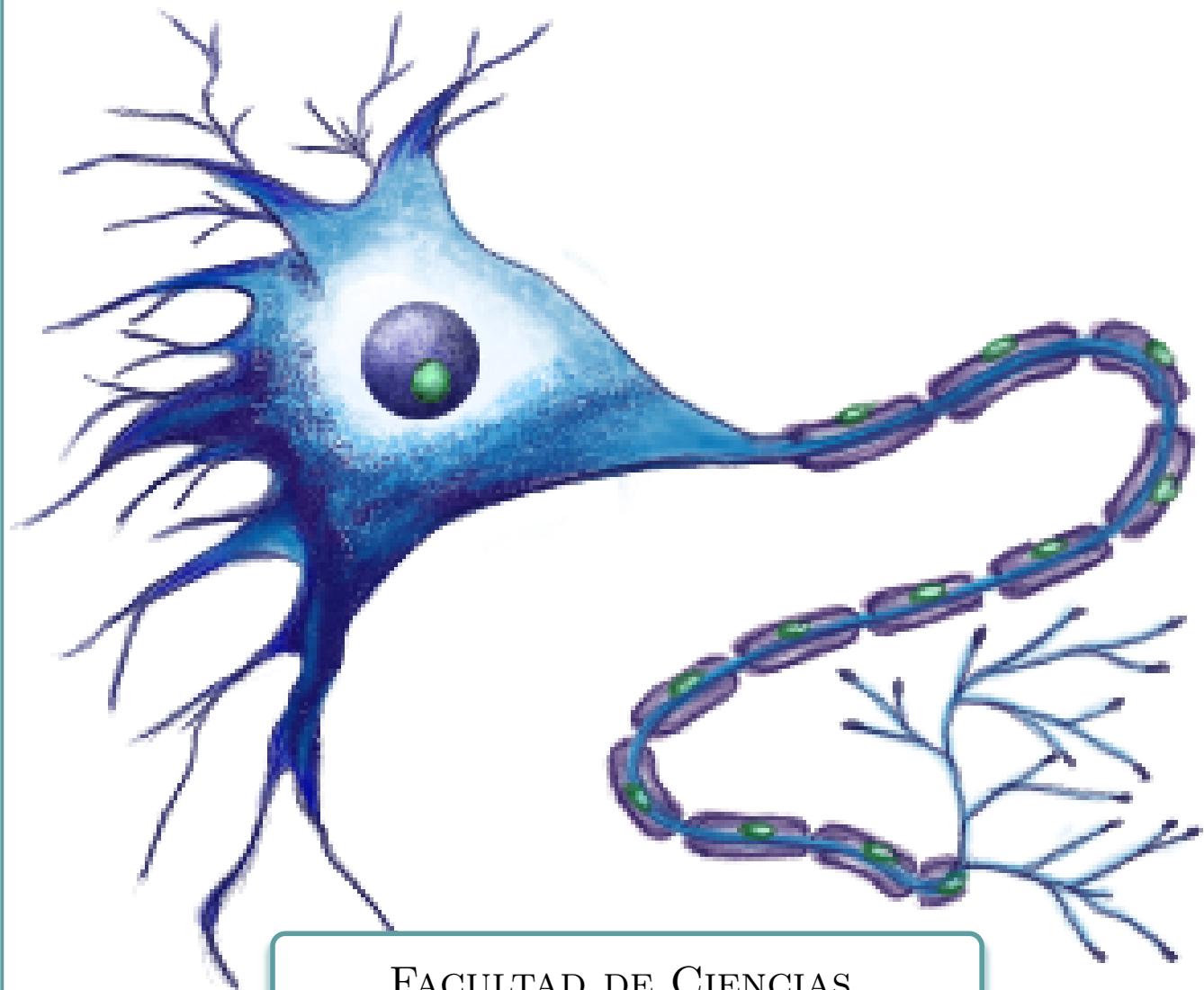


# Redes Neuronales

*Notas de clase*

Karla Fernanda Jiménez Gutiérrez  
Verónica Esther Arriola Ríos



FACULTAD DE CIENCIAS,  
UNAM



# Índice general

<b>Índice general</b>	I
<b>I Antecedentes</b>	2
<b>1 Neurona biológica</b>	3
1.1 Neurociencias computacionales . . . . .	3
1.2 Sistema Nervioso . . . . .	5
1.2.1 Cerebro . . . . .	8
1.2.2 Zonas funcionales . . . . .	10
1.3 Neurona biológica . . . . .	11
1.3.1 La neurona . . . . .	11
1.3.2 Elementos de las neuronas y tipos . . . . .	13
1.3.3 Sinapsis . . . . .	15
1.3.4 Campos receptivos . . . . .	19
1.3.5 Señal eléctrica . . . . .	21
<b>2 Modelo de Hodgkin-Huxley</b>	24
2.1 Introducción . . . . .	24
2.2 Membrana y canal . . . . .	26
2.3 Potenciales de Nerst o de reposo . . . . .	29
2.4 Modelo de la membrana como bicapa de lípidos . . . . .	29
2.4.1 Las conductancias iónicas . . . . .	33
2.5 Modelo de las compuertas iónicas controladas por voltaje . . . . .	35
2.6 Dinámica del voltaje durante un disparo . . . . .	38
2.7 Simulación usando el método de Euler . . . . .	42
2.8 Información condificada en las dendritas . . . . .	44
<b>3 Aprendizaje de máquina</b>	48
3.1 Introducción . . . . .	48
3.2 Espacio de Hipótesis . . . . .	50
3.3 Clasificación de los conjuntos de datos . . . . .	51
3.3.1 Tipos de aprendizaje . . . . .	52

---

## ÍNDICE GENERAL

<b>II Redes dirigidas acíclicas</b>	<b>54</b>
<b>4 Perceptrón simple</b>	<b>55</b>
4.1 Perceptrón . . . . .	55
4.2 Compuertas lógicas con neuronas . . . . .	58
4.3 Funciones de activación . . . . .	60
4.4 Funciones de error . . . . .	60
4.5 Medidas de rendimiento . . . . .	63
<b>5 Perceptrón multicapa</b>	<b>69</b>
5.1 Intro . . . . .	69
5.2 XOR . . . . .	69
5.3 Propagación hacia adelante manual . . . . .	71
5.4 Propagación hacia adelante vectorizada (con matrices) . . . . .	72
5.5 Interpretación matemática del mapeo no lineal . . . . .	75
5.6 Propagación hacia adelante para el perceptrón multicapa . . . . .	76
<b>6 Entrenamiento por retropropagación</b>	<b>78</b>
6.1 Esquema general entrenamiento . . . . .	78
6.2 Función de error: Entropía cruzada . . . . .	81
6.3 Derivada de la función logística . . . . .	83
6.4 Parcial con respecto a los pesos en la penúltima capa . . . . .	85
6.5 Parcial con respecto a los pesos en la última capa . . . . .	88
6.6 Vectorización . . . . .	91
<b>7 Optimización del entrenamiento</b>	<b>96</b>
7.1 Problemas en redes profundas . . . . .	96
7.2 Gradiente desvaneciente (o que explota) . . . . .	96
7.3 Entrenamiento en línea vs en lotes . . . . .	96
7.4 Normalización y normalización por lotes . . . . .	96
7.5 Regularización . . . . .	96
<b>8 Caso de análisis e interpretación</b>	<b>97</b>
8.1 Red Hinton árbol familiar con numpy (entrenamiento) . . . . .	97
8.2 Red Hinton árbol familiar con pytorch . . . . .	97
<b>9 Entrenamiento con genéticos</b>	<b>98</b>
9.1 Algoritmos genéticos . . . . .	98
9.2 Neuroevolución . . . . .	98
9.2.1 Antecedentes: Aprendizaje por refuerzo en videojuegos . . . . .	98
9.2.2 Arquitectura para estimar la función de recompensa . . . . .	98
9.2.3 Entrenamiento . . . . .	98
<b>10 Mapeos autoorganizados</b>	<b>99</b>
10.1 Introducción . . . . .	99

---

ÍNDICE GENERAL

10.2 Aprendizaje no supervisado . . . . .	99
10.3 Mapeos autoo-organizados . . . . .	99
10.4 Kohonen . . . . .	99
<b>11 Redes Neuronales Convolucionales</b>	<b>100</b>
11.1 Convolución . . . . .	100
11.2 Redes Convolucionales . . . . .	100
11.3 Softmax . . . . .	100
11.4 MNIST . . . . .	100
<b>III Redes con ciclos</b>	<b>101</b>
<b>12 Redes Neuronales Recurrentes</b>	<b>102</b>
12.1 Derivadas ordenadas . . . . .	102
12.2 Retropropagación en el tiempo . . . . .	102
12.3 Sistemas dinámicos y despliegue del grafo . . . . .	102
12.4 Arquitectura recurrente universal . . . . .	102
12.5 Función de error . . . . .	102
12.6 Forzamiento del profesor . . . . .	102
<b>13 Atención</b>	<b>103</b>
<b>14 LSTM</b>	<b>104</b>
<b>15 GRU</b>	<b>105</b>
<b>16 Casos de análisis: etiquetado de palabras y conjugación de verbos</b>	<b>106</b>
<b>IV Redes no dirigidas</b>	<b>107</b>
<b>17 Redes de hopfield</b>	<b>108</b>
17.1 Entrenamiento . . . . .	108
<b>18 Máquinas de Boltzman</b>	<b>109</b>
18.1 Entrenamiento . . . . .	109
18.1.1 Partículas y partículas de fantasía . . . . .	109
18.1.2 Máquinas de Boltzman Restringidas . . . . .	109
<b>19 Redes adversarias</b>	<b>110</b>
19.1 GANs . . . . .	110
<b>A Ecuaciones diferenciales</b>	<b>111</b>



# Etc

A lo largo del texto se utilizará la siguiente notación para diversos elementos:

Conjuntos	C
Vectores	x
Matrices	M
Unidades	cm

# **Parte I**

## **Antecedentes**

# 1 | Neurona biológica

## Neurociencias computacionales

Las redes neuronales surgieron completamente inspiradas en los sistemas biológicos. Lo que estamos haciendo los computólogos es tomar una idea de la naturaleza, una idea que ha probado ser sumamente efectiva para procesar información y que logra resolver problemas que nosotros aún no sabemos solucionar con modelos diseñados explícitamente. Los más notorios son:

- Problemas de visión por computadora.
- Procesamiento del lenguaje natural.

A lo largo del texto obtendremos una somera idea de qué hace el sistema nervioso de un ser humano, tomaremos también ejemplos de animales como el calamar gigante y cangrejos; ejemplos que han permitido estudiar biológicamente cómo funcionan las neuronas y cómo funciona su sistema nervioso.

Entonces por un momento pensemos en el sistema nervioso como un todo, lo que realmente está pasando al computar no es el cálculo del proceso de una sola neurona sino de la colección de todas ellas. Lo que sucede con los sistemas biológicos es que son muchísimo más complicados que lo que vamos a ver nosotros como modelos computacionales, sin embargo muchísimas empresas están utilizando estas técnicas. El sistema nervioso como un todo es bastante más complejo, pero conforme han ido evolucionando las redes neuronales computacionales, ya con sus arquitecturas y organizaciones, se están volviendo también más complejas. Varias de las estructuras más exitosas tienen un análogo muy fuerte con un sistema nervioso natural.

Veamos un campo conocido como **neurociencias computacionales** el cual se dedica explícitamente al estudio/modelado de los sistemas biológicos pero ya conjuntando varios campos. Se interesan notablemente en: descripciones y modelos funcionales biológicamente realistas de neuronas y sistemas neuronales. En contraposición, los modelos que veremos en redes neuronales computacionales no necesariamente tienen que ser realistas, lo que nos interesa es que resuelvan los problemas, si se desvían un poco de cómo funcionan los sistemas naturales en un principio no es problema.

## 1. Neurona biológica

Ahora, ¿qué le interesa modelar a las neurociencias computacionales? Se fijan en la fisiología y en la dinámica de estos sistemas, combinando varias ciencias tales como:

- **Biofísica** por el estudio de las propiedades físicas detrás de los sistemas biológicos.
- **Neurociencias tradicionales** con modelos matemáticos.
- **Ciencias de la computación** tanto en la parte del modelado como en la parte de la implementación de estos modelos y la generación de simulaciones computacionales.
- **Ingeniería eléctrica** donde se está diseñando hardware especializado para ejecutar modelos de manera eficiente, algunos de los modelos matemáticos están basados en circuitos eléctricos.
- **Ciencias cognitivas** que tratan de ver qué se está codificando dentro de un sistema nervioso y cómo podemos interpretar esa información que está ahí guardada.

Vamos a ver cómo están influyendo todos estos antecedentes en lo que van a hacer las ciencias de la computación pero con su propio modelo de redes neuronales, pues existe una conexión muy fuerte entre estos dos campos.

Las neurociencias computacionales, como se mencionó anteriormente, estudian modelos del sistema nervioso y clasifican estos modelos en tres tipos:

1. **Modelos descriptivos**, nos limitamos a decir qué está haciendo un sistema; en particular aquí son muy famosos los experimentos con ratones, se está tratando de ver qué puede hacer, que no puede hacer, que puede aprender, que no, pero no se puede explicar “¿cómo?”, simplemente se dice qué es lo que está sucediendo.
2. **Modelos mecanistas**, donde ahora sí nos interesa saber, ¿cómo es que están haciendo las cosas? Aquí vamos a ver cómo los modelos matemáticos precisamente nos están tratando de describir cómo puede ser que se están conectando estas neuronas, cómo pueden estar funcionando las redes de neuronas, cómo podría estarse almacenando la información y transfiriendo de un lado a otro.
3. **Modelos interpretativos**, nos dan una idea del por qué o para qué lo hacen. Se tiene que buscar intencionalidad, razonamiento de más alto nivel.

Cuando trabajemos con redes neuronales computacionales vamos a notar que sí necesitamos adentrarnos un poco en los tipos 2 y 3. Para romper esa traba con nuestras redes neuronales, donde sabemos que aprendieron, pero no estamos ni siquiera seguros de qué aprendieron o porqué lo aprendieron así, vamos a tener que utilizar herramientas matemáticas para tratar de descubrir qué es lo que realmente está haciendo la red entrenada.

Ahora revisemos los **objetivos del modelado** en neurociencias:

(Empezando desde lo más granular que es cada una de las neuronas)

- Las **corrientes** que están pasando a través de las membranas de las neuronas, la influencia que tienen en el paso de la información.
- Las **proteínas** van a jugar un papel importante en la conducción de elementos iónicos no transmisores (acoplamientos químicos).

(El siguiente nivel con combinaciones de varias neuronas)

- Las **oscilaciones de las redes** completas, qué pasa con estas señales, pulsos eléctricos que se están transfiriendo de unas regiones a otras y que empiezan a producir oscilaciones con ciertos períodos, regiones de actividad o regiones que se apagan.
- **Arquitectura topográfica y de columnas** cómo están organizadas estas neuronas, quiénes están conectadas con quiénes, cómo reaccionan dentro de ciertas regiones identificadas, cómo interactúan con otras regiones. Se puede identificar una arquitectura tanto desde el punto de vista fisiológico como del punto de vista funcional. Un caso particular de estas estructuras es la formación de columnas de neuronas que están altamente conectadas y trabajan como una unidad.
- El **aprendizaje**, es decir estamos procesando información, guardando información, recuperándola y eso permite que los seres que cuentan con un sistema nervioso tengan características especiales cuyo comportamiento se puede modificar conforme aprenden.
- La **memoria**, necesitamos almacenar información y recuperarla para procesarla.

## Sistema Nervioso

**¿Qué es un nervio?** Un nervio es una gran colección de axones que están viajando todos juntos en una especie de cable (fibra), pasan vasos sanguíneos por en medio de los nervios. Esto es de lo que está formando el sistema nervioso, se originan desde la médula espinal (31 pares de nervios raquídeos) o encéfalo (12 pares de nervios craneales).

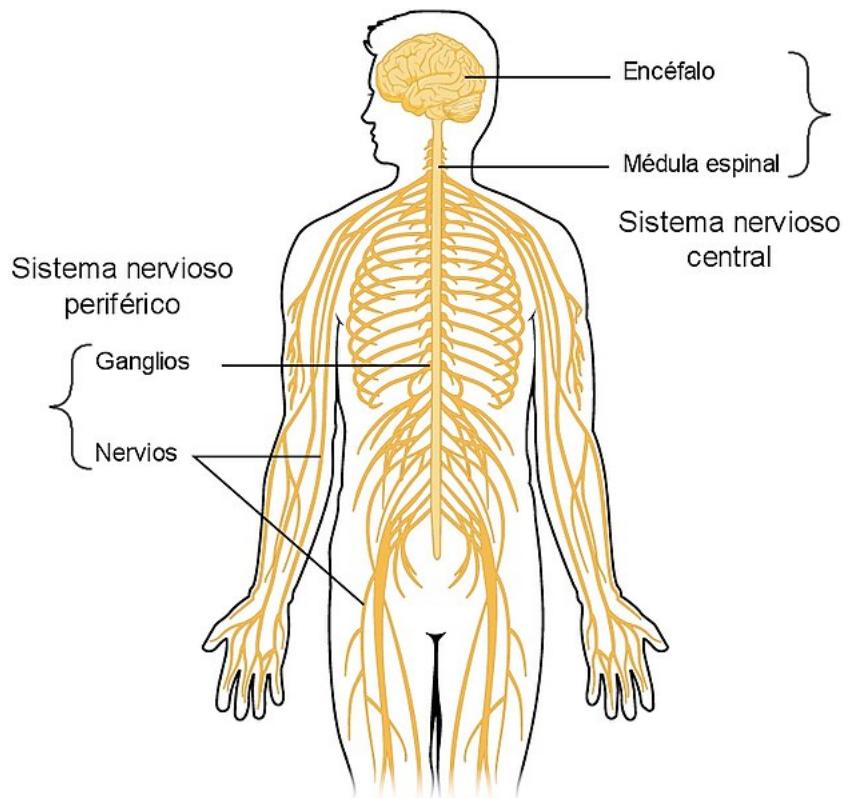
Los nervios son estructuras conductoras de impulsos nerviosos situados fuera del sistema nervioso central, es decir, estamos hablando de todos estos axones que salen desde del cráneo, la médula espinal y están descubriendo el resto del cuerpo. Están formados por un conjunto de axones agrupados cada uno de los cuales procede de una neurona. Pueden ser clasificados como:

- **Motores** salidas, ejecución/acción
- **Sensitivos** entradas

## 1. Neurona biológica

- **Mixtos** son mayoría, tienen tanto fibras sensitivas como motoras

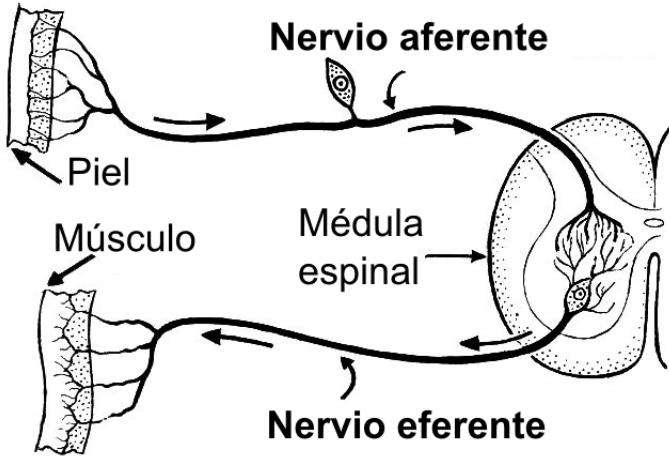
Tenemos dos grandes partes del sistema nervioso, el **sistema nervioso periférico** y el **sistema nervioso central**, como se puede ver en la imagen 1.1.



