

Universidad Nacional del Sur  
Departamento de Cs e Ingeniería de Computación  
Proyecto Final

# Sistemas de trading y Deep Reinforcement Learning

Autor:

Leonardo Jose Caramello

Director:

Diego Martinez

Bahia Blanca, Buenos Aires

Agosto - 2017



# Resumen

Vivimos en un mundo cada vez más digitalizado donde el acceso a la información es cada vez más fácil y abundante, un mundo que en los últimos años se ha transformado radicalmente y que sin duda lo seguirá haciendo en los años venideros. Los avances en machine learning, han permitido automatizar tareas, que hasta hace algunos años parecía imposible, ejemplos de esto son Google self-driving car, sistemas de recomendación como los usados por netflix o mercado libre, o Deep Mind.

Uno de los desafíos más interesantes en este área, es el desarrollo de sistemas de trading que permitan automatizar la comercialización de activos financieros, con el fin de administrar eficientemente un porfolio de inversiones. El siguiente trabajo esta inspirado en Deep Mind, un agente que combinando reinforcement learning y deep learning aprende a jugar juegos de atari. En este papper buscamos investigar la aplicabilidad y efectividad de las estas técnicas, en el desarrollo de sistemas de trading.

# Índice general

|   |           |
|---|-----------|
| Índice de figuras                                     | v         |
| Índice de cuadros                                     | v         |
| Nomenclatura  | vii       |
| <b>1 Introducción</b>                                 | <b>1</b>  |
| 1.1 Alcance del proyecto . . . . .                    | 1         |
| 1.2 Objetivos . . . . .                               | 2         |
| 1.3 Metodología . . . . .                             | 2         |
| 1.4 Contenido . . . . .                               | 3         |
| <b>2 Mercados Financieros</b>                         | <b>5</b>  |
| 2.1 Introducción . . . . .                            | 5         |
| 2.2 Análisis Técnico y Análisis Fundamental . . . . . | 5         |
| 2.3 Teoría Clásica . . . . .                          | 6         |
| 2.4 Fundamentos del análisis técnico . . . . .        | 8         |
| <b>3 Reinforcement Learning y Deep Learning</b>       | <b>9</b>  |
| 3.1 Introducción . . . . .                            | 9         |
| 3.2 Configuración de estados . . . . .                | 10        |
| 3.3 Acciones . . . . .                                | 11        |
| 3.4 El Agente . . . . .                               | 12        |
| 3.5 Q-Network . . . . .                               | 13        |
| 3.6 Algoritmo de aprendizaje . . . . .                | 14        |
| <b>4 Arquitectura y Diseño</b>                        | <b>15</b> |
| <b>5 Evaluación y Desempeño</b>                       | <b>17</b> |
| <b>6 Conclusiones y Recomendaciones</b>               | <b>19</b> |
| <b>Bibliografía</b>                                   | <b>21</b> |

## Índice de figuras

|     |   |    |
|-----|---|----|
| 3.1 | Reinforcement Learning Architecture Overview. . . . . | 10 |
| 3.2 | Gráfico de vela. . . . .                              | 10 |

## Índice de cuadros



# Nomenclatura

|              |   |
|--------------|---|
| $a$          | razón de constantes de tiempo del proceso (modelo).           |
| $\beta$      | factor de peso del valor deseado del controlador.             |
| $C(s)$       | función de transferencia del controlador.                     |
| $C_r(s)$     | función de transferencia del controlador de valor deseado.    |
| $y(t), y(s)$ | ... <i>La nomenclatura debe listarse en orden alfabético.</i> |
| $\zeta$      | Este es un ejemplo de una lista con “descripción”.            |





# 1 Introducción

Este documento forma parte del proyecto final de carrera y se acompaña junto con el código del agente desarrollado, disponible también online. En este documento se pretende dar una presentación formal a los resultados de investigación del proyecto, como así también una descripción del problema a resolver y de la solución adoptada, detallando cuales fueron cada una de las decisiones de diseño adoptadas y por que se adoptaron.

## 1.1 Alcance del proyecto

El proyecto se planteo como un trabajo de investigación que permitiera analizar cuales son las herramientas que brinda el machine learning para desarrollar un agente, capaz de invertir en activos financieros (acciones, bonos, obligaciones negociables, etc), en particular, a través del uso de reinforcement learning.

