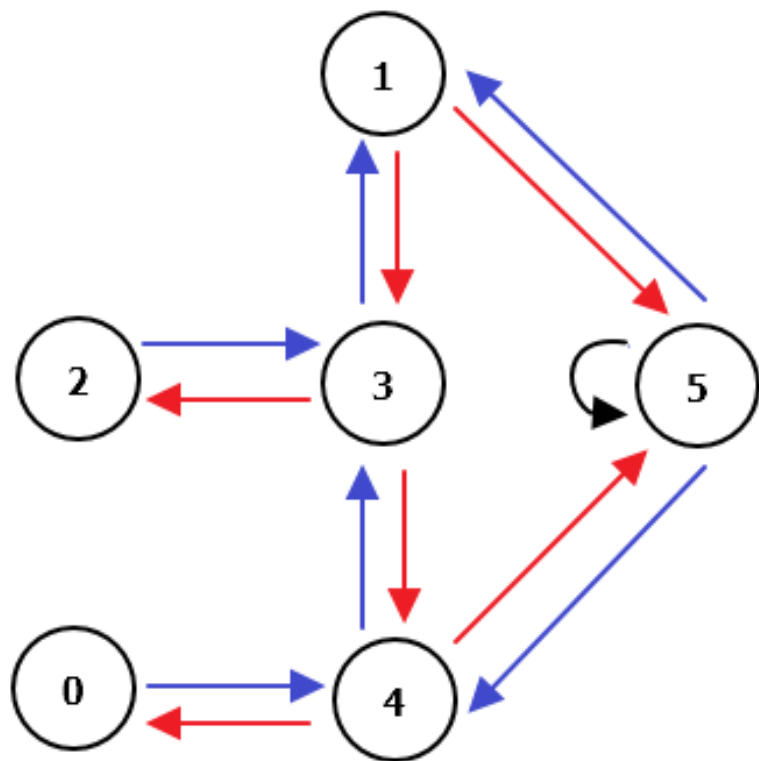


Laberinto de 6 etapas:

- Al llegar a la etapa 5 se gana
- El agente puede iniciar en cualquier etapa
- Los caminos son únicamente los indicados

	0	1	2	3	4	5
0	-1	-1	-1	-1	0	-1
1	-1	-1	-1	0	-1	100
2	-1	-1	-1	0	-1	-1
3	-1	0	0	-1	0	-1
4	0	-1	-1	0	-1	100
5	-1	0	-1	-1	0	100

# APRENDIZAJE POR REFUERZO - EJEMPLO



Actualización del valor  $Q$ : Ecuación de Bellman

$$Q'(s, a) := Q(s, a) + \alpha(r + \gamma \max_{a'} Q(s', a') - Q(s, a))$$

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

# EL JUEGO

## Reglas:

- El agente tiene 3 vidas
- Cantidad máxima de puntos: 200
- Si se cae la pelota pierde 10 puntos
- Si golpea la pelota gana 10 puntos
- Si pierde las 3 vías pierde 20 puntos

## Clase agent:

- Almacena la Q-table
- Tiene las funciones para actualizarla
- Decide que acción tomar a continuación

## Clase enviroment:

- Implementa la lógica del juego y las animaciones
- Da las recompensas o castigos
- Aplica la acción del agente y devuelve un nuevo estado

# Q-TABLE

- Score 70 – Lives: 3 -



El estado queda fijado por tres componentes:

- Posición x de la pelota
- Posición y de la pelota
- Posición de la paleta (solo y)

ACCION:   
RECOMPENSA: 0

$ESTADO(x) = ESTADO(x) - 1$   
 $ESTADO(y) = ESTADO(y) - 1$   
 $ESTADO(p) = ESTADO(p) + 1$



# Q-TABLE

- Score 70 – Lives: 3 -



Para el caso de ejemplo se utiliza una pantalla de 50px x 40px:

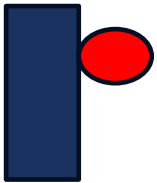
- Cada step mueve de a 5px
- 10 posiciones en X
- 8 posiciones en Y
- 8 posiciones de la paleta