**Figura 1.1** Overview of Nervous System esp, OpenStax, 20 diciembre 2018, WIKIMEDIA COMMONS, [https://upload.wikimedia.org/wikipedia/commons/0/07/1201\\_Overview\\_of\\_Nervous\\_System\\_esp.jpg](https://upload.wikimedia.org/wikipedia/commons/0/07/1201_Overview_of_Nervous_System_esp.jpg), CC BY-SA 4.0

En del sistema nervioso periférico tenemos al:

- **Sistema somático** se controla de forma voluntaria, se conforma de nervios conectados a músculos voluntarios esqueléticos y receptores sensoriales, de los cuales unos son:
  - \* de entrada, **aferentes**
  - \* de salida, **eferentes**
- **Sistema autónomo** funciona de forma involuntaria, se conforma de nervios que se conectan con el corazón, los vasos sanguíneos, los pulmones, el estómago, los intestinos, glándulas



**Figura 1.2** Diagrama explicativo del recorrido eferente y el aferente, Pearson Scott Foresman, 26 August 2010, WIKIMEDIA COMMONS, [https://upload.wikimedia.org/wikipedia/commons/3/3e/Afferent\\_%28PSF%29.es.png](https://upload.wikimedia.org/wikipedia/commons/3/3e/Afferent_%28PSF%29.es.png), CC0

Ahora respecto al sistema nervioso central lo integra:

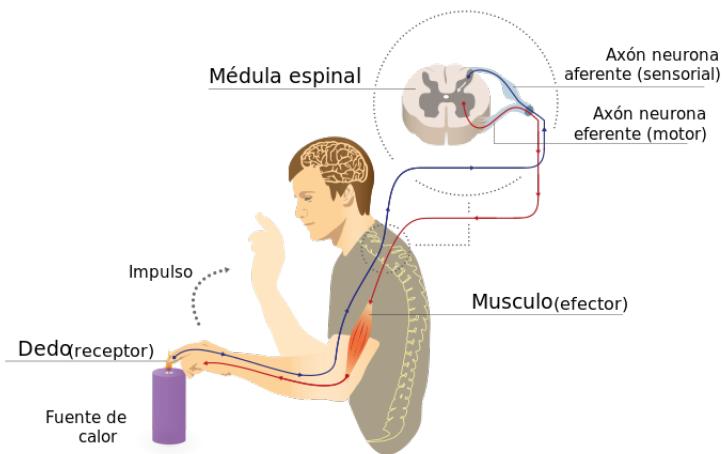
- La médula espinal
  - \* Dentro de esta hay una organización, la presencia de **ciclos de retroalimentación local**, es decir, nuestro sistema va a estar en diferentes etapas son nervios que no necesitan pasar por todo el procesamiento cerebral, las señales simplemente entran llegan a una fase local e inmediatamente reaccionan, ver el ejemplo de la imagen 1.3. Ocurren en un ciclo local y esto también puede convertirse en algo muy importante a la hora de hacer cálculos, no siempre es necesario pasar todo por todas las capas de procesamiento.
  - \* **Señales de control motor descendientes** del cerebro hacia las neuronas motoras, estas son señales que provienen de un campo en una capa mucho más alta de procesamiento y provocan movimientos.
  - \* **Axones sensoriales ascendentes** donde el cuerpo de la neurona está afuera y la información va a viajar hacia arriba, desde los músculos, piel y estas señales viajan hasta el cerebro.
- El encéfalo

Cada colección de nervios que sale de la base del cerebro se asocian con funciones muy específicas (en su mayoría).

Notas:

- Este sistema está hecho en diferentes niveles locales, entradas y salidas

## 1. Neurona biológica



**Figura 1.3** Esquema explicativo del arco reflejo, Marta Aguayo, 18 diciembre 2014, WIKIMEDIA COMMONS, [https://upload.wikimedia.org/wikipedia/commons/c/cb/Imgnotraçat\\_arco\\_reflexo\\_esp.svg](https://upload.wikimedia.org/wikipedia/commons/c/cb/Imgnotraçat_arco_reflexo_esp.svg), CC BY-SA 3.0.

- El procesamiento que esté ocurriendo en el encéfalo puede tener diferentes capas y eso se verá reflejado cuando nosotros definamos arquitecturas para las redes neuronales.
- Las redes neuronales actuales, que han tenido más éxito, se componen de diferentes subunidades o diferentes redes que hacen cosas locales. Es decir esta estructura global que estamos viendo, se está empezando a reproducir/imitar ya con las neuronas computacionales.

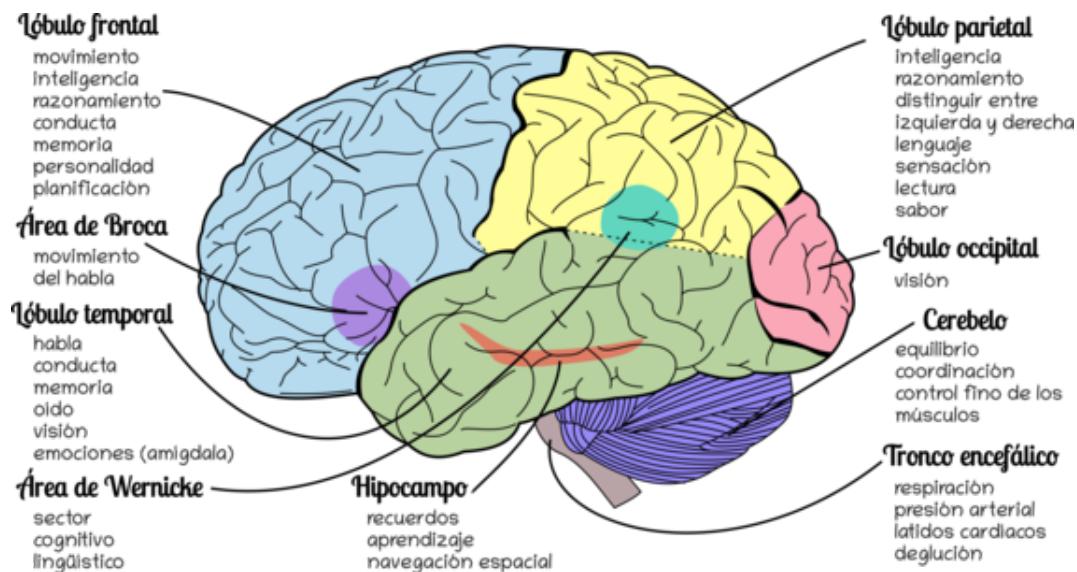
## Cerebro

En esta parte vamos a preocuparnos sobre todo por la parte funcional. Haciendo una breve analogía, vamos a hacer una visión general del “hardware”, para ver qué efectos va a tener en el “software”. En general la arquitectura de cada cerebro es completamente diferente al cerebro de otras personas. Se ha intentado averiguar qué está haciendo cada región con diferentes estudios por ejemplo, ver cuánta sangre se está bombeando en diferentes regiones del cerebro dependiendo de los estímulos que se le presentan a una persona, o si alguna persona tiene un padecimiento se tratan de tomar escaneos para ver qué regiones del cerebro están funcionando y cuáles presentan lesiones. A partir de las lesiones, lo que hacen es que una vez que está identificada la actividad que ya no se puede realizar de forma normal, averiguan qué región era responsable de esa actividad, que ahora está dañada.

Gracias a esos estudios, se ha logrado identificar más o menos en manera general, a qué se dedica cada una de las regiones del cerebro. En ocasiones no se puede decir exactamente qué tan vinculadas están (las regiones) o por qué se están activando otras regiones.

Hay partes funcionales que se comparten entre las diferentes regiones y no están ubicadas en un solo lugar. Otra parte importante a mencionar es, el cerebelo que se considera prácticamente vital, cumple con funciones tales como el equilibrio, la coordinación, el control fino de los músculos, de hecho tiene más neuronas que el cerebro y aun así hay niños que nacen y viven sin cerebelo.

A continuación se mencionan algunas de las diferentes funciones de las regiones, que se han identificado en la imagen 1.4:



**Figura 1.4** Diagrama básico de las regiones del cerebro.

**Lóbulo frontal** se le puede asociar con la parte del raciocinio, la parte de inteligencia, la conducta, la memoria, la personalidad, la capacidad para realizar planes complejos a largo plazo y también es responsable de algunas actividades de movimiento. Dentro de este destaca el área de broca, su principal función es el movimiento del habla, mover los labios, la boca.

**Lóbulo temporal** aquí está otra parte del habla, que tiene que ver más con el uso de símbolos para el lenguaje, la conducta, memoria, aquí se procesa el oído, un poco de visión y emociones. Dentro de este está (compartida entre el lóbulo parietal) el área de Wernicke, trabaja con la parte lingüística, y de cognición. También dentro de este está el **hipocampo** trabaja con recuerdos, aprendizaje y navegación espacial, cómo sabemos cómo llegar de un lado hacia otro.

**Lóbulo parietal** trabaja con la inteligencia, razonamiento, distinguir entre izquierda y derecha, lenguaje, sensación, lectura y sabor.

## 1. Neurona biológica

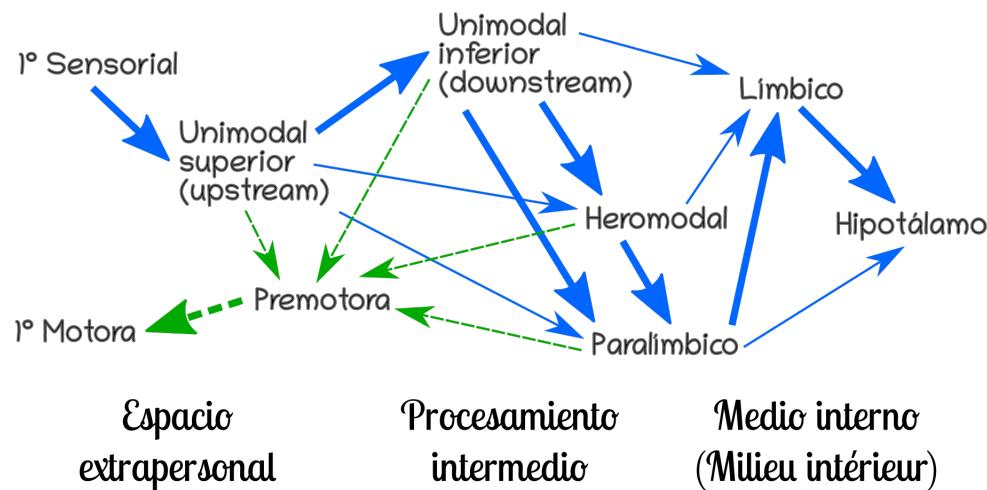
**Lóbulo occipital** se dedica prácticamente solamente a visión, es una región un tanto amplia. En particular en el área de robótica cuando están programando un robot o móvil, los robots tienen dos laptops y una de ellas se dedica prácticamente solo a procesar la visión.

**Cerebelo** se encarga del equilibrio, la coordinación fina de los músculos.

**Tronco encefálico** se encarga de la respiración, presión arterial, latidos cardíacos, de ilusión, conciencia.

## Zonas funcionales

Para visualizar mejor la parte de la arquitectura, que tiene el cerebro para realizar todo lo que se le conoce como, la ruta desde la sensación hasta la cognición, veremos un diagrama de la parte funcional del cerebro.



**Figura 1.5** Diagrama de la arquitectura del cerebro a nivel funcional.

Explicando el diagrama 1.5, en la primera parte (espacio extrapersonal) vamos a pensar en la entrada sensorial, que se enfoca muchísimo en la parte de visión y audio (en general todos los sentidos), notamos qué de las neuronas que están en la parte sensorial, su primera conexión es hacia una capa que se le llama unimodal superior, aquí se procesa la información de cada sentido de manera individual, es decir, las neuronas o solamente están procesando visión o solamente audio, todavía no se mezclan, por ejemplo de visión, se separan colores e intensidad lumínica, se empieza a detectar algunas esquinas, alguna inclinación, la dirección de las luces y las sombras. Notemos que desde aquí hay una rápida conexión a la sección premotora y luego hacia la parte motora, recordando la mención de los circuitos locales y de reflejos, aquí prácticamente lo podemos ver (en este pequeño camino).

Pasando de este primer procesamiento básico entramos al siguiente que es el unimodal inferior, (aquí aún se está trabajando con procesamiento de una sola modalidad) visión sigue siendo visión, audio sigue siendo audio, pero ya son procesamientos un poco más complejos, por ejemplo, reconocimiento de rostros, de objetos. En esta parte tenemos un rápido ciclo de regreso a la parte premotora, por ejemplo la acción de ver a mi mamá y saludarla (aquí aún no se tiene que razonar demasiado).

En la siguiente fase (medio interno), se conecta hacia tres áreas, la **heromodal**, ya se integran diferentes modalidades (audio y visión) ejemplo, oigo que me hablan y volteo a ver, aquí se está juntando ambas cosas, el **límbico** y el **paralímbico** que trabajan con la parte de las emociones y conceptos abstractos.

Finalmente, llegamos al **hipotálamo** que es donde están todas las emociones, en las conexiones entre estas regiones, estarían los procesamientos de alto nivel.

Ahora estas diferentes regiones se replican de cierta manera cuando estamos haciendo los diseños de las arquitecturas modernas para redes neuronales. En algunas ocasiones se comienza con algunas capas de neuronas, haciendo procesamientos con una sola modalidad, extrayendo datos básicos, después se van componiendo en figuras más complejas y después hasta podemos combinar bloques de neuronas, para poder resolver problemas que tomen en cuenta diferentes modelos.

## Neurona biológica

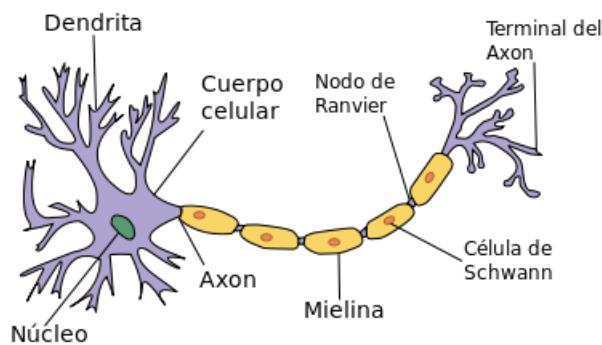
### La neurona

La neurona es un tipo de célula perteneciente al sistema nervioso central, que se comunica tanto por señales eléctricas como por señales químicas. Cada neurona tiene:

- Un cuerpo celular (**soma**) que contiene un núcleo y otros componentes celulares.
- Una zona de recepción denominada **dendritas**.
- Una zona de emisión conocida como **axón**, compuesto de:
  - ★ Cono axonico.
  - ★ Membrana plasmática axonica y citoplasma.
  - ★ Recubrimientos de mielina, interrumpidos con intervalos regulares de nódulos (anillos) de Ranzier.
  - ★ Terminales del axón donde se encuentran los **botones sinápticos**.

Pensemos en la neurona como toda una compuerta, por un lado, está el cuerpo de una neurona típica, en las dendritas tenemos una mezcla de neurotransmisores y iones que

## 1. Neurona biológica



**Figura 1.6** Neurona, Acracia, 14 January 2007, Wikimedia Commons, <https://commons.wikimedia.org/wiki/File:Neurona.svg>, Creative Commons Attribution-ShareAlike 2.5 Generic

pueden moverse a través de la membrana. La forma en que intercambia información es mediante sustancias químicas y iones que se están intercambiando, entre la parte de afuera y de adentro de la neurona. Particularmente en las dendritas, se tienen terminaciones que se pueden conectar con otras neuronas y de esta manera permitir el paso de información.

- Neurona presináptica, transmite una señal.
- Neurona posináptica, recibe una señal.

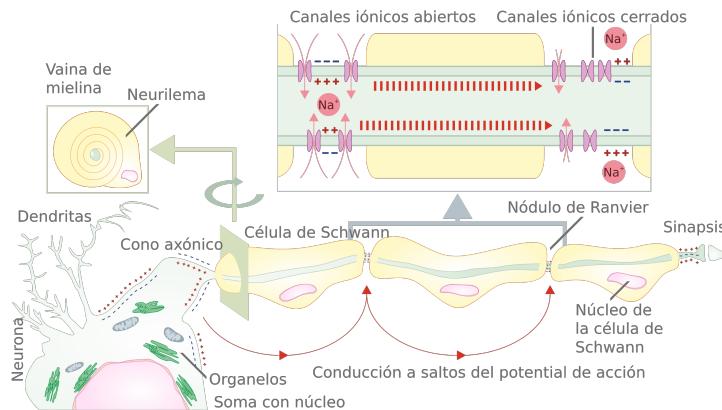
En el interior de la neurona hay una cierta carga eléctrica, en el exterior (el líquido de afuera) hay otra carga eléctrica, es decir, hay una **diferencia de potencial** entre el interior y el exterior de la neurona, por eso se dice que la membrana axónica en sí misma tiene una carga eléctrica. Dado que es porosa, esta membrana va a estar intercambiando partículas con el exterior, esto va a hacer que la polarización de esta membrana vaya cambiando, si en algún momento la diferencia de potencial neta rebasa un cierto umbral.

### *Transmisión de señales y almacenamiento de información:*

1. La neurona desde sus dendritas recibe señales de otras neuronas vecinas.
2. Cada señal se va acumulando en su cuerpo hasta el cono axónico, donde se van a estar sumando la contribución de todos los efectos de cambios de potencial.
3. En el momento que se rebase un cierto valor umbral, la diferencia de potencial se propaga hasta los botones terminales.
4. La neurona entra en un período refractario, donde empieza a cambiar el potencial entre el cono axónico y el axón de la neurona.
5. Se va a transmitir un disparo eléctrico en seguida,

6. La neurona se va a quedar totalmente quieta, durante un breve momento para que la señal pueda viajar hacia el axón.
7. Se va a notar un cambio muy violento en el voltaje, que se va recorriendo a lo largo de todo el axón.

La neurona típica tiene unas células de mielina, que forman nodos que van cubriendo al axón para evitar que se pierda la señal, estos nodos recargan otra vez la señal y permite que avance, al siguiente nodo, donde se recarga nuevamente y avance, hasta que logre llegar al final de axón.



**Figura 1.7** Corriente iónica por el axón (efecto de corto plazo), Helixitta, 1 octubre 2015, Wikimedia Commons, [https://commons.wikimedia.org/wiki/File:Propagation\\_of\\_action\\_potential\\_along\\_myelinated\\_nerve\\_fiber\\_en.svg](https://commons.wikimedia.org/wiki/File:Propagation_of_action_potential_along_myelinated_nerve_fiber_en.svg), Creative Commons Attribution-ShareAlike 4.0 International

Este trayecto puede ser de una neurona a unas pocas neuronas vecinas, hasta unos cuantos metros (ej. esta podría estar en la médula espinal y el axón llegar hasta el dedo). Cuando la señal llega a la terminal de el axón, hay varias terminales que van a reaccionar ante el cambio de electricidad, mediante la liberación de unas vesículas, que contienen **neurotransmisores**.

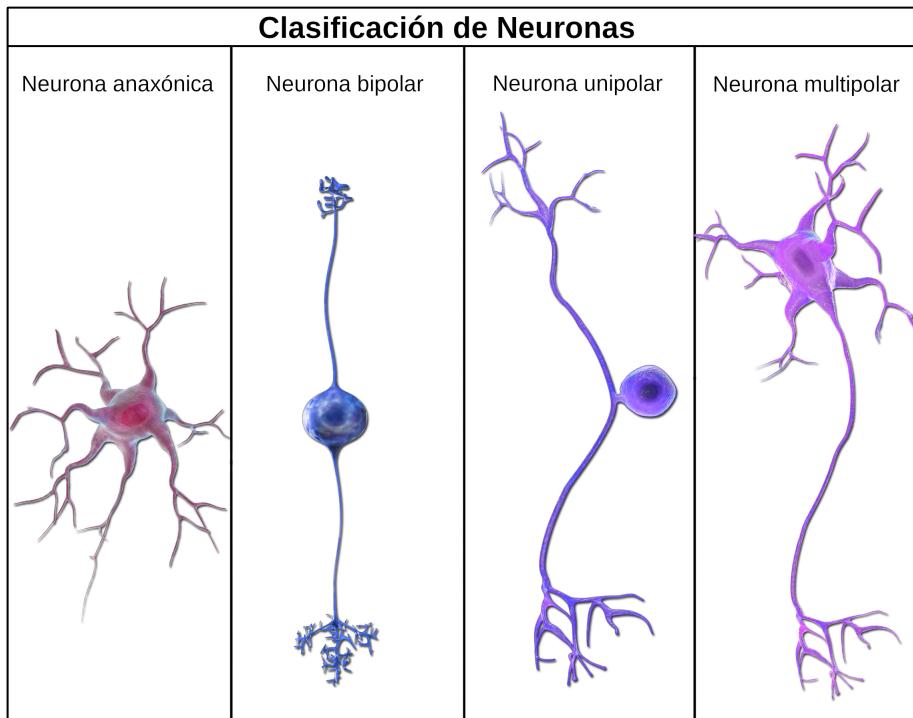
## Elementos de las neuronas y tipos

Elementos durante la transmisión de señales:

- **Neurotransmisores:** son los mensajeros químicos que se comunican entre neuronas adyacentes; La liberación de neurotransmisores de una neurona ayudará a despolarizar o hiperpolarizar (aumentar la magnitud de la carga) la neurona adyacente, lo que hará que sea más o menos probable que ocurra un potencial de acción en la siguiente neurona.