El proyecto pretende ser una prueba de concepto, que permita abordar el desarrollo de sistemas de trading utilizando reinforcement learning. En particular, se buscara identificar cada uno de los componentes que plantea el framework de reinforcement learning, de forma tal, de tratar de encontrar una representacion adecuada.

A su vez, durante el proceso de investigación se tratara de comprender la complejidad y las problemáticas propias del dominio del problema, en este caso, los mercados financieros y el trading en ellos, las cuales sera necesario comprender, para poder modelar una solución mas eficiente.

Queda fuera del alcance de este proyecto lo siguiente:

- Desarrollar un nuevo algoritmo de RL
- Desarrollar un agente capaz generar ganancias
- Desarrollar una comparación de diferentes implementacion del RL.
- Cubrir todas las peculiaridades de un mercado financiero.

## 1.2 Objetivos

A continuación se detallan los objetivos del proyecto

### Objetivo general

Investigar la posible aplicabilidad de reinforcement learning en el desarrollo de sistemas de trading que permitan optimizar y automatizar la toma de decisiones de un potencial inversor.

### Objetivos específicos

- Brindar una descripción general del funcionamiento de los mercados financieros.
- Brindar una descripción general de Reinforcement Learning.
- Brindar una especificación detallada de cómo modelar el problema de trading utilizando los elementos que propone el framework de RL.
- Desarrollar un entorno que permita realizar la simulación de un mercado financiero
- Desarrollar un agente inteligente capaz de percibir su entorno y tomar decisiones de compra o venta de un activo financiero

## 1.3 Metodología

En primer lugar se definirá el tópico central de la investigación, junto con los alcances y limitaciones del proyecto.

En una segunda etapa, se hará una investigación sobre el funcionamiento de los mercados financieros, en particular, se buscará entender el funcionamiento del análisis técnico de los mercados y/o de los principales indicadores bursátiles, de forma tal de poder capturar estos conceptos y poder modelarlos en el desarrollo del agente.

En tercer lugar se llevará a cabo una investigación de los diferentes algoritmos de RL y de cómo modelar el problema de trading utilizando los elementos que propone el framework.

Por último, se proseguirá con el desarrollo del agente inteligente, de forma tal de poder integrar todo el conocimiento adquirido en las etapas previas.

## 1.4 Contenido

Por último, concluimos este capítulo introductorio con una breve descripción de cada uno de los capítulos siguientes

- Capítulo 2 - Mercados Financieros: En este capítulo se presenta el problema del trading, se dará una breve explicación del funcionamiento de un mercado financiero y de como operan los inversores en el, abarcado conceptos como acciones, precio, análisis técnico vs análisis fundamental, drivers, indicadores bursátiles, tendencias, etc.
- Capítulo 3 - El Problema: En este capítulo se introducirá brevemente el framework de reinforcement learning y cada uno de sus componentes, para luego esbozar la representación que se adoptara junto con un pseudo algoritmo que mostrara una primera aproximación a la solución elegida
- Capítulo 4 - Arquitectura y Diseño: Aquí se hará una descripción detallada de la arquitectura y el diseño del sistema implementado, se mostraran detalles y demás peculiaridades de implementación, en particular parámetros de configuración, funcionalidades, etc. Se mostraran algunas partes mas relevantes del código del agente.
- Capítulo 5 - Evaluación y desempeño: En esta sección analizaremos el desempeño y eficacia de la solución implementada, mostrando y analizando los diferentes resultados obtenidos durante la ejecución del agente.
- Capítulo 6 - Conclusiones y Recomendación: Como punto final, esta sección estará destinada a comentar diferentes conclusiones que pudimos establecer, junto con algunas recomendaciones y/o observaciones de posibles mejoras a futuro.
- Capítulo 7 - Bibliografía y Referencias: Listado de las diferentes referencias que formaron parte de la investigación.



## 2 Mercados Financieros

### 2.1 Introducción

En economía, un mercado financiero es un espacio (físico, virtual o ambos) en el que se realizan los intercambios de instrumentos financieros y se definen sus precios. Los mercados financieros están afectados por las fuerzas de oferta y demanda. La clave del éxito está en saber predecir el futuro y actuar en consecuencia. Quedarse largo en una posición (comprar) si se piensa que el mercado va a subir, o deshacer posiciones o quedarse corto si se piensa que el mercado va a bajar. Si se consigue hacer esto en forma reiterada se podrán obtener ganancias y con ello incrementar el capital inicial.