ACCION:   
RECOMPENSA: 0

$ESTADO(x) = ESTADO(x) - 1$   
 $ESTADO(y) = ESTADO(y) - 1$   
 $ESTADO(p) = ESTADO(p) - 1$

# Q-TABLE

- Score 80 – Lives: 3 -



ACCION:   
RECOMPENSA: 10

Además, la acción  o 

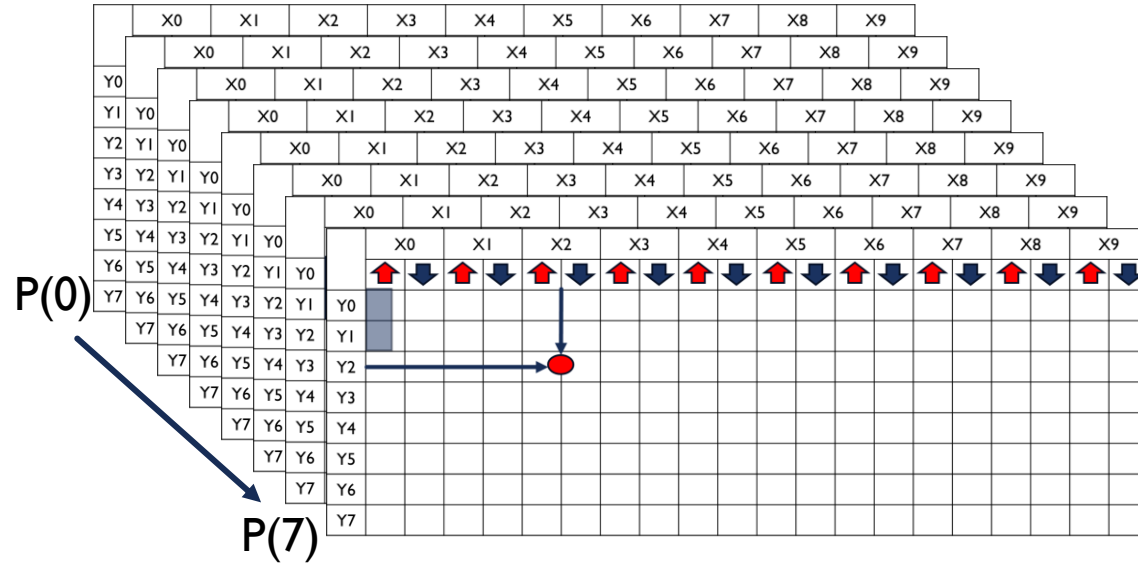
- Es lo que modifica el estado P
- Los otros dos estados se modifican con la actualización de cada frame del juego

$ESTADO(x) = ESTADO(x) + 1$   
 $ESTADO(y) = ESTADO(y) - 1$   
 $ESTADO(p) = ESTADO(p) + 1$

# Q-TABLE

[illegible]

# Q-TABLE



ESTADO[2][2][7][UP] = Q\_VALUE

Q\_VALUE = Q\_VALUE +  $\alpha(R + \lambda * \text{MAX}(\text{ESTADO}[I][I][6]) - Q\_VALUE)$

# CLASE AGENT

```
JS main.js JS agent.js X JS enviroment.js
public > js > JS agent.js > Agent
1  export default class Agent {
2  >   constructor(enviroment, policy, discountFactor = 0.1, learningRate = 0.1, explorationRatio = 0.9) { ...
15  }
16
17  >   getShape(policy) { ...
25  }
26
27  >   createQTable(position) { ...
46  }
47
48  >   getNextStep(state, enviroment) { ...
55  }
56
57  >   randomChoice(opt, state) { ...
76  }
77
78  >   update(enviroment, oldState, action, reward, newState, end) { ...
93  }
94
95  >   printPolicy() { ...
110 }
111
112 >   getPolicy() { ...
114 }
115 }
```

```
constructor(enviroment, policy, discountFactor = 0.1, learningRate = 0.1, explorationRatio = 0.9) {

    if(policy)
        this.q_table = policy;
    else{
        this.position = this.getShape(enviroment.policy);
        this.position.push(2);
        this.q_table = this.createQTable(this.position);
    }

    this.discountFactor = discountFactor;
    this.learningRate = learningRate;
    this.explorationRatio = explorationRatio;
}
```

# CLASE ENVIROMENT

```
JS main.js JS agent.js JS enviroment.js X
public > js > JS enviroment.js > Enviroment
1 import {width_px, height_px, max_lives, max_score, mov} from "../main.js";
2
3 export default class Enviroment {
4   > constructor(myDiv) { ...
36   }
37
38   > reset() { ...
46   }
47
48   > async step(action, animate = false) { ...
56   }
57
58   > async applyAction(action, animate) { ...
74   }
75
76   > advancePlayer() { ...
82   }
83
84   > advanceFrame() { ...
116   }
117
118   > async plotFrame() { ...
203   }
204
205   > detectColission() { ...
209   }
210 }
```

```
constructor(myDiv) {
  this.actions = ['up', 'down'];
  this.state = [0, 0, 0];
  this.totalScore = 0;

  this.policy = [];
  for (let x = 0; x < height_px/mov; x++) { ...
  }

  this.lives = max_lives;
  this.penalty = 0;
  this.paddleHeight = height_px/4
  this.score = 0;
  this.player = height_px/2

  this.x = width_px/2
  this.y = height_px/2
  this.dx = mov
  this.dy = mov
  this.radio = 2.5;
  this.myDiv = myDiv;

  this.interval = null;
}
```

# PLAY

```
for(let playedGames = 0; playedGames < rounds; playedGames++) {
  let state = enviroment.reset();
  let reward = null;
  let done = null;
  let iteration = 0;
  let contador = 0;

  while(done !== true && iteration < 3000 && enviroment.totalScore < max_score) {
    let oldState = state.slice();
    let nextAction = agent.getNextStep(state, enviroment);

    ({state, reward, done} = await enviroment.step(nextAction, animate));

    if(animate){
      await enviroment.plotFrame();
    }

    if(rounds > 1) {
      agent.update(enviroment, oldState, nextAction, reward, state, done)
    }
    iteration += 1;
  }

  totalReward += enviroment.totalScore;

  if(enviroment.totalScore > max) {
    max = enviroment.totalScore;
    firstMax = playedGames;
  }

  if(playedGames%500 === 0 && playedGames > 1 && !animate) {
    contador += 1;
    console.log(`-- Games: [${playedGames}] -- AvgScore: [${(totalReward/(500*contador)).toFixed(2)}] -- AvgSteps: [${iteration}] -- MaxScore: [${max}] -- GameOfFirstMax: [${firstMax}]`);
    max = 0;
    firstMax = 0;
    totalReward = 0;
  }
}
```



GRACIAS