## 1. Neurona biológica

- **Impulsos eléctricos:** potenciales de acción que son, cambios de voltaje que van a ir ocurriendo a lo largo del axón. Sucede una vez que se acumularon demasiadas señales a través de las dendritas, entonces la neurona puede disparar un impulso eléctrico, a través del axón, que va a provocar que su terminal libere más químicos, estos químicos son los que hacen los efectos pequeños en cada uno de los cuerpos de las neuronas posinápticas.
- **Plasticidad:** modificación a largo plazo de las conexiones entre neuronas. En el cerebro las neuronas pueden cambiar de manera permanente, perder canales (que permiten el intercambio de nuevos transmisores sin impulsos eléctricos), formar más canales o incluso pueden crear protuberancias. Por ejemplo, cuando un cerebro aprende está transformando su arquitectura, es decir, los aprendizajes de largo plazo, modifican el cerebro y en consecuencia va a pensar y reaccionar distinto, que antes del aprendizaje.

*Clasificación de tipos de neuronas:*

**Figura 1.8** Representación de la clasificación de neuronas, BruceBlaus, 26 junio 2017, Wikimedia Commons, [https://commons.wikimedia.org/wiki/File:Neuron\\_Classification.png](https://commons.wikimedia.org/wiki/File:Neuron_Classification.png), Creative Commons Attribution-ShareAlike 4.0 International

- **Neuronas sin axones**, nunca dispara, pero si tiene intercambios de neurotransmisores en las dendritas

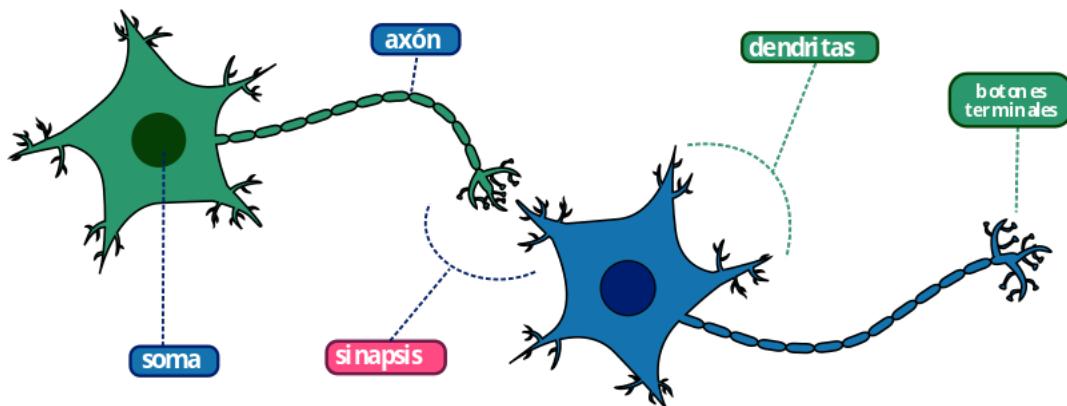
- **Neuronas bipolares**, tienen dos axones.
- **Neurona unipolar**, solamente hay una conexión entre el cuerpo y el axón, pero el axón tiene dos ramas, cuando dispara va a disparar hacia los dos lados, haciendo llegar su señal a diferentes regiones.
- **Neurona multipolar**, la más conocida, empieza con un cuerpo con dendritas y luego un largo axón que va a terminar con varias terminaciones axónicas.

Cuando modelamos redes neuronales lo típico es modelar, una neurona con dendritas, su disparo y su axón, que se conecta con las siguientes dendritas, pero aquí ya estamos viendo que la naturaleza nos dice que hay que pensar más y plantear cómo hacer las representaciones de estas conexiones que nos presenta la naturaleza, un poco diferente pero tal vez con resultados más satisfactorios.

## Sinapsis

Aquí veremos más a detalle cómo una neurona recibe o transmite información a otras neuronas, donde para un solo disparo están participando un montón de elementos que veremos más adelante.

El momento en que dos neuronas transmiten información se llama **sinapsis** y es mediante conexiones que se dan en las terminales del axón (vesículas sinápticas) de la neurona presináptica hacia la postsináptica. Es importante notar que estas neuronas no tienen contacto anatómico, sino que están separadas por un espacio muy pequeño, **la brecha sináptica**. Lo que sucede en estas conexiones es un intercambio electroquímico que produce cambios de polaridad a lo largo la membrana.



**Figura 1.9** Part of neurons in Spanish, Dana Scarinci Zabaleta, 24 febrero 2019, Wikimedia Commons, [https://commons.wikimedia.org/wiki/File:Part\\_of\\_neurons\\_in\\_Spanish.svg](https://commons.wikimedia.org/wiki/File:Part_of_neurons_in_Spanish.svg), OpenStax, CC0

Clasificación de sinapsis, las terminales del axón de la neurona presináptica puede hacer contacto con la neurona postsináptica en:

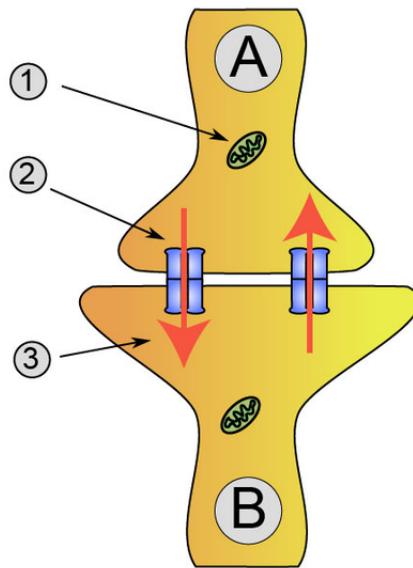
## 1. Neurona biológica

1. su dendritas, axodendrítica.
2. su cuerpo (soma), axosomática.
3. su axón, axoaxónica.

Distingamos entre dos tipos de sinapsis:

**Sinapsis eléctrica:** las membranas de las células pre y posinápticas se unen en la brecha sináptica por una unión tipo gap, o unión comunicante, que son pequeños canales que permiten el paso de iones.

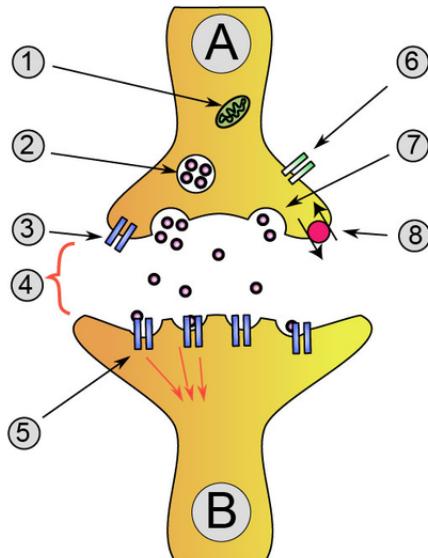
1. Posee una transmisión bidireccional de los potenciales de acción.
2. Sincronización en la actividad neuronal, lo cual hace posible una acción coordinada.
3. Los potenciales de acción pasan a través del canal proteico directamente sin necesidad de la liberación de los neurotransmisores, por tanto, es más rápida.



**Figura 1.10** Synaptical transmission (electrical). Neurona A transmisora, Neurona B receptora, 1. Mitocondria, 2. Uniones gap formadas por conexinas, 3. Señal eléctrica, Nrets commonswiki, 23 September 2005, Wikimedia Commons, [https://commons.wikimedia.org/wiki/File:Synapse\\_diag2.png](https://commons.wikimedia.org/wiki/File:Synapse_diag2.png), Inkscape 0.42, CC-BY-SA 3.0

**Sinapsis química:** la neurona libera moléculas neurotransmisoras a otra neurona adyacente en un pequeño espacio (la brecha sináptica) ver 1.11. Se puede resumir en cuatro etapas principales:

1. Un potencial de acción llega al botón terminal proveniente desde cono axónico.
2. Los neurotransmisores contenidos en las vesículas que están él los botones terminales, son liberados en la brecha sináptica y se dispersan.
3. Cada neurotransmisor se une a su receptor ubicado en la membrana de la neurona postsináptica.
4. El exceso de neurotransmisores que queda en el espacio sináptico es degradado o recaptado.



**Figura 1.11** Synaptical transmission (chemical). Neurona A transmisora, Neurona B receptora, 1. Mitocondria, 2. Vesícula sináptica llena de neurotransmisor, 3. Autorreceptor, 4. Brecha sináptica, 5. Receptor de neurotransmisores, 6. Canal de calcio, 7. Neurotransmisor liberador de vesículas fusionadas, 8. Bomba de recaptación de neurotransmisores, Utilisateur:Dake, 23 September 2005, Wikimedia Commons, [https://commons.wikimedia.org/wiki/File:Synapse\\_diag1.png](https://commons.wikimedia.org/wiki/File:Synapse_diag1.png), Inkscape 0.42, CC-BY-SA 3.0

Detallando el proceso de la sinapsis química, la neurona postsináptica está recibiendo un montón de señales por la liberación de neurotransmisores tanto de sus vecinos, como lo que ella misma va intercambiando, una vez que están generando el efecto completo de cambiar la polarización de la membrana, van a provocar que la neurona haga un disparo eléctrico. En el cuerpo están llegando estos intercambios de iones que se suman en el cono axónico, empiezan a viajar a través del axón, en las vainas de mielina (donde se refuerza la señal). Aquí hacemos mención por primera vez de los iones positivos: sodio y potasio, estos iones lo que hacen es, que *la membrana tenga una cierta carga la mayor parte del tiempo*. Cuando salen tres sodios entran dos potasios, entonces siempre hay

## 1. Neurona biológica

más positivos afuera que en el interior de la neurona, es decir, por lo general *tiene una carga más negativa que su entorno* (ver 1.17). Cuando ocurre un disparo de la neurona y se da el cambio de polarización en la membrana, se abren sus poros/ canales. El hecho que los canales abran o cierren depende de varios cambios que puedan estar ocurriendo alrededor de la neurona, en particular los que transmiten el disparo eléctrico, reaccionan ante el cambio de potencial que ocurrió en la membrana de la neurona.

La señal va pasando por los nodos de Rainvier, se refuerza y pasa por los canales iónicos ya abiertos, hasta finalmente llegar a la sinapsis a esto le llaman la **conducción a saltos** (ver 1.7). Ahora lo que ocurre al final del recorrido es que, el cambio de electricidad otra vez provoca que unas vesículas, que están en el interior de la neurona, que contienen neurotransmisores, se peguen a la membrana axónica y se liberen esos neurotransmisores nuevamente a otra neurona.

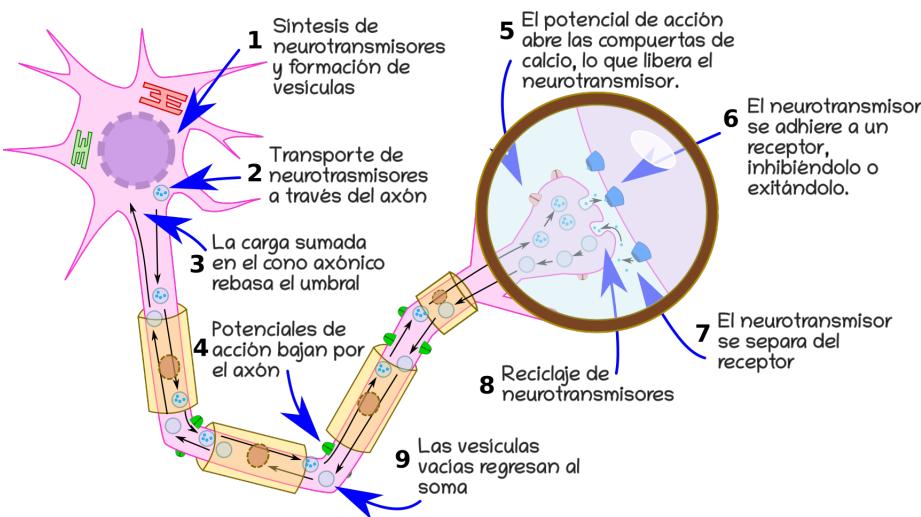
Entonces la información le va a llegar a la neurona vecina, en la forma de neurotransmisores que fueron liberados (en lo que sea que lo haya recibido, típicamente son dendritas, pero podría ser su cuerpo o su axón), eso es lo que va a percibir la otra neurona y otra vez esta otra neurona va a empezar a sumar los efectos de estos neurotransmisores, para que en algún momento decida a lo mejor disparar y otra vez provocar que se liberen neurotransmisores a su final e influir con otras neuronas.

### Neurotransmisión

Cuando la neurona no está mandando señales eléctricas, tiene un potencial de reposo, su diferencia de potencial entre el interior y el exterior de la neurona, es más negativo en el interior y más positivo (o menos negativo) en el exterior. En el caso de las sinapsis químicas, llega un disparo y se altera el potencial de la membrana, entrarán las células de calcio y entra la participación de las vesículas para liberar neurotransmisores. Los neurotransmisores están flotando en la brecha sináptica, viajan hasta adherirse a los receptores de la neurona posináptica, en ese momento están alterando el intercambio normal que existe entre iones en el interior y en el exterior de la célula y van a cambiar las cargas netas que hay adentro y afuera. Este es un cambio local que está ocurriendo en las espinas dendríticas (pequeñas prolongaciones citoplasmáticas).

este cambio en sí es una especie de transferencia de información, pero muy local. Se puede distinguir entre dos efectos en la membrana:

- **El efecto excitatorio**, despolariza la membrana postsináptica es decir ahora va a ser más propensa a disparar porque ya le cambió la diferencia de potencial que tenía.
- **El efecto inhibitorio**, hiperpolariza la membrana postsináptica, es decir, va a incrementar la diferencia de potencial entre el exterior e interior, pero de tal manera que ahora ya no va a querer disparar esta neurona.



**Figura 1.12** Esquema detallado de una neurotransmisión.

De estos efectos también va a darse el efecto de la **plasticidad**, que es cuando dos neuronas tienden a excitarse juntas, después de esta conexión se va a tender fortalecer, sí más bien tienden a inhibirse lo que va a suceder después es que estos canales empiezan a encoger, haciendo que se reduzcan y ya no dispare.

Ejemplos de neurotransmisores: serotonina, dopamina, oxitocina, endorfinas, adrenalina.

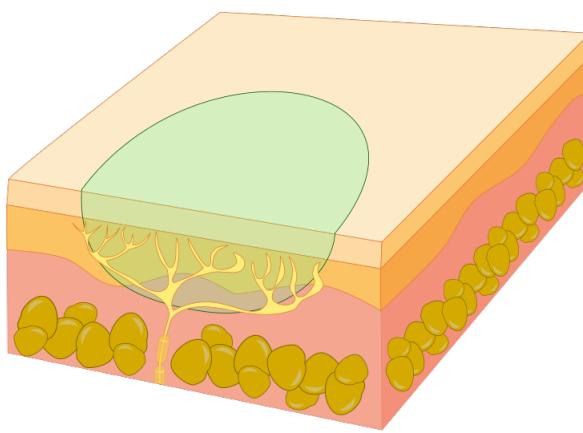
## Campos receptivos

Aquí lo que nos interesa es, en qué región puede ser afectada una neurona. Se define un campo receptivo, como la región en la periferia sensorial dentro de la cual los estímulos pueden, influir la actividad de las células sensoriales (ver 1.13). Hay diferentes niveles donde pueden aparecer los campos receptivos tanto cerca de la piel, cerca del gusto, el olfato, donde las neuronas van a estar asociadas con otras células que les pueden ayudar, que son sensibles a los cambios correspondientes, a veces la misma neurona va a tener alguna protuberancia especializada. También podemos encontrarlos más hacia adentro del nivel de procesamiento, no necesariamente todos van a estar pegados a la parte sensorial física.

Comprenden a los receptores sensoriales que alimentan a las neuronas sensoriales, pueden ser:

- Receptores específicos en una neurona como protuberancias especializadas.
- Conjuntos de receptores capaces de activar una neurona mediante conexiones sinápticas.

## 1. Neurona biológica

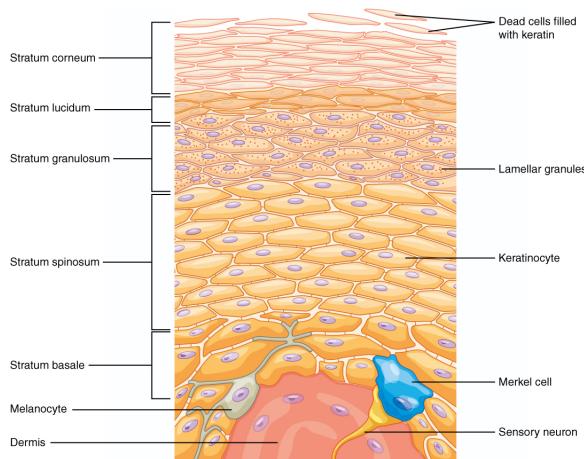


**Figura 1.13** Campo receptivo. Se muestra una región que está bajo cierto estímulo, las terminales de la neurona están recibiendo información y transmitiéndola.

- Describen la ubicación donde debe estar presente un estímulo sensorial para licitar una respuesta desde una célula sensorial.

Ejemplos:

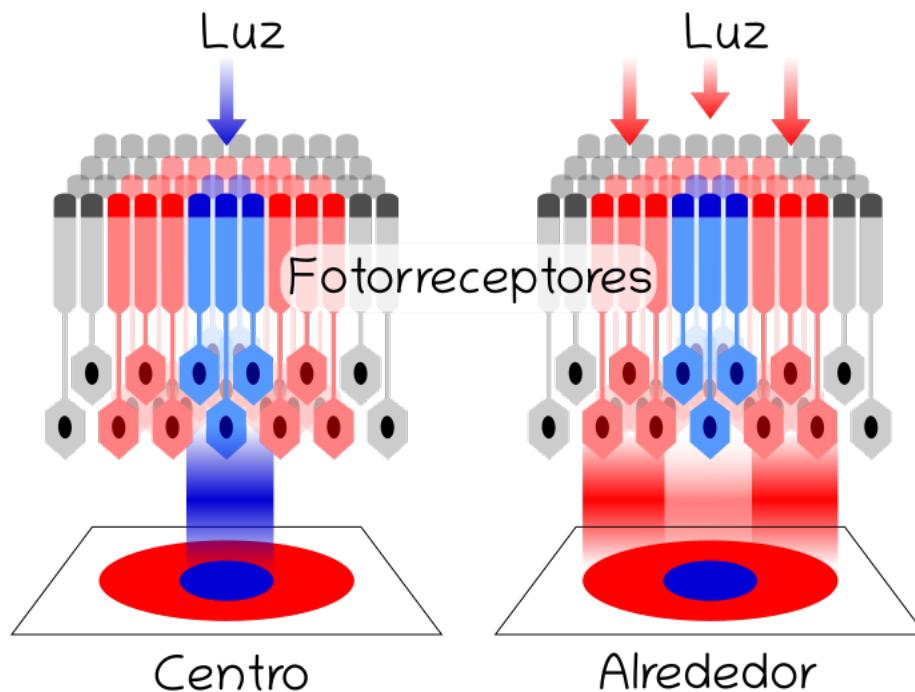
- En la piel tenemos células que nos están protegiendo en la epidermis, una de las células auxiliares **la célula de Merkel** que es sensible a la presión. Esta puede estar muy cerca a una neurona, sus terminales se activan de acuerdo a las acciones de la célula de Merkel y va a pasar la información (ver 1.14).



**Figura 1.14** Capas de la epidermis y en azul la célula de Merkel.

- El ojo, para procesamiento visual, actualmente se utiliza una de las redes neuronales más famosas que son las redes convolucionales, que están inspiradas en el ojo,

nosotros tenemos campos receptores donde hay unos fotorreceptores en los conos y los bastones que son sensibles a luces de diferentes colores a cambios de intensidad de la luz y que pueden detectar, por ejemplo, en una cierta región física si está llegando luz o por ejemplo, si llega en la periferia entonces va a inhibir el disparo de estos elementos, por otro lado, tenemos también su complemento que permite ser estimulado por las señales que llegan, como que en la parte de afuera de un círculo y más bien se inhiben con un estímulo en la parte de afuera (ver 1.15). Esta especie de células que tienen una posición física y geométrica relevante van a determinar cuando disparan o no las neuronas. Los siguientes niveles del cerebro se van a encargar de interpretar mejor el cambio de sombras, como a una persona que pasó corriendo, un auto que se está moviendo cerca o reconocer algún tipo de alimento.