### 2.2 Análisis Técnico y Análisis Fundamental

Para intentar saber cómo va a estar un valor en el futuro, se distinguen tradicionalmente dos corrientes bien diferenciadas: los que siguen el análisis técnico y los que siguen el análisis fundamental. Los fundamentales se basan en que el valor de una acción está dado por los beneficios futuros de la empresa. Ni más, ni menos. Lo que intentan es determinar cuáles serán esos beneficios futuros, y para ello tratan de conocer diferentes detalles de la empresa: noticias que les afecten, posibles movimientos societarios, estrategias, competidores, nuevos productos, etc. Toda la información micro económica tiene impacto en dichos beneficios futuros. Así como también la macro económica: cómo evoluciona el entorno general de la empresa, el entorno regulatorio, el entorno político, etc. Se trata, en definitiva, de analizar la mayor cantidad de información posible, y de convertir esa información en cuentas de resultados provisionales que se puedan descontar para hallar el valor actualizado de la acción. El análisis técnico, por el contrario, se basa en que el precio de la acción lo descuenta todo. Es decir, todos los factores relevantes a la inversión, cualquiera que ellos sean, pueden ser reducidos al nivel de precios de la acción y volumen transado. El precio de mercado representa el total conocimiento de los inversionistas respecto de cualquier activo dado en un momento particular. Además, refleja todas las noticias sobre el mercado así como la suma de conocimientos de los participantes en éste. Aquí se habla de tendencias alcistas o bajistas, de líneas de soporte (cotizaciones donde se cree que la acción dejará de bajar y tendrá un rebote”), líneas de resistencia (cotizaciones en las que el valor de

la acción se atascará y que le costará romper”), etc. La lógica nos dice que el análisis fundamental es el que tiene más sentido. Sin embargo, en la realidad esto no siempre es así. Lo cierto es que cuanto más gente crea en los análisis técnicos, más probabilidad tendrán de ser reales sus predicciones (ya que la gente actuará como si fueran reales, contribuyendo a su efectiva realización).

Este proyecto se basa fundamentalmente en las ideas del análisis técnico y las herramientas que este brinda, las cuales, serán los pilares fundamentales sobre los que el agente tomara sus decisiones. Para comprender un poco mas acerca de este y visualizar como es posible tomar decisiones acertadas solamente apoyándonos en el análisis técnico, debemos mirar mas en detalles, algunos aspecto de la teoría clásica financiera.

## 2.3 Teoría Clásica

Gran parte de la teoría financiera clásica parte del principio fundamental que los inversores son racionales y que los precios del mercado reflejan en todo momento y de manera instantánea el valor fundamental de los títulos. Este principio fundamental establece que la competencia entre los distintos participantes que intervienen en el mismo, conduce a una situación de equilibrio en la que el precio de mercado de un activo constituye una buena estimación de su precio teórico, es decir, que los precios que se negocian en el mercado reflejan toda la información existente y se ajustan total y rápidamente a los nuevos datos que puedan surgir. La consecuencia de este principio es que un inversor racional no puede hacer nada para “ganar” al mercado.

De acuerdo con este principio de racionalidad económica, lo que debe hacer un inversor es intentar maximizar su riqueza final. Para lograr esto, lo mejor que puede hacer este inversor racional es invertir en el mercado de manera diversificada de una forma igual a la del mercado y permanecer en esta misma cartera salvo por necesidades de liquidez o debido a variaciones de su situación actual o de cambio en sus necesidades futuras.

A pesar de la cantidad de libros de texto y de artículos que sostienen los principios anteriores sobre la forma en que deben comportarse los inversores, lo cierto es que la evidencia empírica nos dice que las cosas no suceden de la forma en la que deberían suceder según este principio, o por lo menos, no enteramente. Para poder explicar este fenómeno deberemos detenernos en algunos aspectos fundamentales en el proceso de decisión de los inversores.