**Figura 1.15** Capas de la epidermis y en azul la célula de Merkel.

## Señal eléctrica

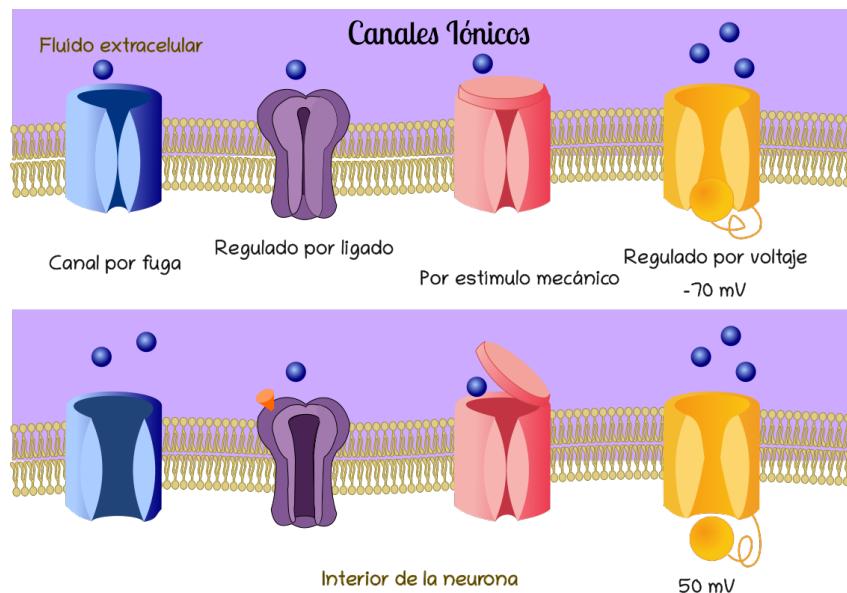
Veamos que pasa en los canales de iones y el paso de la señal eléctrica primero diferenciamos los tipos de compuertas iónicas:

- **Canal por fuga:** Estos se abren y cierran aleatoriamente, todo el tiempo están activos en la neurona, intercambiando por ejemplo: sodio y potasio.

## 1. Neurona biológica

- **Canal regulado por ligado:** Aquí se hace presente un neurotransmisor que es el que va a provocar que se abran o al revés impedir que se abran.
- **Canal por estímulo mecánico:** Permiten que pasen más iones o menos iones dependiendo, si se ejerció una presión, por ejemplo, con las neuronas cerca de la piel, las células de merkel.
- **Canales regulados por el voltaje:** Tienen el rol protagónico en la transmisión del pulso eléctrico (que se ha estado mencionando) describiendo uno de ellos, este canal tiene una pequeña compuerta abajo, que la puede cerrar independientemente del hecho de que el canal se abre o se cierre. Existen varias variantes de este tipo de canales regulados por el voltaje, la forma en que se están activando y desactivando sus compuertas, es lo que permite el paso del pulso.

Existen realmente una buena cantidad de iones presentes en el cerebro, pero los más protagonistas son precisamente el **potasio**, el **sodio**, el **cloro** y son los que vamos a utilizar para un modelo matemático de las neuronas.



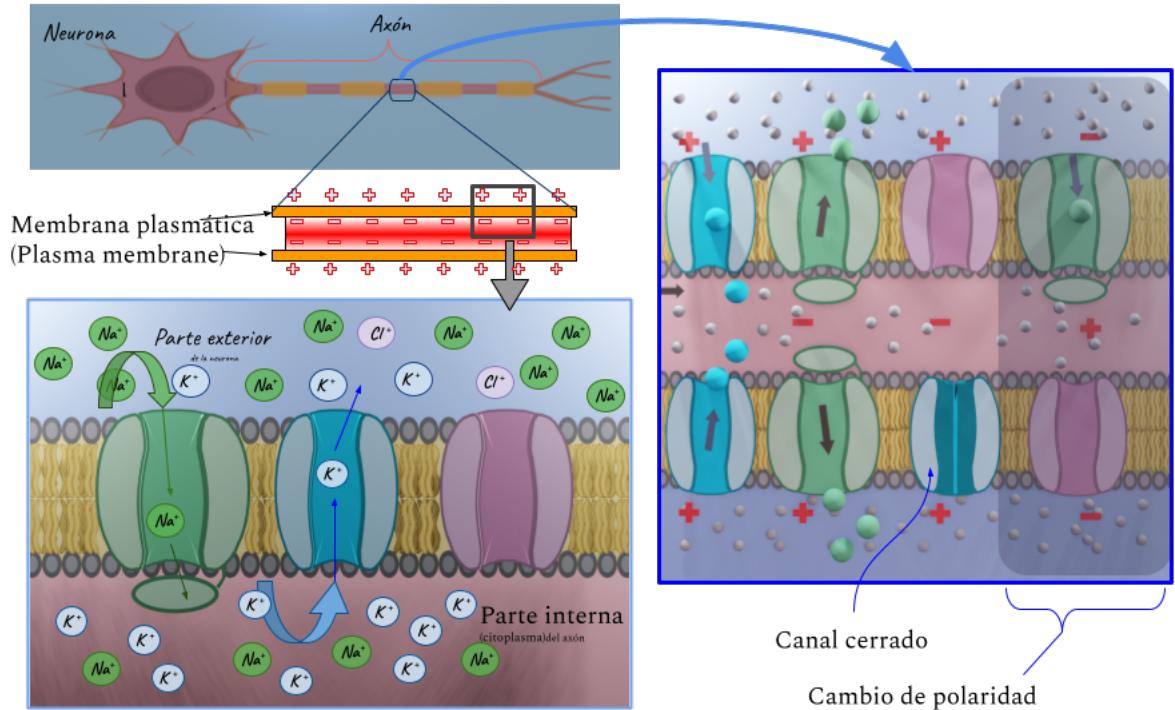
**Figura 1.16** Representación de la clasificación de los canales iónicos (más representativos).

Por último veamos brevemente **la neuroplasticidad**, es lo que nos permite el aprendizaje a largo plazo en el cerebro, es un mecanismo de aprendizaje del cerebro en el cual:

Cuando las neuronas se activan simultáneamente con frecuencia la conexión entre ellas se fortalece.

Este mecanismo constituye la principal inspiración para el diseño de las redes neuronales artificiales, concretamente en esto se inspiran los algoritmos de entrenamiento. Lo que

se hace es calcular, qué conexiones debemos reforzar y cuáles debemos de debilitar para que nuestras redes neuronales calculen las funciones que a nosotros nos interesan.



**Figura 1.17** Representación de la membrana axónica en potencial de reposo en la parte inferior izquierda, y en la parte derecha con un estímulo que genera el cambio de polaridad en la misma, así como el cierre de canales y paso de iones.

## 2 | Modelo de Hodgkin-Huxley

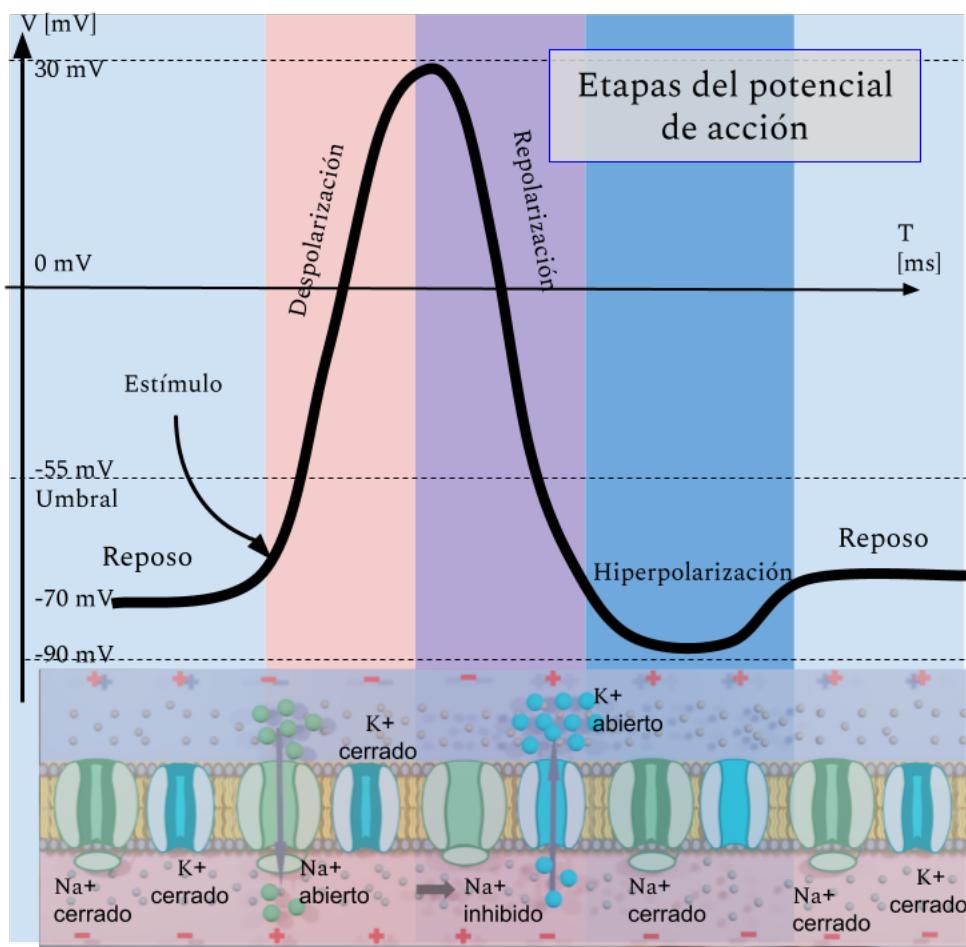
### Introducción

Esta sección se enfocará a la parte de transmisión de información y tipo de operaciones lógicas matemáticas que ocurren para que un cerebro pueda realizar cálculos. Específicamente se detallará la mecánica de los disparos de las neuronas, siendo estos una de las características más relevantes a la hora de modelar las redes neuronales artificiales. Si en algún momento de su vida han visto temas relacionados con: compuertas digitales, arquitectura de computadoras o diseño electrónico digital, les será más fácil abstraer el concepto, pues vamos a ver los procesos de paso de información a través de compuertas pero en un sistema biológico (de la naturaleza).

Notemos primeramente un impulso nervioso, recordemos que éste es una onda que avanza desde el cono axónico de la neurona hasta la neurona postsináptica. Esta onda electroquímica ocurre dada la diferencia de potencial entre la parte interna y externa de neurona, esta diferencia se da a consecuencia de las distintas concentraciones de iones en ambos lados de la membrana plasmática. Los estados en la membrana plasmática (del axón) se pueden diferenciar en, potenciales neuronales:

- **Potencial de reposo:** Es la diferencia de cargas en la membrana y está polarizada a -70mV. Es positiva por fuera ( $\text{Na}^+$ ) y negativa por dentro por  $\text{Cl}^-$  y proteínas- y no transmite señal.
- **Potencial de acción o membrana:** Un estímulo umbral de 55 mV, despolariza la membrana y abre los canales del  $\text{Na}^+$  y  $\text{K}^+$  y avanza la señal nerviosa, es un cambio muy rápido en la polaridad de la membrana de negativo a positivo y vuelta a negativo.

Retomando la sinapsis eléctrica, donde participan los canales iónicos y las entradas de las neuronas (dendritas) están siendo alteradas poco a poco, hasta que ocurre la suficiente carga (diferencia de potencial) en sus dendritas y en el cuerpo de la neurona, para que desde el cono axónico se dé un disparo o potencial de acción (spike), transmitiendo la información gracias a la apertura y cierre de ciertos canales de iones cargados. Este



**Figura 2.1** Representación gráfica de la respuesta de los canales iónicos de sodio ( $Na^+$  en verde) y potasio ( $K^+$  en azul) ante un estímulo de voltaje, dando como resultado un potencial de acción que viajará a lo largo de todo el axón.

cambio brusco de la diferencia de potencial, se nota en forma de un pulso eléctrico (ver 2.1), para saber más a detalle qué está ocurriendo en esta rápida elevación en la diferencia de potencial, se contará de dónde salió este modelo y por qué toma la forma que tiene.

Los primeros científicos que estudiaron el potencial de acción y dieron un modelo (de la unión sináptica eléctrica) fueron Alan Lloyd Hodgkin y Andrew Fielding Huxley alrededor de 1952, obteniendo un modelo matemático <sup>1</sup>, que intenta explicar qué es lo que estaba pasando en las neuronas. Ellos trabajaron con un calamar gigante (que puede medir hasta 4 metros de largo) que dado su gran tamaño, tiene un axón también bastante gigantesco, que recorre casi la mitad del cuerpo del calamar y su grosor es de medio milímetro, considerando el tamaño estándar de un axón de una neurona (1-20  $\mu\text{m}$ ). El axón del calamar gigante es tan grande que les permitió introducir dispositivos para medir

<sup>1</sup>El texto original de este experimento se puede encontrar en la siguiente url: <https://physoc.onlinelibrary.wiley.com/doi/pdf/10.1113/jphysiol.1952.sp004764>

## 2. Modelo de Hodgkin-Huxley

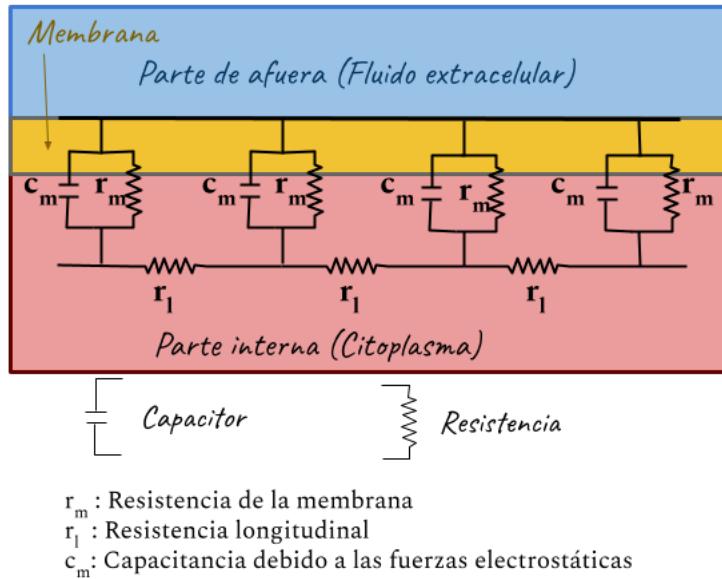
el voltaje, es decir, la diferencia de potencial entre, el interior de la neurona y la parte de afuera (el ambiente externo de la neurona). Con estas mediciones experimentales que lograron obtener, se pudo determinar qué pasaba con las cargas eléctricas tanto en el interior como en el exterior y así estudiar cómo se lograba la transferencia de electricidad cuando disparaba este pulso. Se dan cuenta de que podían modelar este comportamiento como un circuito eléctrico donde están corriendo estas corrientes, si bien aún no sabían todavía cuál era exactamente el mecanismo biológico por detrás, si observaron que había dos elementos protagonicos que serían el sodio y el potasio. Notaron que estos existen en diferentes concentraciones, en la parte de afuera y en la parte de adentro de las neuronas. Con esto nosotros podemos aprender también el por qué es importante consumir algo de sal y nunca estar bajos de potasio, pues estos dos elementos son indispensables para que las neuronas puedan transmitir sus señales.

### Membrana y canal

Hodgkin y Huxley se dedicaron a estudiar qué pasaba con las concentraciones de estos iones (sodio y potasio) en la parte de afuera o en la parte de adentro cuando empezaban a fluir las corrientes. El sistema parecía una especie de circuito eléctrico, se lo imaginaron como una especie de membrana porosa (lo cual es bastante cercano a lo que después se descubrió con la microscopía) y la forma en que lo vieron fue como un circuito eléctrico donde *la membrana está funcionando como un capacitor* que almacena ligeramente las cargas cuando están tratando de pasar de un lado hacia el otro y además con la cualidad que tenía de veces dejar pasar más iones y a veces no (semipermeable), modelan esto como una especie de *resistencias variables*. Bajo ciertas condiciones de voltaje de la diferencia de potencial entre la parte de afuera y la parte de adentro, estos canales permiten pasar más de estos iones (ya sean sodio, potasio o calcio) o, por el contrario, impiden su paso (ver 2.2).

Ahora se necesitan más detalles de la representación de los canales y toman en cuenta que el comportamiento de estas resistencias viene acompañado con un voltaje de reposo, en estos voltajes particulares cada tipo de ion (de la resistencia modelada) se estabiliza y ya no va a cambiar esta resistencia (ver 2.3).

Lo que observan es que el **ion de sodio ( $\text{Na}^+$ )** y su resistencia va a variar dependiendo del voltaje, a esto se le llama un **canal transitorio** porque en ciertos voltajes si puede pasar; si es muy bajo, no puede pasar y si rebasa un cierto umbral entonces se vuelve a tapar y ya no puede pasar. Lo que sucede con el **ion de potasio ( $\text{K}^+$ )** es que, puede salir si el voltaje está más allá de un cierto valor, si no, no pasan y va variando un poquito que tanto puede pasar, a esto se le llama **canal persistente**. Por estas características de que el potasio es un intervalo dentro de la recta y el sodio es a partir de cierto valor, por tanto, se les modelan de maneras ligeramente diferentes. Más adelante se descubrió porque tenían este comportamiento, básicamente el canal de potasio es una puerta hecha de cuatro subpuertas por donde los elementos pasan o no pasan, el canal de sodio es



**Figura 2.2** Un primer modelo de la membrana axónica modelada como circuito eléctrico. La parte amarilla es la membrana.

como una compuerta que está hecha de tres subpuertas que se pueden abrir y tiene aparte un tapón extra, que hace que aunque estas tres están abiertas bloquee toda la compuerta. Las neuronas están trabajando con muchos más iones aparte de estos dos, uno que destaca bastante es el caso del cloro ( $\text{Cl}^-$ ) que tiene carga negativa. Se tienen canales para intercambio aleatorio de otros iones, **L** un canal aleatorio (leaky).

Entonces con lo que ellos midieron experimentalmente, notaron cómo se estaban comportando estas resistencias dependiendo del voltaje o la diferencia de potencial que había entre ambos lados de la membrana. A partir de estas pudieron describir matemáticamente y simular los disparos que se conocen como potenciales de acción. Vamos a ver cuáles fueron estos conceptos de electricidad que se están utilizando para el modelo tenemos este concepto de potenciales eléctricos.

- Potenciales eléctricos  $E$  ó  $V$ ; resultan de la separación de cargas opuestas. Se mide en  $mV$ .

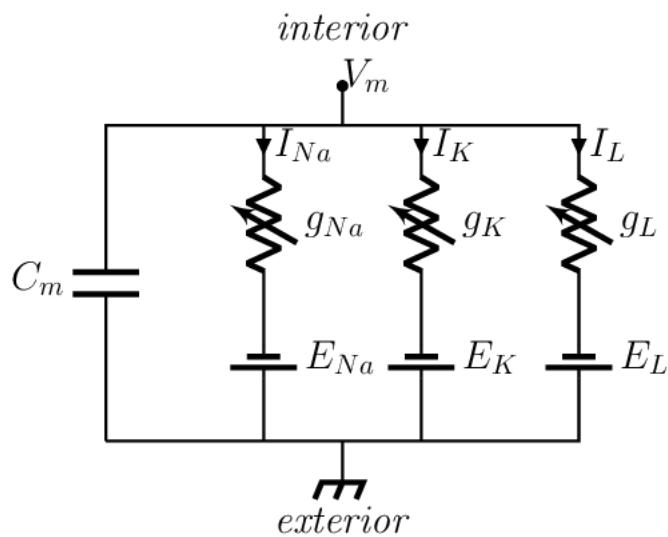
\*  $E_{(\text{Na}, \text{K}, \text{L})}$  voltaje en reposo para los iones de Na, K y L, o también conocido como potencial de inversión iónico, es el potencial de membrana en el que no hay flujo neto (total) de ese ion en particular de un lado de la membrana al otro.

\*  $V_m$ , el potencial eléctrico de la membrana.

- Corriente  $I$ ; Movimiento de cargas. Se mide en  $\mu\text{A}$ .

\*  $I_{(\text{Na}, \text{K}, \text{L})}$  corriente entrante a los canales de Na, K o L.

## 2. Modelo de Hodgkin-Huxley



**Figura 2.3** Modelo de la membrana axónica modelada como circuito eléctrico, con los distintos canales presentes y su voltaje de reposo.

- Resistencia  $R$ ; Medida de la oposición al movimiento de las partículas cargadas.
- Capacitancia o capacidad eléctrica  $C$ . Cantidad de energía eléctrica almacenada en un capacitor para una diferencia de potencial eléctrico dada.
  - \*  $C_m$  la capacitancia de la membrana.
- Conductancia  $g$ ; Inverso de la resistencia  $\frac{1}{R}$ , es decir, facilidad de transmisión de las partículas cargadas.
  - \*  $g_{(Na,K,L)}$  la conductancia del canal de sodio, potasio, cloro y la L también refiriéndose a otros canales de iones.

Lo que está pasando en los **potenciales eléctricos** es que hay mucho sodio en la parte externa de la membrana por ej. tres iones de sodio que son cargas positivas por dos iones de potasio que hay en la parte interna, entonces hay muchas más cargas positivas en la parte de afuera que las que hay en la parte de adentro y eso es lo que provoca la diferencia de cargas que es lo que estamos viendo como un potencial eléctrico.

La capacidad eléctrica o **capacitancia** es la que estamos utilizando para modelar la membrana conformada por lípidos, que es una capa de grasa y esa es la cantidad de energía eléctrica almacenada en un capacitor para una diferencia de potencial eléctrico dada. Éste comportamiento bastante interesante porque las cargas quedan almacenadas un momento, pero se van liberando poco a poco y se va descargando ese capacitor.

Durante el experimento con el axón, se le dieron cargas eléctricas directamente al axón y gracias a eso lograban ir midiendo que era lo que estaba pasando con las concentraciones de cargas afuera y adentro en el caso de las neuronas reales, esto en un ambiente

no alterado ocurre cuando entran en juego los neurotransmisores y provocan que haya cambios, en estas corrientes. Entonces Hodgkin y Huxley jugaron el rol que tendrían que jugar usualmente los *neurotransmisores* para abrir otras compuertas. Nosotros en la manera en la que lo vamos a simular es precisamente con estas corrientes que son las que se están poniendo en el experimento y vamos a ver cómo reacciona el axón.

## Potenciales de Nerst o de reposo

Los Potenciales de Nerst o de reposo son los potenciales a los cuales el flujo neto de iones a través de los canales abiertos es cero. Aquí vemos precisamente porque estamos utilizando la *E* generalmente la vamos a utilizar para referirnos a la diferencia de potencial entre la parte de afuera de la célula y la parte de adentro, las vamos a utilizar para representar a aquellos voltajes donde cada una de las compuertas encontrarán su equilibrio. Estos voltajes son distintos para cada una de las compuertas, esto va a provocar precisamente la dinámica de la de la neurona, por ejemplo:

- $E_{Na} \ 50mV$
- $E_{Ca} \ 150mV$
- $E_K \ 80mV$
- $E_{Cl} \ 60mV$

Aquí vemos que el sodio estaría su equilibrio en un valor positivo, el calcio que es el que va a jugar un rol de que se activen los neurotransmisores y se transmita el disparo, observamos que el voltaje tendría que ser bastante positivo. El potasio que es el que usualmente está trabajando intercambiándose casi todo el tiempo en la neurona, veremos que el punto de equilibrio usual de la neurona anda por los  $-76mV$  y el del cloro. Cada uno de estos canales pues está tratando de jalar la dinámica hacia su potencial de equilibrio y no hay precisamente un acuerdo entre ellos y eso es precisamente lo que hace que las neuronas cobren "vida".

## Modelo de la membrana como bicapa de lípidos

La membrana de una neurona es modelada como un elemento de un circuito con capacitancia  $C_m$  y potencial  $V$ , las corrientes que fluye a través de la bicapa lipídica están regidos por las siguientes ecuaciones:

$$I_m = C_m \frac{dV_m}{dt} \quad (2.1)$$

## 2. Modelo de Hodgkin-Huxley

Esta sería la ecuación principal (2.1) donde  $\frac{dV_m}{dt}$  está representando el cambio voltaje en la membrana respecto al tiempo.

$$C_m \frac{dV_m}{dt} = -g_{Na}m^3h(V - E_{Na}) - g_Kn^4(V - E_K) - g_L(V - E_L) + I_{ext} \quad (2.2)$$

Cada una de las partes del lado izquierdo de la ecuación 2.2 corresponde a las compuertas de los canales y la corriente de un estímulo externo que pueda influir a la membrana (este estímulo siempre será desde el exterior hacia el interior).

Retomando lo escrito anteriormente el canal de sodio es una compuerta compuesta de **tres** subpuertas y una subpuerta que actúa como tapón y el canal de potasio es una compuerta compuesta de **cuatro** subpuertas iguales, con esto podemos notar claramente que las conductancias sean representadas como:

- $\frac{1}{R_{Na}} = g_{Na} * m^3 * h$  donde  $g_{Na}$  es una constante que representa el valor de la conductancia máxima,  $m$  es la proporción de los canales de sodio abiertos (representa la concentración de sodio) y nos indica la activación (subpuertas abiertas) del canal,  $h$  es el “tapón” de la compuerta que puede impedir el paso de iones independientemente de las otras tres subpuertas, es decir la inactivación (compuerta bloqueada). Los movimientos combinados de  $m$  y  $h$  son los que controlan la compuerta de sodio.
- $\frac{1}{R_K} = g_K * n^4$  donde  $g_K$  es una constante que representa el valor de la conductancia máxima,  $n$  es la proporción de los canales de potasio abiertos (representa la concentración de potasio) y nos indica la activación del canal de potasio.
- $g_L$  es una constante, de los canales por fuga, que representa la concentración de los demás iones que pasan por la membrana.

Ahora  $m$ ,  $n$  y  $h$ , son variables de activación que describen la probabilidad de que los canales iónicos estén abiertos, se puede describir mediante las siguientes ecuaciones diferenciales ordinarias:

$$\frac{1}{\gamma(T)} \frac{dn}{dt} = \alpha_n^\infty(V)(1 - n) - \beta_n(V)n = \frac{n(V) - n(t)}{\tau_n(V)} \quad (2.3)$$

$$\frac{1}{\gamma(T)} \frac{dm}{dt} = \alpha_m(V)(1 - m) - \beta_m(V)m = \frac{m^\infty(V) - m(t)}{\tau_m(V)} \quad (2.4)$$

$$\frac{1}{\gamma(T)} \frac{dh}{dt} = \alpha_h(V)(1 - h) - \beta_h(V)h = \frac{h^\infty(V) - h(t)}{\tau_h(V)} \quad (2.5)$$

## 2.4 Modelo de la membrana como bicapa de lípidos

Donde la ecuación 2.3 representa al canal de potasio y las ecuaciones 2.4 y 2.5 representando al canal de sodio tomando en cuenta que tiene dos tipos de subpuertas.

Las expresiones de  $\alpha$  y  $\beta$  están dadas por las siguientes ecuaciones:

$$\begin{aligned}\alpha_n &= \frac{0.01(10 - V)}{\exp\left(\frac{10 - V}{10}\right) - 1} & \beta_n &= 0.125 \exp\left(-\frac{V}{80}\right) \\ \alpha_m &= \frac{0.01(25 - V)}{\exp\left(\frac{25 - V}{10}\right) - 1} & \beta_m &= 4 \exp\left(-\frac{V}{18}\right) \\ \alpha_h &= \frac{0.07}{\exp\left(-\frac{V}{20}\right)} & \beta_h &= \frac{1}{1 + \exp\left(\frac{30 - V}{10}\right)}\end{aligned}$$

Los factores  $\alpha$  y  $\beta$  se denominan como constantes de velocidad de transición.  $\alpha$  es el número de veces por segundo que se abre una puerta que está en estado cerrado, mientras que  $\beta$  es el número de veces por segundo que se cierra una puerta que está en estado abierto. Si la membrana tiene una carga negativa,  $\alpha$  debe aumentar y la  $\beta$  debe disminuir, cuando la membrana esté despolarizada.

Hasta ahora sabemos que en la bicapa de lípidos, una pequeña carga está pasando entre sus capas de grasa. También sabemos que la carga es almacenada por un breve periodo de tiempo, dando como resultado que la bicapa se comporte como un **capacitor**. Esta membrana también está con cierta resistencia al paso de corriente. Con esto tenemos el siguiente diagrama <sup>2</sup> 2.4

Tenemos dadas las siguientes ecuaciones:

$$I_C + I_R - I_{ext} = 0 \quad (2.6)$$

$$C \frac{dV}{dt} + \frac{V}{R} - I_{ext} = 0 \quad (2.7)$$

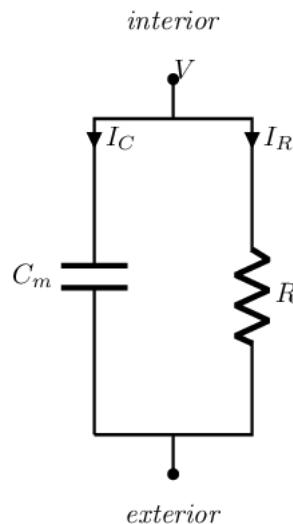
$$C \frac{dV}{dt} = -\frac{V}{R} + I_{ext}$$

Ahora por la ley de corriente de Kirchhoff <sup>3</sup> tenemos que la suma de las corrientes del capacitor y la resistencia debe ser cero por la conservación de corriente y si consideramos

<sup>2</sup>Otra explicación más profunda de las ecuaciones dadas partir del diagrama 2.4la podemos encontrar en [https://neurowiki.case.edu/wiki/Action\\_Potential\\_IV:\\_Hodgkin-Huxley\\_Equations\\_and\\_Other\\_Conductances](https://neurowiki.case.edu/wiki/Action_Potential_IV:_Hodgkin-Huxley_Equations_and_Other_Conductances)

<sup>3</sup>La ley de la corriente de Kirchhoff dice que la suma de todas las corrientes que fluyen hacia un nodo es igual a la suma de las corrientes que salen del nodo

## 2. Modelo de Hodgkin-Huxley



**Figura 2.4** Modelo de la bicapa de lípidos donde,  $V$  son los cambios de voltaje en la membrana que es el potencial eléctrico,  $I_C$  es la corriente del capacitor,  $I_R$  es la corriente de la resistencia,  $C_m$  es la capacitancia de la membrana,  $R$  es la resistencia.

un factor adicional de una corriente externa aplicada o administrada a la neurona, tenemos la ecuación 2.6.

Después tenemos la relación entre la diferencia de potenciales, que almacena energía y la carga eléctrica que guarda, donde:  $C$  es la capacidad, medida en faradios,  $Q$  la carga eléctrica almacenada, medida en culombios,  $V$  la diferencia de potencial medida en voltios. Entonces  $C = Q/V$ , despejando a  $Q$  tenemos  $Q = CV$  y derivando de ambos lados respecto al tiempo y considerando que  $C$  es una constante al ser una propiedad de la membrana  $\frac{dQ}{dt} = C \frac{dV}{dt}$ . Como la definición de corriente es el cambio de carga en el tiempo tenemos que  $I_C = C \frac{dV}{dt}$ . Notemos finalmente la corriente de la resistencia  $I_R$ , recordando la ley de Ohm <sup>4</sup> la podemos rescribir como  $\frac{V - V_{rest}}{R}$ . Sustituyendo de lo anterior en la ecuación 2.6 se obtiene la siguiente ecuación:

$$C \frac{dV}{dt} + \frac{V - V_{rest}}{R} - I_{ext} = 0 \quad (2.8)$$

$$C \frac{dV}{dt} = -\frac{V - V_{rest}}{R} + I_{ext}$$

<sup>4</sup>La ley de Ohm establece que la diferencia de potencial  $V$  que aplicamos entre los extremos de un conductor determinado es directamente proporcional a la intensidad de la corriente  $I$  que circula por el conductor, es decir  $V = R * I$ . Notemos también que  $V = V_m - V_{rest}$

Ahora multiplicando todo por R:

$$RC \frac{dV}{dt} = -V + (V_{\text{rest}} + RI_{\text{ext}}) \quad (2.9)$$

Denotando  $RC$  como la constante de tiempo  $\tau$  y tomando en cuenta que en cuanto se aplica la corriente va a empezar a cambiar el voltaje poco a poco hasta establecerse en un voltaje de equilibrio (ahí se va a quedar quieta). Entonces cuando el voltaje ya no está cambiando con el tiempo quiere decir que su derivada con respecto al tiempo es cero. Observemos que  $dV/dt = 0$  significaría que  $V$  es igual a infinito y que el voltaje en el estado estacionario cuando,  $dV/dt = 0$  depende del potencial de reposo y del producto entre la resistencia con la corriente externa suministrada, entonces tenemos que:

$$V_{\infty} = V_{\text{rest}} + RI_{\text{ext}} \quad (2.10)$$

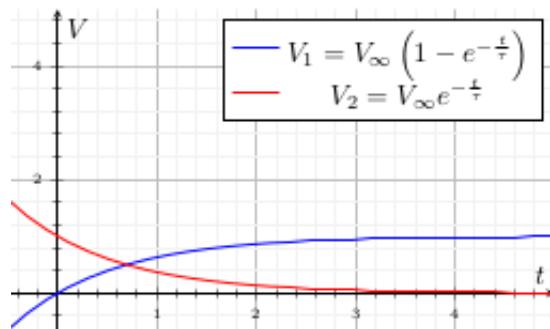
Sustituyendo con 2.10 y la constante, en 2.9 tenemos que:

$$\tau \frac{dV}{dt} = -V + V_{\infty} \quad (2.11)$$

## Las conductancias iónicas

Nuestro objetivo aquí es encontrar ecuaciones que describan las conductancias con precisión razonable y lo suficientemente simples para el cálculo teórico de *el potencial de acción* y *el período refractario*.

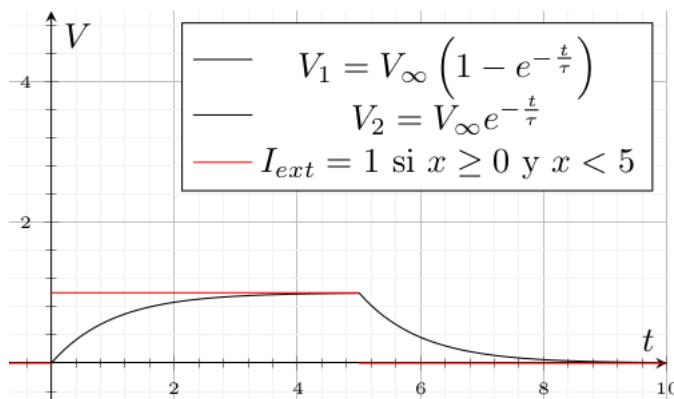
Si tomamos las ecuaciones diferenciales anteriores notamos que las soluciones tienen este tipo de forma 2.5:



**Figura 2.5** Soluciones para el pulso.

Donde si estamos aplicando una corriente externa lo que sucede es lo que estamos viendo en azul, un exponencial que va creciendo y que tiende hacia un cierto valor límite que sería de infinito. Si dejamos de aplicar la corriente externa entonces ahora tendremos

## 2. Modelo de Hodgkin-Huxley

**Figura 2.6** Soluciones para el pulso escalón.

un exponencial, pero que tiende hacia el cero y se va a estabilizar en cero, lo que vemos en rojo.

Simulando lo que hicieron Hodgkin y Huxley que fue al axón de repente darle un toque, siendo en el origen de la gráfica (que visualizamos en la figura 2.6) la parte en la que le están dando el toque al axón, momentos antes estaba quieta la neurona de repente le aplican una cierta cantidad de electricidad y va a empezar a cambiar el comportamiento de los canales la porosidad de la membrana, vamos a ver que empieza a incrementarse la diferencia de potencial hasta que llegan a un nuevo equilibrio (alrededor de  $t = 3$ ) y si siguieran dándole el toque en esa cantidad se quedaría ahí la neurona, ya no veríamos más cambios lo que va a suceder entonces es que, retiramos las pinzas (se le deja de dar el toque) y los canales otra vez van a empezar a regresar a la normalidad y vamos a ver un descenso en adelante.

Entonces hasta aquí ya tenemos la idea de cómo va a reaccionar la neurona ante cierto estímulo; sin embargo, esto que acabamos de ver en las gráficas sería como si tuviéramos un solo tipo de canal, ahora qué pasa si consideramos que tenemos diferentes tipos de canales pasando iones, en condiciones distintas. Aquí es donde va a importarnos el hecho de que existen diferentes tipos de canales con voltajes de equilibrio diferente. Retomando a los potenciales de Nerst  $E_{Na}$ ,  $E_K$ ,  $E_L$  notemos que están dados por:

$$E = \frac{k_B T}{zq} \ln \frac{[\text{adentro}]}{[\text{afuera}]} \quad (2.12)$$

Estos potenciales están relacionados con las características termodinámicas, en la ecuación anterior  $k_B$  la constante de Boltzman,  $q$  es la carga del ion, y  $z$  es el número de iones. El logaritmo natural representando el promedio de cuántos elementos tenemos en la parte de adentro con respecto a cuántos elementos tenemos en la parte de afuera.

Considerando los diferentes puntos de equilibrio en los cuales se puede encontrar la diferencia de potencial en la membrana, vamos a distinguir entre tres estados de esta (también se puede ver en 2.1):

**1. Polarizada** en su estado de reposo con  $V < 0$  ( $V \approx -70\text{mV}$ ).

- Su estado de reposo, cuando la neurona no está haciendo nada simplemente están corriendo los sodios y entran los potasios.

**2. Despolarizada** cuando  $V \geq 0$ .

- Cuando en sus dendritas y en el cuerpo de la neurona se acumula una carga muy grande (un voltaje eléctrico, disparo), se abre la compuerta de sodio y van a empezar a entrar el sodio, esta diferencia de potencial que existía entre lo fuera y lo adentro se va a reducir de hecho se puede llegar a reducir bastante dependiendo de la carga que se le aplique.
- Iones positivos entran a la membrana.
- Valores positivos en la diferencia de potencial.

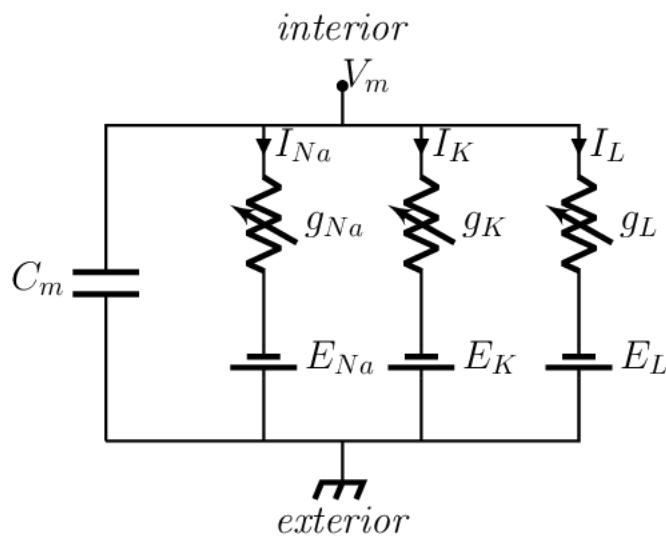
**3. Hiperpolarizada** cuando la diferencia de potencial incrementa su magnitud  $V \ll 0$ .

- En cuanto se despolarice van a empezar interacciones entre los diferentes tipos de canales que lo que van a intentar hacer es regresar a la neurona en su estado normal.
- Los canales de potasio abren sus compuertas, provocando que salga el potasio que está dentro del axón (en el citoplasma).
- Iones positivos salen de la membrana.
- Si antes estaba quieta a los  $-70\text{mV}$ , ahora va a quedar todavía más abajo alrededor de  $-90\text{mV}$ . Esto va a permitir un fenómeno que se le conoce como *el periodo de refracción* y ese periodo sirve para que simplemente se lance un disparo y que el comportamiento eléctrico no se rebote otra vez en dirección contraria en la neurona, va a quedar muy quieta la neurona durante un rato y después regresará otra vez a su estado de equilibrio.

## Modelo de las compuertas iónicas controladas por voltaje

Retomando el modelo del circuito eléctrico modelando la membrana, junto con los canales y los iones, volvamos a verlo ahora en la 2.7

## 2. Modelo de Hodgkin-Huxley



**Figura 2.7** Modelo de la membrana axónica modelada como circuito eléctrico, con los distintos canales presentes y su voltaje de reposo.

Recordemos brevemente las definiciones de los dos tipos de canales protagonistas en el modelo:

### Definición 2.1

*Canal persistente* Tiene un solo tipo de compuerta y dos estados posibles:

1. **Activado**
2. **Desactivado**

### Definición 2.2

*Canal transitorio* Tiene compuertas de activación e inactivación, y tres estados:

1. **Activado** Ambas compuertas abiertas.
2. **Desactivado** Compuerta de activación cerrada, inactivación abierta.
3. **Inactivada** Compuerta de inactivación cerrada.