Un claro ejemplo de esto podemos verlo si analizamos el proceso de decisión de venta, de acuerdo con la teoría clásica, los precios de cualquier activo si-

guen un movimiento aleatorio. Esto significa que la mejor predicción sobre el precio futuro es la que se tiene hoy. La consecuencia inmediata de esto es que no tiene ningún sentido vender, para a continuación, volver a comprar éste mismo activo u otro diferente. Dado que la expectativa de ganancia debido a la diferencia entre los precios venta y de compra sería nula. Solo tendríamos una pequeña pérdida debido al coste de la transacción. En otras palabras, el mercado no es predecible y, en consecuencia, no es posible obtener un beneficio realizando trading. Sin embargo, veamos algunos conceptos que contradicen esto:

- *La creencia de los inversores en la reversión a la media*, es el principio según el cual existe un valor medio de cada acción al cual se acaba volviendo en algún momento. Así, si un valor tiene un precio que el inversor cree que está por debajo del que le corresponde (su valor “medio”) , tarde o temprano, el precio de esta acción subirá hasta llegar a ese precio. Y, en consecuencia, recuperará las pérdidas que está teniendo en este momento. Lo mismo puede decirse cuando el valor está por encima del valor medio. Así, si el inversor tiene una serie de valores que entiende que están infravalorados por el mercado, tiende a mantenerlos esperando que vuelvan a su valor “medio”.
- *La aversión a la pérdida*. Si por alguna razón, tenemos alguna predicción confiable que nos dice que el precio de un activo va a bajar, lo racional es vender, independientemente de que con el activo a vender se obtenga un beneficio por su venta o no. La realidad nos muestra que esto no siempre es así, y que en general, los inversores tienden a no realizar las pérdidas, es decir, a no efectuar la venta real, en valores en los que pierden y, en cambio, vender antes de tiempo aquellos en los que tienen ganancias. Una gran parte de los inversores venden valores ganadores y mantienen los perdedores.
- *El efecto disposición*. El comportamiento debido a la aversión a la pérdida es una parte de un comportamiento más general de los inversores, llamado efecto disposición según el cual, los inversores mantendrían demasiado tiempo activos en pérdidas y venderían demasiado pronto activos con ganancias. Esto se da debido a que los inversores son mucho más sensibles a las pérdidas que a las ganancias en el sentido de que las primeras influyen el doble que las segundas.
- *El efecto atención*. Los inversores tienden a invertir en aquellos valores que llaman más su atención por la razón que sea, incluso aunque esta atención sea debida a noticias negativas, y como consecuencia, gran parte de los inversores diversifica mucho menos de lo que debería.

Estos argumentos y otros, nos indican que el inversor no siempre se comporta de manera esperada. Tal vez, el principal error de la teoría clásica sea partir de que los inversores son entes racionales y que no operan según sentimientos, euforia, miedo, avaricia o codicia.

## 2.4 Fundamentos del análisis técnico

Si bien es cierto que existen diversos argumentos en contra del análisis técnico, como por ejemplo:

- La profecía del auto cumplimiento
- El pasado no sirve para predecir el futuro
- El paseo aleatorio
- Mercados Eficientes

Lo interesante aquí es entender que el análisis técnico no trata de predecir el valor futuro del precio de un activo, sino mas bien de responder la siguiente pregunta:

*¿cómo cree el conjunto de inversionistas que evolucionará el precio en el futuro?*

Es importante notar, que lo que importa no es la evolución del precio, sino lo cree *la masa de inversores sobre la evolución del precio*, y ¿Por que esta pregunta es relevante?. Pues por que en el corto y mediano plazo, los precios se mueven mas por cuestiones psicológicas de los inversores que por variables financieras.

El análisis técnico nos dará las herramientas necesarias poder analizar desde un punto de vista estadístico y probabilístico, como fue el comportamiento de la masa de inversores en el pasado, ante una situación similar, para luego poder tomar la decisión mas conveniente.

En otras palabras lo que se busca es reducir el nivel de incertidumbre que tiene el inversor a la hora de decidir si compra o vende un activo, haciendo uso de diferentes herramientas matemáticas como las medias móviles, líneas de tendencias o patrones con el objetivo de determinar cual es el comportamiento mas probable de la masa de inversores.



## 3 Reinforcement Learning y Deep Learning

### 3.1 Introducción

En esta sección buscaremos dar una descripción mas detallada, tanto del problema a resolver, como así también de como se aplicaron los conceptos de reinforcement learning y deep learning a las soluciones adoptas durante el desarrollo del mismo. Esta sección no pretende ser una guía o una introducción a reinforcement learning, ni a deep learning, sin embargo es necesario contar con conocimientos básicos de ellos para poder comprender los detalles de la solución propuesta. Dejamos en manos del lector la capacitación en estos a través de los diferentes recursos que pueden hallarse en linea. (ver Capitulo 7 - Bibliografía).