Y retomando la primera ecuación diferencial donde tenemos, por un lado, la corriente que está pasando a través del capacitor, por otro lado, vamos a tener las corrientes que están circulando a través de los diferentes canales,

$$C_m \frac{dV_m}{dt} = -g_{Na} m^3 h(V_m - E_{Na}) - g_K n^4 (V_m - E_K) - g_L (V_m - E_L) + I_{ext} \quad (2.2)$$

Las capacitancias y variables del lado izquierdo están explicadas en la sección [Modelo de la membrana como bicapa de lípidos](#), aquí vamos a retomar las ecuaciones 2.3,2.4,2.5 de esa misma sección, (recordemos que estas ecuaciones describen la probabilidad de que los canales iónicos estén abiertos) que son las siguientes:

$$\frac{1}{\gamma(T)} \frac{dn}{dt} = \alpha_{n^\infty}(V)(1-n) - \beta_n(V)n = \frac{n(V) - n(t)}{\tau_n(V)} \quad (??)$$

$$\frac{1}{\gamma(T)} \frac{dm}{dt} = \alpha_m(V)(1-m) - \beta_m(V)m = \frac{m^\infty(V) - m(t)}{\tau_m(V)} \quad (??)$$

$$\frac{1}{\gamma(T)} \frac{dh}{dt} = \alpha_h(V)(1-h) - \beta_h(V)h = \frac{h^\infty(V) - h(t)}{\tau_h(V)} \quad (??)$$

Ahora notemos los elementos en estas ecuaciones anteriores con  $\text{ion}$  pudiendo denotar las compuertas del potasio  $n$  o del sodio, ya sea  $m$  o  $h$ :

- $\frac{1}{\gamma(T)}$  Este es el coeficiente de escala temporal, dependiente de la temperatura los por eso está apareciendo aquí una  $t$ . Para las simulaciones que nosotros vamos a hacer vamos a pensar que estamos en una temperatura fija.
- $\alpha_{\text{ion}}(V)$  probabilidad de que una compuerta transite de cerrada a abierta.
- $\beta_{\text{ion}}(V)$  probabilidad de que una compuerta transite de abierta a cerrada.
- $\text{ion}^\infty(V)$  probabilidad de compuerta abierta en el equilibrio cuando  $t \rightarrow \infty$ .
- $(\text{ion})$  Probabilidad de que cada compuerta ( $n, m, h$ ) esté abierta.
- $(1 - \text{ion})$  Probabilidad de que cada compuerta ( $n, m, h$ ) esté cerrada.
- $\tau_{\text{ion}}(V)$  Tiempo que toma llegar al equilibrio.

Lo que vamos a ver es que forma de escribir la ecuación depende precisamente del número de compuertas que tenían para poder abrirse y cerrarse. Reescribir la ecuación de esta manera lo que nos permite es medirlo en términos de estas probabilidades de que se abran y cierren las compuertas que serían

Esta probabilidad se empieza a alterar conforme cambiamos el voltaje, pero no va a llegar a su valor de equilibrio sino hasta después de pasado un cierto periodo.

## 2. Modelo de Hodgkin-Huxley

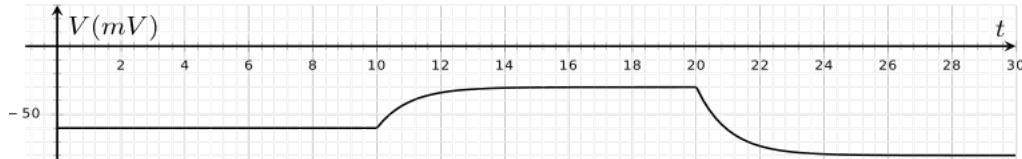
## Dinámica del voltaje durante un disparo

Ahora qué pasa cuando tomamos en cuenta que todas las compuertas están reaccionando al mismo tiempo. Notemos primeramente como está reaccionando la membrana ante un impulso eléctrico o voltaje, en la siguiente imagen 2.8.



**Figura 2.8** En la primera parte los canales están en un estado de reposo y la membrana está polarizada. En la segunda parte los canales han sido afectados por un impulso recibido desde el cono axónico, las compuertas de sodio se abren permitiendo el paso de iones de sodio al interior de la membrana y dejando a la membrana despolarizada. Momentos después la membrana llega a un estado de hiperoxialación donde intentará regresar al estado de equilibrio que tenía previamente, para esto la subpuerta de inactivación de sodio cerrará hasta no permitir el paso de sodio y el canal de potasio abrirá sus compuertas para dar salida a los potasios (iones positivos) que fluyan hacia el exterior de la membrana, así dejando el voltaje de la membrana incluso más negativo de lo que tenía durante su estado polarizado.

Ahora lo que vemos en la imagen 2.8 se grafica de manera un poco más apegada a lo que pasa en los experimentos de la siguiente forma en la imagen 2.9.

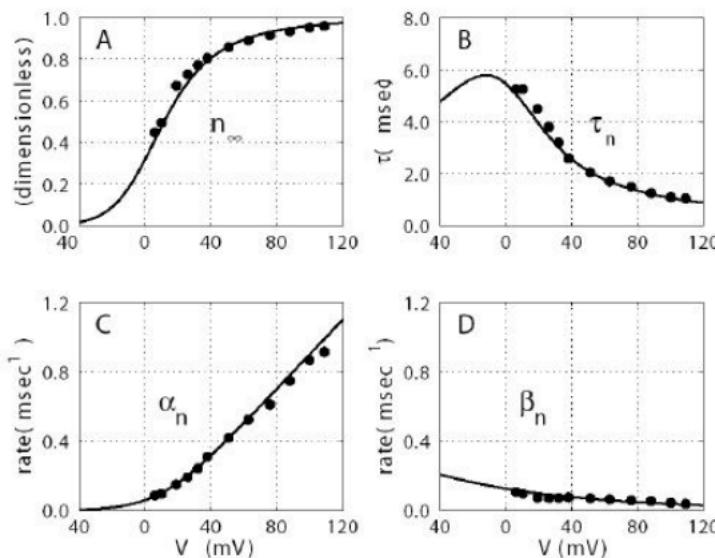


**Figura 2.9** Gráfica del cambio de voltaje en la membrana dado un impulso a través del tiempo. Cuando se rebasa el voltaje umbral, los canales de  $\text{Na}^+$  y  $\text{K}^+$  interactúan para producir una rápida despolarización de la membrana provocando una elevación del voltaje, para luego hiperoxializarla. Durante el momento de despolarización la membrana puede llegar a valores positivos en este caso en particular el estímulo no es tan grande y se queda en valores negativos.

Ahora lo anterior está en términos de lo deseado, veamos entonces que pasó en las mediciones de Hodgkin y Huxley con el axón en las siguientes gráficas 2.10, donde nos está mostrando como las subpuertas  $n$  y los factores  $\tau$ ,  $\alpha$  y  $\beta$  del canal de potasio se comportaron antes, durante y después del estímulo del voltaje. Recodemos que  $n$  nos indica la probabilidad que las compuertas de potasio estén abiertas, este es un factor

adimensional.  $\tau$  es el tiempo que tarda en llegar a un estado de equilibrio.  $\alpha$  y  $\beta$  las probabilidades que las compuertas del canal de potasio pasen de cerrados a abiertos y viceversa.

Entonces notamos que en el mismo periodo de tiempo, la membrana está en reposo y conforme va recibiendo el voltaje incrementa la probabilidad que las compuertas de potasio estén abiertas, es decir  $n$  va incrementando conforme al estímulo, mientras que el estado de reposo es claramente alterado provocando que el valor  $\tau$  disminuya considerablemente. La probabilidad que las compuertas de potasio pasen de un estado cerrado a uno abierto durante el proceso ( $\alpha$ ) aumenta prácticamente de forma exponencial, mientras que la probabilidad de que pasen de abierto a cerrado disminuye poco a poco. Con esto cumpliendo lo esperado en la dinámica del voltaje, al notar como está reaccionando el canal de potasio durante la polarización y despolarización (pulso).



**Figura 2.10** Medición experimental de los parámetros y ajuste manual de curvas. Imagen de Nelson 2004

Con estas medidas experimentales ellos dan con curvas paramétricas ajustadas a los factores  $\alpha$  y  $\beta$  expresadas de la siguiente forma:

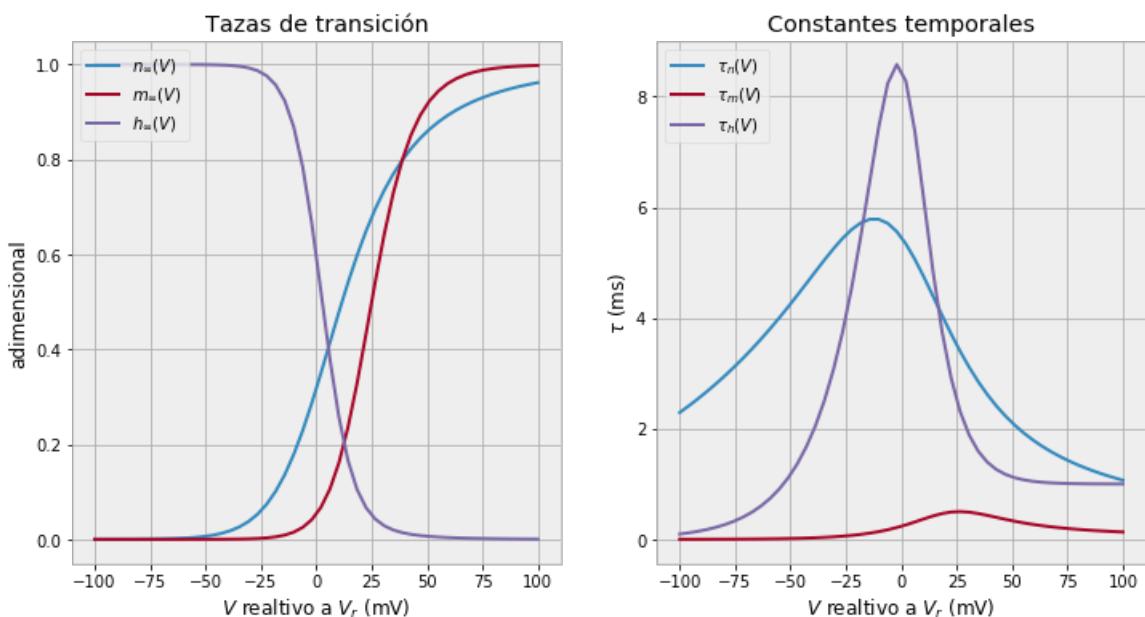
$$\begin{aligned} \alpha_n &= \frac{0.01(10-V)}{e^{\left(\frac{10-V}{10}\right)} - 1} & \beta_n &= 0.125e^{-\frac{V}{80}} \\ \alpha_m &= \frac{0.01(25-V)}{e^{\left(\frac{25-V}{10}\right)} - 1} & \beta_m &= 4e^{-\frac{V}{18}} \end{aligned}$$

## 2. Modelo de Hodgkin-Huxley

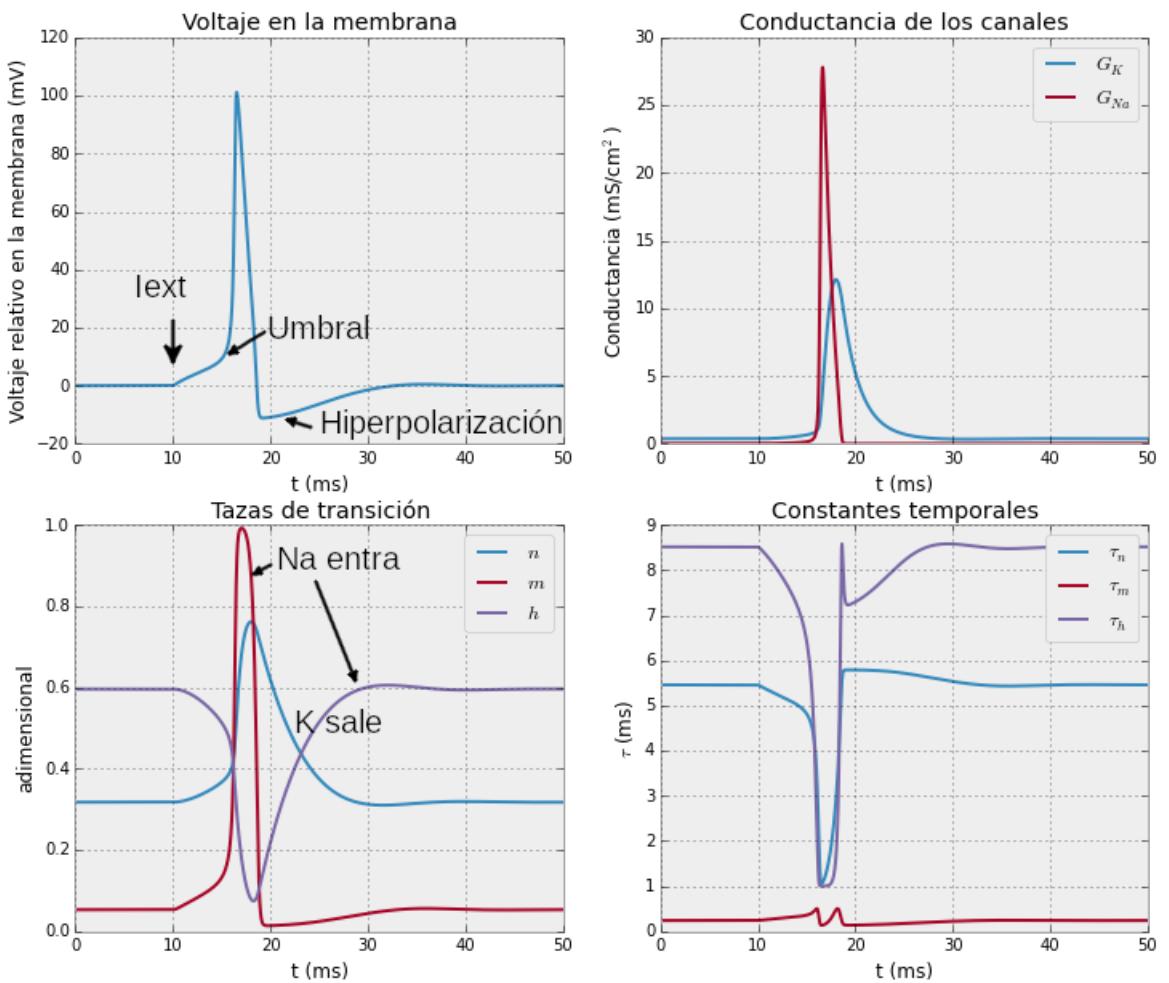
$$\alpha_h = 0.07e^{-\left(\frac{V}{20}\right)}$$

$$\beta_h = \frac{1}{e^{\left(\frac{30-V}{10}\right)} + 1}$$

Ahora veamos las gráficas de la dinámica del disparo pero ya con los las compuertas del potasio y sodio interactuando al mismo tiempo, en las figuras 2.11 y 2.12



**Figura 2.11** Activación e inactivación de los canales. Del lado izquierdo vemos que alrededor de los -50mV aumenta la probabilidad que las compuertas de potasio abran y los potasios salgan, el sodio tarda un poco más en reaccionar para dejar entrar a los sodios, y la probabilidad que quede no bloqueado disminuye, hasta bloquear por completo a los iones de sodio alrededor de los 50mV. Del lado derecho vemos como  $\tau$ , va interactuando conforme al voltaje, el  $\text{Na}^+$  reacciona rápidamente ante él impuso pues regresa rápidamente a su estado de equilibrio, esto indica porque aún que el sodio abre sus compuertas, la compuerta de bloqueo se cierra rápidamente, dejando inactivado el canal, luego el  $\text{K}^+$  reacciona más lentamente para regresar a su estado de equilibrio dejando más tiempo activado el canal y salgan los potasios.



**Figura 2.12** Medidas experimentales del estímulo y la reacción de todos los componentes de la membrana. El origen en estas gráficas realmente representa el punto de equilibrio que es alrededor de los -70mV. Del lado superior izquierdo mostrando como cambia el voltaje de la membrana respecto al tiempo, del lado superior derecho el comportamiento de la conductancia de los canales mostrando que aunque la conductancia del sodio reacciona más, es más breve, mientras que el potasio reacciona en menor medida pero durante más tiempo. En la parte inferior derecha la taza de transición de las compuertas  $n, m$  y  $h$ , mostrando que en cuando llega el impulso el primero en reaccionar es el sodio permitiendo que entre muy brevemente el sodio e inactivándose casi inmediatamente después, mientras que en ese momento los potasios se están liberando hacia la parte exterior de la neurona. Finalmente del lado inferior izquierdo mostrando como se comporta las constantes temporales, mostrando que las compuertas de activación de sodio no son tan afectadas comparado a la compuerta de bloqueo, notando algunos cambios bruscos en esta, mientras que el potasio si se nota más afectado y durante más tiempo pero recuperándose sin cambios bruscos.

## Simulación usando el método de Euler

En esta sección se listará como estamos resolviendo las ecuaciones diferenciales, para tener una simulación numérica y se trata del algoritmo de integración<sup>5</sup>. Las gráficas de la sección anterior se obtuvieron con el método de Euler que se describe más adelante.

---

**Algoritmo 1** Algoritmo de integración de Euler [Wells pp51].

---

```

1: function INTEGRADISPARO( $T, \Delta T, V_0, I_{ext}(t)$ )
2:   Inicializar arreglos de logitud  $T[]$ ,  $V[]$ ,  $n[]$ ,  $m[]$ ,  $h[]$ ,  $G_{Na}[]$ ,  $G_K[]$ ,  $\tau_n[]$ ,  $\tau_m[]$ ,  $\tau_h[]$  ←
   arreglo[numeroDePasos]
3:    $V[0] \leftarrow V_0$ 
4:   for  $t = 0$  at  $= T$  cada  $\Delta t$  do
5:     Calcular  $\alpha_n, \beta_n, \alpha_m, \beta_m, \alpha_h, \beta_h$  utilizando  $V(t)$ .
6:     Calcular las tres  $\tau_x$  y  $x^\infty$  apartir de las anteriores.
7:     Calcular las probabilidades de las compuertas  $n, m, h$ , utilizando las ecuacio-
   nes en diferencias en su forma matricial  $\Pi(t + \Delta t) = A_\pi \Pi(t) + B_\pi$ .
8:     Calcular  $G_{Na} = g_{Na} m^3 h$  y  $G_K = g_K n^4$ .
9:     Almacenar los resultados de este paso en los arreglos
    $T[], V[], n[], m[], h[], G_{Na}[], G_K[], \tau_n[], \tau_m[], \tau_h[]$ 
10:     $I_{ext} \leftarrow I_{ext}(t)$ 
11:    Calcular  $V_m(t + \Delta t)$ 
12:   end for
13:   Devolver los arreglos  $T[], V[], n[], m[], h[], G_{Na}[], G_K[], \tau_n[], \tau_m[], \tau_h[]$  con los re-
   sultados para los tiempos  $[0, T]$ 
14: end function
```

---

Comenzamos con un valor inicial y a partir de ahí empleamos las ecuaciones, para calcular las tangentes, aproximamos a la curva con su tangente.

Para la función INTEGRADISPARO necesitamos cuatro valores de entrada que van a provocar diferentes comportamientos:

1.  $T$  nos indica durante cuánto tiempo queremos correr la simulación.
2.  $\Delta T$  nos indica que tan finos queremos que sean los pasos recordemos que vamos a aproximar la función con segmentos de recta siguiendo la tangente. Si hacemos pasos demasiado pequeños nos vamos a tardar demasiado en hacer el cómputo.
3.  $V_0$  es el voltaje inicial en donde empieza nuestra simulación donde estaba en nuestra neurona cuando empezamos a trabajar.

<sup>5</sup>Hodgkin-Huxley Simulation Using Euler's Method lo puedes encontrar en la siguiente liga <https://webpages.uidaho.edu/rwells/techdocs/Biological%20Signal%20Processing/Chapter%2003%20The%20Hodgkin-Huxley%20Model.pdf>

4.  $I_{ext}(t)$  es la corriente externa, de qué magnitud fue el toque que le estamos dando en este momento al axón

A partir de los elementos iniciales proporcionados ya podemos calcular lo demás. Vamos a querer guardar lo que está ocurriendo para todos los tiempos desde cero hasta  $t$  en cada delta, y lo vamos a hacer en forma de arreglos donde en la posición [0], está la primera medición en  $t=0$  y la posición  $t$  está la medición en el tiempo  $t$ . Vamos a tener toda una serie de puntos donde estamos guardando estos pasos para inicializarlo.

Sabemos que necesitamos es un primer valor a partir del cual vamos a calcular la tangente y vamos a ir aproximando lo demás, entonces para eso queríamos *el voltaje inicial*, como nosotros sabemos donde estaba en reposo nuestra célula originalmente, vamos a poder guardar ese voltaje como el primer valor para nuestra simulación. Ahora todos los demás elementos como  $\alpha_s$  y  $\beta_s$  y etc. se pueden calcular si ya conocímos ese voltaje inicial, entonces a partir de este momento podemos repetir el mismo ciclo tantas veces como sea necesario para cubrir, el intervalo. Desde el tiempo inicial hasta el tiempo  $t$ , brincando de delta  $T$  en delta  $T$ . Entonces dado un voltaje vamos a calcular las diferentes  $\alpha_s$  que son las que se median experimentalmente originalmente utilizando las ecuaciones de las curvas paramétricas ajustadas ?? (paso 5).

Ya calculadas las alfas y betas ahora si se puede calcular las  $\tau$ ,  $n^\infty$ ,  $m^\infty$ ,  $h^\infty$  (paso 6). Teniendo estas entonces podemos calcular las probabilidades para las compuertas  $n$ ,  $m$ ,  $h$  usando las ecuaciones en diferencias en forma matricial, estas matrices se pueden expresar como 2.13.

$$\begin{bmatrix} n(t + \Delta t) \\ m(t + \Delta t) \\ h(t + \Delta t) \end{bmatrix} = \begin{bmatrix} (1 - \Delta t/\tau_n) & 0 & 0 \\ 0 & (1 - \Delta t/\tau_m) & 0 \\ 0 & 0 & (1 - \Delta t/\tau_h) \end{bmatrix} \begin{bmatrix} n(t) \\ m(t) \\ h(t) \end{bmatrix} + \begin{bmatrix} (\Delta t/\tau_n)n_\infty \\ (\Delta t/\tau_m)m_\infty \\ (\Delta t/\tau_h)h_\infty \end{bmatrix}.$$

**Figura 2.13** Forma matricial para el cálculo de las probabilidades de las compuertas  $n$ ,  $m$ ,  $h$ .

Ahora esta parte fue importante sobre todo por el asunto de las pérdidas numéricas, recordemos que la computadora tiene una representación en punto flotante, lo cual quiere decir que un número real no se puede representar en la computadora, esto es importante por el problema del truncamiento, cada vez que se truncan dígitos al momento de calcular estamos perdiendo precisión, por esto es importante la forma en la que se procede a hacer los cálculos porque puede que se llegue a resultados un tanto distintos a los deseados aunque los cálculos sean correctos.

Ya teniendo estos datos podemos fácilmente calcular las conductancias, simplemente sustituyendo los valores (paso 8).

Es bueno una vez que ya terminamos de calcular estos términos que se van a necesitar

## 2. Modelo de Hodgkin-Huxley

en la ecuación principal 2.2 que es la del voltaje de la membrana, ir almacenando en los *resultados temporales* dentro de nuestros arreglos en la casilla que les corresponda, para ese momento (paso 8). En el paso 9, es donde ya vamos a utilizar *la corriente externa* para meterla en la ecuación diferencial para el voltaje. Una vez que tengamos esto tenemos que repetirlo para cada paso, se está calculando el tiempo t, entonces conforme avance el tiempo, va a dar el siguiente valor del voltaje, ya teniendo el siguiente valor del voltaje se va a poder usar para el siguiente paso ("como valor inicial nuevo") y así nos vamos a seguir todo el tiempo. Terminando el tiempo asignado a la simulación, tenemos ya los arreglos llenos de los datos que nos van a poder permitir graficar, qué fue lo que sucedió en cada tiempo, con cada canal y es precisamente así que se puede graficar las mediciones presentadas en la sección anterior.

## Información codificada en las dendritas

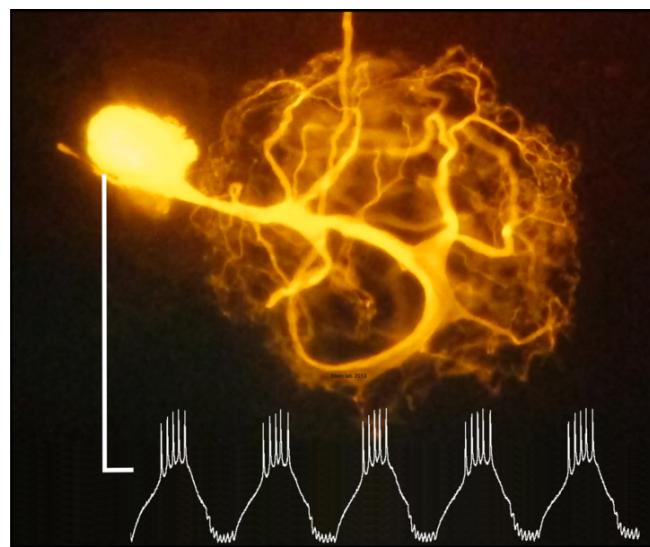
En esta última sección de este capítulo, se va a tratar el cómputo en las neuronas, que vamos a simplificar para poder modelar las neuronas artificiales. Ya vimos detalladamente que pasa a lo largo del axón de la neurona, se ha mencionado que son quienes conforman la región que recibe la información (disparos / pulsos). Estos mandados desde los botones de las terminales de las neuronas presinápticas y son las dendritas las encargadas de enviar estas señales hasta el soma de la neurona y hacer que estas señales químicas se transformen o no en impulsos eléctricos. Ahora la siguiente cuestión es ¿Podemos hacer que se dispare más de una vez (con el mismo estímulo)?.

Los disparos que produzca una neurona depende fuertemente de la interacción con sus neuronas vecinas. Recordemos que durante el período refractario lo que sucede, es que cuando se acumula información en las periferias del cuerpo de la neurona a través de las dendritas y de su cuerpo es porque; está recibiendo información de neuronas vecinas y eso va a hacer que dispare.

Algunas neuronas tienen el período de refractario muy pequeño y disparen frecuentemente, por las conexiones que tienen con sus vecinas, son neuronas que mandan frecuentemente información y a veces pueden lanzar una serie de pulsos muy seguidos, en otras ocasiones puede haber interrupciones entre estas series de pulsos (ver imagen 2.14, en caso de los seres humanos podemos referirnos a las neuronas motoras que están en constante recepción y transmisión de información).

Por otro lado, hay otras neuronas que son muchísimo más pasivas y disparen muy rara vez (neuronas que están a niveles más altos, como reconocimiento de olores, colores, imágenes). Y también nos encontramos con casos donde aparentemente una neurona no reacciona ante ningún estímulo, hasta que pasa un estímulo muy específico<sup>6</sup> y está

<sup>6</sup>Neuronas individuales que forman conceptos abstractos, responden por ejemplo, al nombre de un ser humano. Es así como se descubrió la neurona 'Jennifer Aniston', que disparaba cada vez que el retrato de la actriz se mostraba a los sujetos. Estas neuronas, que responden ante la presentación de



**Figura 2.14** Disparos de una neurona, Wstein, 14 September 2013, Wikimedia Commons. Esta foto muestra la neurona de cangrejo que se tiñó mediante la inyección intracelular de un colorante fluorescente. El recuadro muestra una grabación de las oscilaciones rítmicas del potencial de membrana. Las grabaciones fueron realizadas por Christopher Goldsmith en el laboratorio de Wolfgang Stein en la Universidad Estatal de Illinois, [https://upload.wikimedia.org/wikipedia/commons/c/ca/PD\\_neuron\\_staining\\_and\\_recording.png](https://upload.wikimedia.org/wikipedia/commons/c/ca/PD_neuron_staining_and_recording.png), CC BY-SA 3.0.

dispara (neuronas conceptuales o aisladas). La comunidad científica aún no se atreve a afirmar si este tipo de neurona solo dispara ante ese estímulo específico, o simplemente se trata de un evento que cumple con todas las condiciones para que esta neurona reaccione y mande un pulso. Cabe mencionar que estas neuronas se encuentran en las capas más altas de abstracción de procesamiento del cerebro<sup>7</sup>. Este tipo de diferencias entre neuronas te será más fácil recordar si regresas a la última sección del primer capítulo donde se encuentra un diagrama de los procesos de información en los que puede invertir una neurona, este diagrama es 1.5.

En este momento notamos la gran importancia de la codificación de la información en las redes neuronales de un cerebro. Donde la frecuencia de potenciales de acción nos muestran, ciertos patrones que nos indican el nivel de abstracción al que está respondiendo la neurona y el tipo de información que ayuda codificar o decodificar.

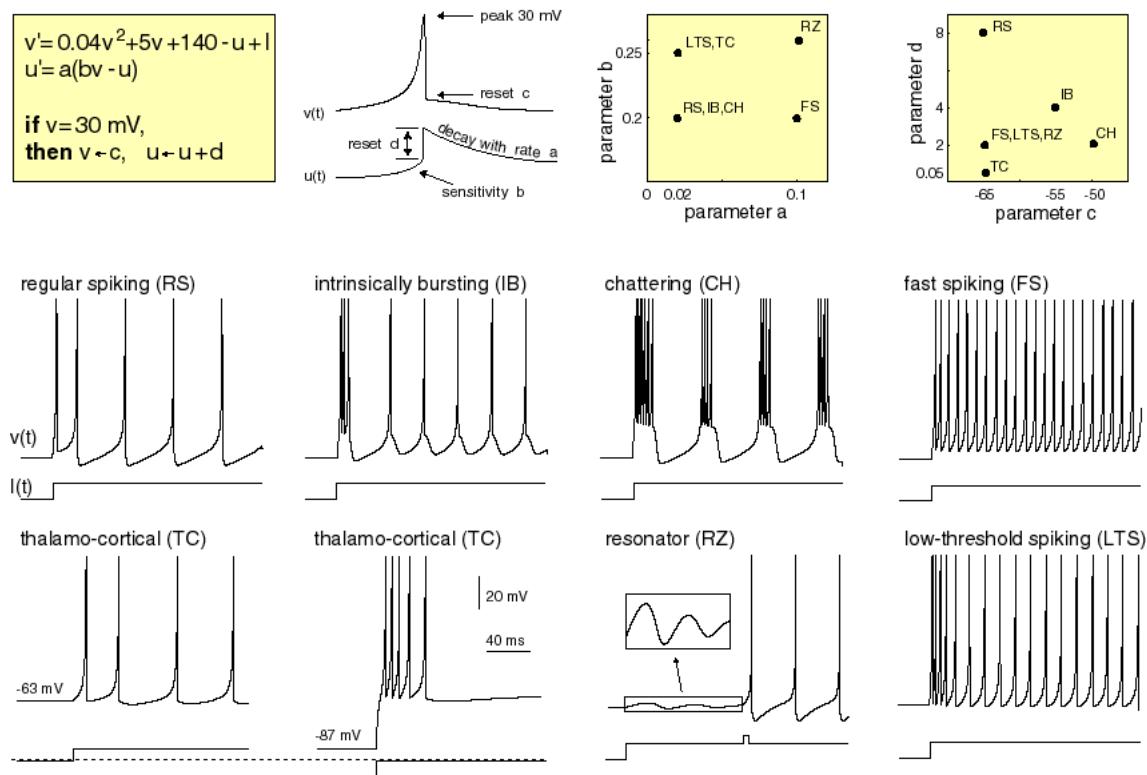
Teniendo ahora noción del papel que juegan los patrones de la frecuencia de disparos, se han hecho numerosas investigaciones acerca de estos. En un intento de entender mejor como es que un cerebro recibe la información desde el exterior, procesa y finalmente provoca ciertas reacciones ante el estímulo. En la siguiente imagen 2.15 que es el resultado

algunas imágenes recibieron la denominación de “células abuelas”.

<sup>7</sup>En la siguiente liga se puede encontrar información más detallada acerca de la investigación con estas neuronas: <https://www.nature.com/articles/s41598-020-64466-7>

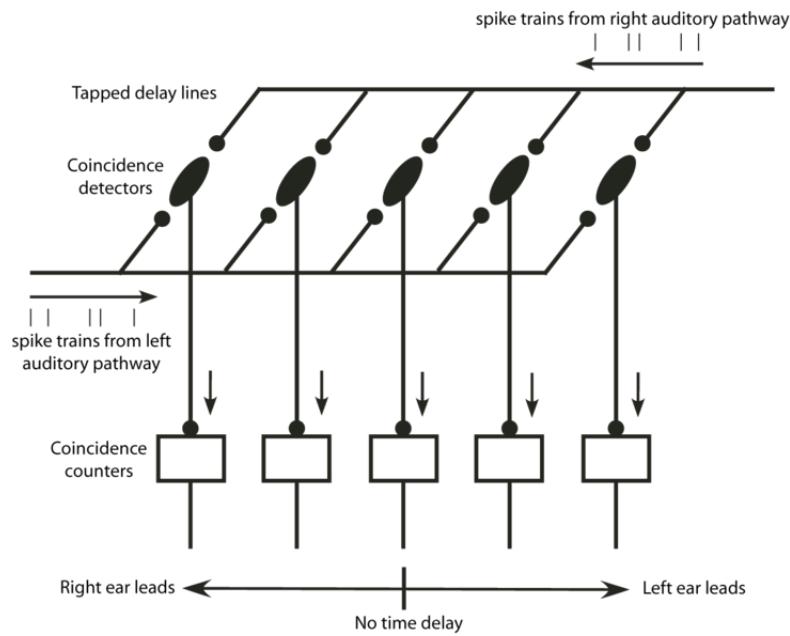
## 2. Modelo de Hodgkin-Huxley

de un artículo donde se utilizaron diferentes métodos para simular disparos de neuronas y notamos la clasificación de estos patrones que se dan ante ciertos estímulos. Cada uno de estos patrones, están codificando cosas distintas, donde una misma neurona podría alternar entre diferentes patrones dependiendo de cuál es el estímulo que está recibiendo. Si bien ha habido algunas propuestas donde se tratan de hacer neuronas artificiales que tomen en cuenta la frecuencia de disparos aún hay un amplio campo de investigación en este tema.



**Figura 2.15** Diferentes patrones de disparo (simulados), Eugene M. Izhikevich, 2003, The Neurosciences Institute, <http://www.izhikevich.org/publications/spikes.html>

Por último mostremos la situación del sistema auditivo con la siguiente imagen 2.16:



**Figura 2.16** Localización del sonido por disparidad, Ian Stevenson, 7 enero 2008, El sistema auditivo registra y analiza pequeñas diferencias en el tiempo de llegada de los sonidos a los dos oídos para estimar la dirección desde la cual el sonido es emitido, <http://www.scholarpedia.org/w/images/4/4d/JeffressFig1.png>

La imagen muestra cómo computan las neuronas, este ejemplo se trata del sistema de audio humano. Cuando nosotros oímos podemos tratar de inferir aproximadamente desde donde está siendo emitido el sonido y eso se logra gracias a que tenemos neuronas conectadas al área auditiva (tanto con el oído derecho como con el oído izquierdo). La señal auditiva va a tardar más en llegar a un oído que al otro, dependiendo de su posición en el espacio, por lo que el cerebro siempre tratar de calcular ese desfase. En el esquema se representa como dentro del espacio de una persona se encuentra una fuente de sonido. Dependiendo de la dirección que provenga (izquierda o derecha), las neuronas receptoras de cada oído harán un cálculo rápido de la distancia a la que se percibió la señal, estás llegaran al cerebro y en algún punto enlazaran con neuronas donde estás cuentan las diferencias de distancias percibidas hasta llegar a dar con el resultado, dándonos a saber cuáles receptores (izquierda o derecha) están más cerca de la señal, y así ubiquemos desde donde proviene el sonido.

# 3 | Aprendizaje de máquina

## Introducción

En este capítulo se desarrolla el procesos que pasa una maquina para que "aprenda", para esto notemos el concepto de aprender. En lo seres humanos se denota como el hecho de adquirir el conocimiento de algo mediante el estudio o la experiencia a partir de ejemplos específicos en nuestro medio, entonces aquellos problemas que inicialmente no pueden resolverse, puedan ser resueltos después de obtener más información acerca del problema. Desde pequeños empezamos por aprender por palabras (conceptos) que asociamos con algo específico, para después relacionar que varios objetos pertenecen a un tipo de conjunto y otros no. Como que un muñeco, una pelota y unos bloques de construcción de plastico, pertenecen a un conjunto denotado como "juguetes", y que plato, taza, y tazón no pertenecen a este conjunto sino al conjunto "vajilla". Entonces para organizar los conceptos que vamos aprendiendo hacemos uso de una función booleana, donde la entrada es el concepto, la pregunta es, ¿Este objeto pertenece a un cierto conjunto de objetos con características similares? y la salida es falso o verdadero. A este proceso, se le conoce como *aprendizaje de conceptos* y en este curso lo simplificamos como *función booleana de aproximación mediante ejemplos*.

Ahora lo que denotamos como el hecho que una maquina aprenda, lo vemos como cualquier programa que mejore su desempeño en alguna tarea mediante la experiencia. Con más formalidad se denota como:

### Definición 3.1

**Aprendizaje maquina:** Se dice que un programa de computadora aprende, si su desempeño en T, medido por P , mejora con la experiencia E. Tal que:

- *T* es un tipo de tarea, como:
  - ★ Jugar un juego de mesa.
  - ★ Clasificar varios tipo de hojas.

- ★ Reconocer una voz en particular.
- *P* es una medida de desempeño, como:
  - ★ Porcentaje de juegos ganados en las partidas.
  - ★ Porcentaje de hojas correctamente clasificadas.
  - ★ Porcentaje de reconocimiento de timbre de la voz.
- *E* la experiencia con ejemplos (entrenamiento), como:
  - ★ Jugar juegos de práctica.
  - ★ Una secuencia de imágenes etiquetadas.
  - ★ Una secuencia de audios etiquetados.

, <sup>a</sup>.

<sup>a</sup>Machine Learning, Mitchell 1997, pag. 14.

A partir de ahora nos dedicaremos a definir correctamente tareas que nos interese que un programa aprenda, para entender la forma más abstracta del problema y así proponer los algoritmos que nos ayuden a resolverlo.

Consideremos también que los sistemas de redes neuronales artificiales son un tipo de algoritmo para la representación del proceso de aprendizaje. Un problema de aprendizaje bien definido requiere una tarea bien especificada, medidas de desempeño y datos para obtener experiencia.

El aprendizaje maquina se apoya de disciplinas, como la inteligencia artificial, probabilidad, estadística, complejidad computacional, psicología, neurobiología y filosofía.

Para proponer un algoritmos de aprendizaje automático necesitamos, elegir el tipo de experiencia de entrenamiento, definir la función objetivo a aprender y un algoritmo para aprender la función objeto a partir de ejemplos de entrenamiento.

Los algoritmos de aprendizaje maquina han sido utilizados ampliamente por la industria bancaria, por gobiernos y por su puesto por el área de la salud. En la industria bancaria por ejemplo, donde es necesario aprender las reglas generales para determinar la solvencia crediticia, a partir de las bases de datos. Por los gobiernos para el reconocimiento de rostros humanos a partir de imágenes. En el área de salud para a partir de bases de datos de pacientes descubrir automáticamente regularidades implícitas en los resultados de tratamientos.

## Espacio de Hipótesis

El aprendizaje automático, es utilizar datos disponibles para, aprender una tarea mediante una función que mejor mapee entradas a ciertas salidas. A esto se le llama aproximación de función, en el que aproximamos una función de destino desconocida (que suponemos que existe) que puede asignar mejor las entradas a las salidas en todas las observaciones posibles del dominio del problema.

Una función de un modelo que se aproxima a la función objetivo y realiza asignaciones de entradas a salidas se denomina hipótesis.

Ahora estas funciones pueden tener formas muy generales en el aprendizaje de máquina pueden tener forma, por ejemplo, de estructuras de datos, como los árboles de decisión, donde cada nodo pregunta si o no, pertenece una clasificación. pueden ser también funciones matemáticas como el caso de las redes neuronales, entonces la forma que tomen estas hipótesis en general puede abarcar muchos métodos y estructuras.

Entonces el aprendizaje consiste en, explorar un espacio de posibles hipótesis para encontrar la hipótesis (una función) que mejor encaje, de acuerdo a lo se obtuvo en los ejemplos de entrenamiento, y predecir alguna característica de salida deseada. Usualmente se denotan como sigue:

- $h$  (hipótesis): una sola hipótesis, por ejemplo una instancia o modelo candidato específico que asigna entradas a salidas, se puede evaluar y se usa para hacer predicciones.
- $H$  (conjunto de hipótesis ): Un espacio de posibles hipótesis para mapear entradas.

Una breve ejemplo para denotar un espacio de hipótesis sería el problema es saber los días que nos conviene ir al cine, donde nuestra tarea  $T$  es aprender a predecir el conjunto de días que nos conviene ir al cine, basado en los atributos de los días, donde cada hipótesis la representaremos apartir de un conjunto de atributos de las instancias (días), entonces cada hipótesis es un vector con tres atributos, *tiene2x1, esEstrenoDePelícula, actoresConocidos*. Para cada atributo de la hipótesis tendría uno de los siguientes valores; Si, No, ?. Donde ? indica que cualquier valor es valido para ese atributo.