Recordemos que el objetivo es desarrollar un agente inteligente, que a partir de un capital inicial, sea capaz de realizar compras y ventas de un activo financiero. El agente no poseerá ningún conocimiento previo acerca del funcionamiento del mercado, ni de la empresa sobre la que opera, ni ningún otro tipo de información. Solo conocerá el precio actual de la acción, y su evolución durante los días previos, junto un conjunto de indicadores bursátiles.

El agente va a interactuar con un simulador de un mercado bursátil, esta interacción se llevará a cabo siguiendo una serie de acciones, observaciones y recompensas. Cada interacción será llevada a cabo en episodios que tendrán una duración de  $m$  días

En cada instante  $t$ , el agente tendrá que seleccionar una acción  $a_t$  de un conjunto válido de acciones  $A = a_1, a_2, \dots, a_n$ . La acción será pasada al entorno, el cual modificará su estado interno y como respuesta a esta interacción el agente recibirá una recompensa  $R_{t+1}$ , la cual será calculada por el entorno. El entorno será estocástico, es decir, su comportamiento será no determinístico. Cada uno de los estados contendrá información relevante sobre el papel, como el precio, el volumen, y diferentes indicadores bursátiles, como medias móviles, medias móviles exponenciales, etc.

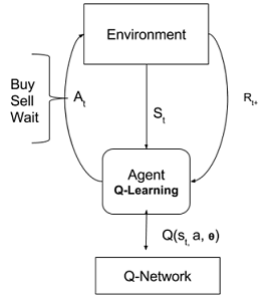


Figura 3.1: Reinforcement Learning Architecture Overview.

### 3.2 Configuración de estados

Ya sea de que se trate de un inversor experimentado o un agente inteligente, difícilmente sea posible tomar una decisión acertada sobre, la compra o venta de un activo, solamente mirando lo que sucede con el precio actual del activo, es necesario considerar una ventana de tiempo de tamaño  $n$ , para poder analizar la variación que ha tenido y así poder determinar la situación actual del activo, es decir, si se encuentra en una tendencia alcista o bajista, si se encuentra realizando una corrección de precio o no, si se encuentra en sobre compra o sobre venta, etc. Una manera eficiente de ver esto, es a través, de un gráfico de vela, donde cada vela, representa un periodo, el cual puede ser, un año, una semana, o un mes.



Figura 3.2: Gráfico de vela.

La representación que vamos a adoptar para cada uno de los estados que recibirá el agente, consiste en tomar esta idea de time frame y plasmarla en la estructura de los estados, es decir, el agente será capaz de ver una ventana de tiempo de tamaño  $n$  hacia atrás y de diferentes tipos de periodos.

Teniendo en cuenta esto, vamos a definir a cada estado  $S_t$  compuesto de 3 capas o layers, de la siguiente forma:

$$S_t = L_{días}, L_{semanas}, L_{meses} \quad \text{donde} \quad L_i = P_1, P_2, \dots, P_n$$

y en donde cada periodo  $P_i$  tendrá la siguiente estructura:

$P_i$  = Capital  
 Actual Position  
 Open Price  
 Close Price  
 Min Price  
 Max Price  
 Volume  
 ATR  
 MA8  
 EMA20  
 EMA50  
 EMA200  
 RSI  
 DMI  
 MACD  
 BollingerBands

### 3.3 Acciones

En cada instante  $t$  el agente podrá decidir comprar, vender o esperar, cada vez que decida comprar o vender, la cantidad de acciones operadas, se calculara en base a una estrategia de entrada y salida del activo pre configurada en el agente. El agente no tendrá limitación alguna en cuanto a la cantidad de acciones a comprar o vender. Si bien, en la realidad esto no necesariamente seria posible dado que para que pueda comprar una cantidad  $n$  a un precio  $p$ , debe existir algún vendedor dispuesto a vender  $n$  acciones a un precio  $p$  cada una. Ademas cada operación de venta o compra, tendrá un costo asociado  $c$ , el cual sera descontado de su capital. El agente no podrá ejecutar la acción elegida si su capital no alcanza para cubrir el costo total de la operación.

### 3.4 El Agente

Recordemos que el objetivo principal de agente es seleccionar acciones que maximicen su recompensa futura, es decir, deberá de maximizar su capital inicial, para esto vamos a asumir que las recompensa futuras están descontadas por un factor  $\gamma$ , por cada instante de tiempo  $t$ , es decir, vamos a valorizar los rewards más cercanos en el tiempo por un factor  $\gamma$ . Definimos así el, future discounted reward:

$$R_t = \sum_{t'=t}^T \gamma^{t'-t} r_{t'} \quad (3.1)$$

Donde  $T$  es el instante en el agente termina de invertir (o bien por qué perdió todo su capital inicial o bien porque se acabaron los datos de simulación). También vamos a definir una función estado-acción óptima,  $Q^*(s_t, a_t)$ , como la máxima recompensa esperada que podemos alcanzar estando en el estado  $s_t$ , y seleccionando la acción  $a_t$ , y luego continuando con una estrategia  $\pi$

$$Q^*(s, a) = \max_{\pi} \mathbb{E}[R_t | s_t = s, a_t = a, \pi] \quad (3.2)$$

Esta función tiene una propiedad importante, conocida como Bellman equation, que intuitivamente se basa en lo siguiente: Si conociéramos el valor óptimo de  $Q^*(s', a')$  para el próximo estado  $s'$  y para cada posible acción posible  $a'$ , entonces la estrategia óptima de selección de una acción para el estado actual  $s$  que maximice la recompensa esperada está dada por

$$Q^*(s, a) = \max_{\pi} \mathbb{E}_{s' \approx \epsilon} [r + \gamma \max_{a'} Q^*(s', a') | s, a] \quad (3.3)$$

La idea general de nuestro algoritmo de aprendizaje va a ser estimar esta función  $Q^*$  usando la ecuación de bellman en forma iterativa, esto es:

$$Q_{i+1}^*(s, a) = \max_{\pi} \mathbb{E}_{s' \approx \epsilon} [r + \gamma \max_{a'} Q_i^*(s', a') | s, a] \quad (3.4)$$

donde  $i \rightarrow \infty$  y  $Q_i \rightarrow Q^*$

Esta es la idea detrás del algoritmo de Q-learning que implementará el agente, la cual en principio, parecía suficiente para lograr que el agente aprenda a invertir, sin embargo en la práctica, veremos que en general los estados del mercado no van a repetirse en forma idéntica y que además la cantidad de estados posibles hara que sea prácticamente inviable lograr la convergencia hacia  $Q^*$ .

### 3.5 Q-Network

Para poder solucionar este inconveniente, vamos a utilizar una función de aproximación, cuyo objetivo será obtener una generalización de los estados, lo cual permitirá en teoría, lograr más rápidamente la convergencia a  $Q^*$ .

Para implementar esta función utilizaremos un red neuronal con con una matriz de pesos  $\theta$ , inicializada aleatoriamente:

$$Q(s, a, \theta) \approx Q^*(s, a) \quad (3.5)$$

a la cual denominaremos como Q-network. Dicha Q-network será entrenada optimizando una función de pérdida L y usando el algoritmo del gradiente descendente.

$$L_i(\theta) = \mathbb{E}_{s, a \sim p(\cdot)} [(y_i - Q(s, a, \theta_i))^2] \quad (3.6)$$

$$\text{donde } y_i = \mathbb{E}_{s \sim \epsilon} [r + \gamma \max_{a'} Q(s', a', \theta_{i-1}) | s, a] \quad (3.7)$$

Además se usará experience replay para entrenar la Q-Network, es decir, los datos de ejemplos serán generados a partir de la propia experiencia que vaya adquiriendo el agente. Para lograr esto, en cada instante de tiempo t, el agente deberá guardar en un replay memory D, cada experiencia  $\epsilon = (s_t, a_t, r_t, s_{t+1})$ , observada. Transcurridos n días el agente deberá ejecutar la fase de entrenamiento de su Q-Network, para esto se generará un mini batch de experiencias pasadas, obtenidas de su replay memory y con ella se procederá a realizar una última fase actualización de la función de estimación de  $Q(s, a, \theta_i)$ , usando la regla de actualización de Q-learning y la nueva matriz de pesos  $\theta_{i+1}$ .

### 3.6 Algoritmo de aprendizaje

A continuación detallamos el pseudo algoritmo de aprendizaje de nuestro agente.

---

$D \leftarrow \text{new\_memory\_replay}(n)$     $\theta \leftarrow \text{random}()$     $Q^*(s, a) \leftarrow Q(\theta)$     $\text{episode} \leftarrow \text{generate\_episode}() \neq \text{null}$

---





## 4 Arquitectura y Diseño



## 5 Evaluación y Desempeño



## **6 Conclusiones y Recomendaciones**



## Bibliografía