Cuando alguna instancia  $x$  cumpla con todos los atributos de una  $h$ , entonces  $h(x) = 1$  y  $x$  es un ejemplo positivo. Entonces para representar la hipótesis, que nos conviene ir solo los días con 2x1, y que hay películas donde los actores son conocidos, la escribimos como  $h(<\text{Si}, ?, \text{Si}>) = 1$ , la hipótesis que cualquier día nos conviene ir al cine la denotamos como  $h(<?, ?, ?>) = 1$ , nuestra función objetivo la denotamos como una función booleana  $c : X \rightarrow \{0, 1\} | X$ , el conjunto de los 365 días del año, entonces  $c(x) = 1$  cuando en los datos nos dicen que con la instancia  $x$  conviene ir al cine,  $c(x) = 0$  en caso que no. Por tanto para aprender la tarea  $T$ , necesitamos *una hipótesis  $h$  en  $H$  tal que  $h(x) = c(x)$*

*para todas las  $x$  en  $X$* . La tarea de aprendizaje del concepto  $c$  requiere aprender el conjunto de instancias que lo satisface, describiendo este conjunto mediante una conjunción de restricciones sobre los atributos de la instancia.

Estas hipótesis (funciones) pueden llegar a ser sumamente complejas y tener que mapear datos de entrada con muchas formas ej. imágenes, trayectorias, etc. En el caso de las redes neuronales, el espacio de hipótesis está determinado por la arquitectura de la red. Vamos a definir el espacio de hipótesis, cuando decidimos qué neuronas vamos a poner en nuestro sistema, como las conectamos entre sí y cómo van a transferirse información de una a la otra y cuántas neuronas van a ser. Lo que veremos a lo largo del curso son diferentes arquitecturas y el impacto que tiene hacer diferentes modificaciones así como las matemáticas que existen detrás de estas.

## Clasificación de los conjuntos de datos

La experiencia  $E$  para aprender la vamos a obtener mediante un conjunto datos, llamados datos de entrenamiento, estos se separan en tres bloques:

- **Entrenamiento:** Datos con los cuales se ajustan los parámetros de la hipótesis (del 50 % al 80 % de los datos). En este bloque se escoje que función del espacio fue mejor para el aprendizaje.
- **Validación:** Datos utilizados para ajustar los parámetros (hiperparametros) del algoritmo de entrenamiento, que puedan afectar qué hipótesis es seleccionada (del 25 % al 10 % de los datos y no deben ser usado durante el entrenamiento). Un ejemplo de un hiperparámetro para redes neuronales son el número de nodos ocultos en cada capa.
- **Prueba:** Datos utilizados para evaluar la posibilidad de que la hipótesis aprendida generalice<sup>1</sup> a datos no vistos anteriormente. Esta porción que se mantiene aparte. Con estos se evalúa el modelo, se reporta la eficacia del modelo según los resultados en este conjunto (del 25 % al 10 % de los datos).

*El conjunto de datos de entrenamiento se usa para aprender una hipótesis y el conjunto de datos de prueba para evaluarla.*

---

<sup>1</sup>Se desea que nuestro modelo de aprendizaje, una vez entrenado con datos que ya hemos visto, se pueda usar con datos nuevos. Para ello debemos asegurarnos que el modelo no ha simplemente memorizado las muestras de entrenamiento, sino que ha aprendido propiedades del conjunto.

## Tipos de aprendizaje

**Aprendizaje Supervisado**, el modelo usa datos etiquetados a una respuesta específica(labeled data), durante el entrenamiento se intenta encontrar una función que aprenda a asignar los datos de entrada (input data) con los datos en el etiquetado. Para después predecir una relación, dado un dato totalmente nuevo para el modelo. Los modelos pueden ser:

- Regresión: Un modelo de regresión busca predecir valores de salida continuos. Por ejemplo, en predicciones meteorológicas, de expectativa de vida, de crecimiento de población.
- Clasificación: En un problema de clasificación se desea predecir una salida discreta. Por ejemplo, identificación de dígitos, diagnósticos.

**Aprendizaje no supervisado**, es usado cuando no se tienen datos “etiquetados” para el entrenamiento. Solo sabemos los datos de entrada. Por tanto, únicamente podemos describir la estructura de los datos, para intentar encontrar algún tipo de organización que simplifique un análisis. Por ello, no se tienen valores correctos o incorrectos (es utilizado para aprender de una manera autoorganizada).

**Aprendizaje por refuerzo**, inspirado en la psicología conductista; donde el modelo aprende por sí solo el comportamiento a seguir basándose en *recompensas y penalizaciones*. Este tipo de aprendizaje se basa en mejorar la respuesta del modelo usando un proceso de retroalimentación (*feedback*). Su información de entrada es el feedback que obtiene del mundo exterior como respuesta a sus acciones. Aprende a base de ensayo-error.

Mientras que el aprendizaje supervisado y el no supervisado aprenden a partir de datos obtenidos en el pasado, el aprendizaje por refuerzo aprende desde cero, es decir, con un estado inicial y su ambiente, va aprendiendo a futuro, mediante posibles penalizaciones o recompensas. El *aprendizaje por refuerzo* es usado en videojuegos porque cada vez que se realizan las acciones correctas se ganan puntos y entonces se entrena a la gente para que pueda conseguir la mayor cantidad de puntos. En este siempre hay: un agente, un ambiente definido por estados, acciones que el agente lleva a cabo (que le llevan de un estado a otro), y recompensas o penalizaciones que el agente obtiene.

En cada acción, el agente solo conoce el estado en el cual se encuentra y las acciones posibles que puede elegir a partir de ese estado. No sabe si llegando al siguiente estado, obtendrá mejores o peores recompensas, irá aprendiendo en cada estado qué acciones lo llevarán a obtener una mayor recompensa a largo plazo, y qué el valor de las acciones en ese estado puedan subir. *Se enfoca en que el agente aprenda una política óptima para alcanzar el objetivo*. El agente siempre está en fases de *exploración* y *explotación*, en la fase de exploración el agente toma una acción de manera aleatoria, y en la de explotación va a tomar acciones basándose en cuán valiosa es realizar una acción a partir de un estado dado.

En plataforma de ventas en línea es donde podemos encontrar este tipo de modelo que están entrenados con este tipo de aprendizaje, donde al iniciar la sesión no conoce nada del usuario, solamente tiene un ambiente dado por los productos de la plataforma y su estado inicial es cero, para hacer individual la experiencia del usuario y que compre más. El algoritmo realiza la acción de mostrar ciertos productos (algún estado) si el usuario da clic a estos productos, el agente recibirá un punto de recompensa, por lo cual pasará a otro estado donde ofrecerá productos del mismo estilo donde pueda maximizar una venta, así se irá adaptando a cada usuario.

## **Parte II**

# **Redes dirigidas acíclicas**

# 4 | Perceptrón simple

## Perceptrón

El perceptrón fue la primera red neuronal artificial (o ANS, Artificial Neural Systems) descrita algorítmicamente. En las décadas de los 60's y 70's, lo popularizó el psicólogo Franck Rosenblatt, en su libro llamado Principios de neurodinámica, donde presentó varios modelos de perceptrones, en el Laboratorio Aeronáutico de Cornell en Estados Unidos, originalmente estaba diseñado para ser una máquina, en vez de un algoritmo. Estaba diseñado específicamente para el reconocimiento de imágenes donde, cada peso era un cable físico por pixel de entrada, este era una matriz de 200 x 200, conectados aleatoriamente a las "neuronas", las actualizaciones de los pesos se realizaron mediante motores eléctricos.

El perceptrón es en sí, es la representación de una sola neurona, este se ocupa para la clasificación de patrones en un conjunto de datos multivariados, con ese se obtienen fronteras lineales en el plano, mediante un algoritmo de aprendizaje que veremos más adelante.

Recordando, una neurona es una célula elemental que a partir de un vector de entrada procedente del exterior o de otras neuronas (estímulo), proporciona una única respuesta (si activo el potencial de acción o no), ver figura 4.1. Los elementos que actúan en una neurona los podemos listar como:

- **Entradas:**  $x_j(t)$ . Las variables de entrada y salida pueden ser binarias (digitales) o continuas (analógicas) dependiendo del modelo de aplicación.
- **Pesos sinápticos:**  $w_{ij}$ . Representan la intensidad de interacción entre cada neurona presináptica  $j$  y la neurona postsináptica  $i$ .
- **Regla de propagación:**  $h_i(t) = \sigma(w_{ij}x_j(t))$ . Proporciona el valor del potencial postsináptico, de la neurona  $i$  en función de sus pesos y entradas.
  - \*  $h_i(t) = \sum_{j=0}^n w_{ij}x_j$ , Es una suma ponderada de las entradas con los pesos sinápticos. Así, si la entrada es positiva, dependiendo de los pesos podemos

#### 4. Perceptrón simple

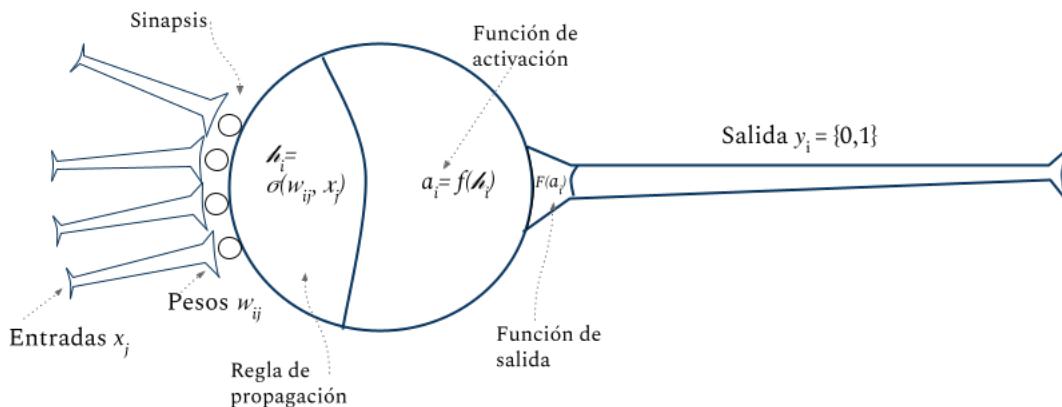
saber si fue una sinapsis excitadora (pesos positivos) o inhibidora (pesos negativos).

- **Función de activación o de transferencia:**  $a_i(t)$  Proporciona el estado de activación actual, de la neurona  $i$  en función de su estado anterior,  $a_i(t - 1)$  y de su potencial postsináptico actual.

- $\star a_i(t) = f_i(a_i(t - 1), h_i(t))$ , es la que usualmente se usa.
- $\star a_i(t) = f_i(h_i(t))$ , en algunos modelos solo se considera que el estado actual no depende del tiempo anterior.

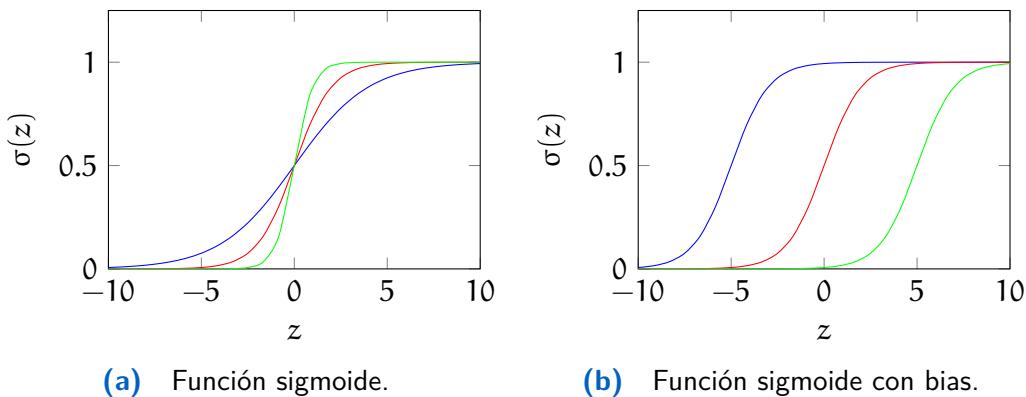
- **Función de salida:**  $F_i(a_i(t))$  Da la salida actual,  $y_i(t)$ , de la neurona  $i$  en función de su estado de activación actual. El estado de activación de la neurona se considera como la propia salida.

- $\star y_i(t) = F_i(a_i(t))$
- $\star y_i(t) = F_i(f_i(a_i(t - 1), \sigma(w_{ij}, x_j(t))))$



**Figura 4.1** Neurona vista como un modelo artificial (perceptrón).

Un perceptrón toma un vector de entradas de números reales, calcula una combinación lineal de estas entradas, luego emite un 1 si el resultado es mayor que algún umbral y -1 de lo contrario. Es decir, dadas las entradas  $x_1 \dots x_n$ , la salida  $o(x_1, \dots, x_n)$  calculada por el perceptrón es 1 si  $w_0 + w_1x_1 + w_2x_2 + \dots + w_nx_n > 0$  y -1 de lo contrario, donde cada  $w$  es una constante  $\mathbb{R}$ , un peso, que determina la contribución de la entrada  $x$  a la salida del perceptrón. La constante  $w_0$  es un *umbral* (bias) que la suma de las entradas con los pesos debe superar para que el perceptrón emita un 1. En otras palabras es un peso que va a actuar junto con una entrada de valor 1, que vamos a poder ajustar para que nuestra función de activación se mueva de derecha a izquierda en el plano para ayudarnos a ajustar nuestros resultados, provocando un gran impacto en el aprendizaje.



**Figura 4.2** Comportamiento de la función de activación (sigmoide) de un perceptrón con una sola entrada, (b) el perceptrón sin el uso de bias, (b) con el uso del bias, donde apesar de estar representandos con la misma entrada , el uso del bias afecta en los resultados de salida. Así si quisieramos que este perceptrón nos diera  $y = 0$  con una entrada  $x = 2$  sin el uso, ni ajuste del bias sería imposible, pues en (a) apesar que la gráfica azul está la entrada está ajustada con el peso  $w = 0.5$ , la roja con el  $w = 1$ , y la verde con el  $w = 2$  solo lo logramos alargarla un poco, haciendo que entradas que antes eran correctas ahora caigan 0 también. Entonces lo que necesitamos es más bien "mover" la gráfica, esto lo logramos con la gráfica (b) donde la entrada (única) está sumada con un bias (umbral)  $x_0 = 1$ , ajustado en azul con peso  $w_0 = 5$ , en rojo con  $w_0 = 0$ , en verde con  $w_0 = -5$  y el peso  $w_1 = 1$ . Donde con  $w_0 = -5$  logramos nuestro objetivo de tener una salida  $y = 0$  con  $x = 2$ . El bias nos permite mover la función fuera del origen.

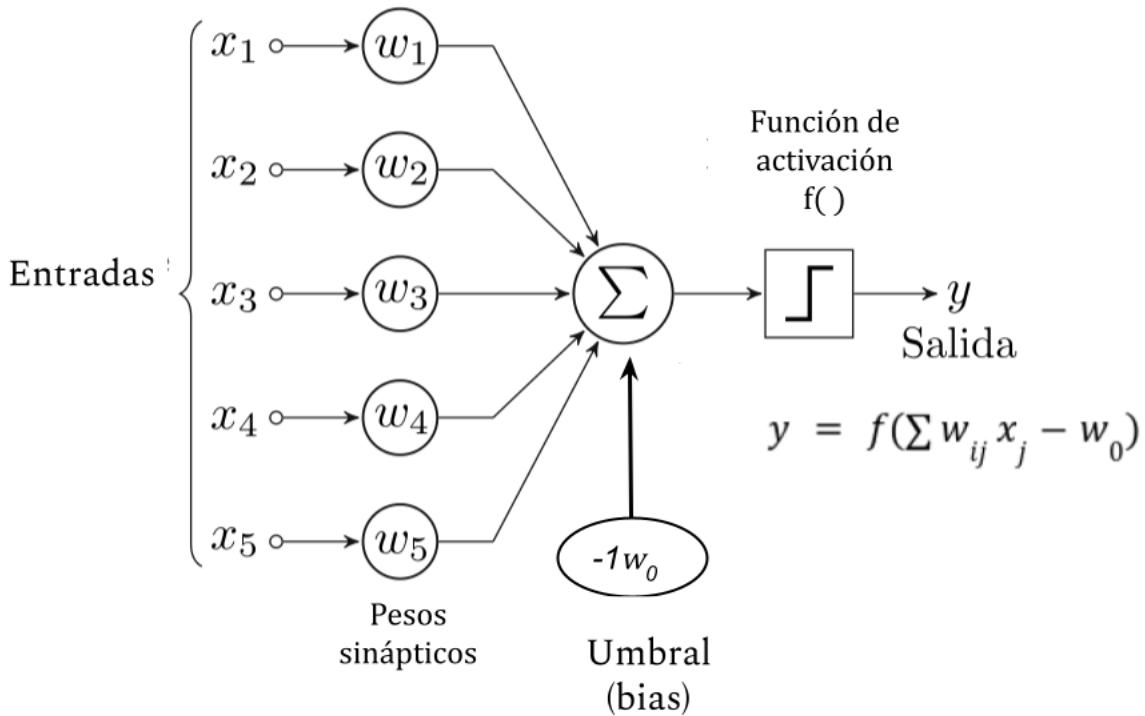
Esto se muestra en las siguientes graficas 4.2b. Más adelante hablaremos de su regla de entrenamiento (Training rule).

El hecho de que un perceptrón aprenda implica elegir valores para los pesos denotados también por  $\theta$ . Ahora, el espacio  $H$  de las hipótesis candidatas consideradas en el aprendizaje del perceptrón, es el conjunto de todos los posibles vectores de pesos.

$$H = \{w | w \in \mathbb{R}^{n+1}\}$$

Si bien en el momento que se publicó los logros con el modelo del percetrón las expectativas eran bastante altas, a medida de los años los científicos Marvin Minsky and Seymour Papert desestiman en gran medida los alcances que realmente se puede tener con el perceptrón, al mostrar<sup>1</sup> que no puede predecir operaciones lógicas que no sean linealmente separables, tal es el caso de la función XOR, que no es separable linealmente, siendo imposible que pueda aprender esta función. Esto y el gran costo que representaba procesar todos los elementos que implicaba el entrenamiento, causa que por un buen

<sup>1</sup>En 1969 publican Marvin Minsky y Seymour Papert que perceptrones de una sola capa (simples) solo son capaces de aprender a distinguir patrones linealmente separables, en el libro "Perceptrons".



**Figura 4.3** Modelo estandar de un perceptrón.

tiempo se desetime el uso del perceptrón. Es hasta después de unas decadas que vuelve a tener relevancia con la propuesta de un perceptrón multicapa usando retropagación (Feedforward), siendo estos capaces de resolver la función XOR

## Compuertas lógicas con neuronas

Aquí se muestra como se puede utilizar un perceptrón para simular compuertas lógicas tales como el or, not, and.

Para simular la compuerta *not*, como está es una función booleana de  $B \rightarrow B$ , tal que  $\text{not}(x) = -x$  entonces en un plano de dos dimensiones, la podemos representar con dos puntos, el  $p_1 = (0, 0)$  y  $p_2 = (1, 0)$  donde  $p_1$  representa cuando  $\text{not}(0) = 1$ ,  $p_1$  representa cuando  $\text{not}(0) = 1$ . Teniendo el espacio de la función definido lo que nos toca es, separar el plano para clasificar las entradas, este claramente se puede separar con una linea vertical, o con lineas con pendiente 1 o  $-1$ . Para esta función solo necesitamos de una entrada y un sesgo (bias), donde la entrada la combinaremos con un peso, este peso lo asignaremos a tanteo (por la sencillez de la operación). Así el peso  $w_1 = -1$  y el peso asignado al bias sera  $w_0 = 0.5$ , ahora con esto datos podemos:

- Hacer la función de propagación donde  $h(x) = (x * -1) + (0.5) * 1 = 0.5 - x$

- Hacer la función de activación escalón  $a(x) = \text{sgn}(h) = \text{sgn}(0.5 - x)$
- Dar la salida donde  $s(1) = \text{sgn}(0.5 - 1) = 0$  y  $s(0) = \text{sgn}(0.5 - 0) = 1$ , en este caso la salida es la identidad de la activación.

x	h	s
0	0.5	1
1	-1.5	0

Algo similar va a pasar con la compuerta *and* y *or* donde al necesitar de dos entradas para la compuerta, asignaremos dos entradas para el perceptrón igualmente y las representaremos en el plano con cuatro puntos, donde cada punto representa una instancia y se le asigna valor positivo o negativo en el plano, así pues para el *and* tenemos los puntos  $p_1 = (0, 0)$ ,  $p_2 = (0, 1)$ ,  $p_3 = (1, 0)$ , negativos y  $p_4 = (1, 1)$  el único positivo. Así nos damos cuenta que necesitamos una recta con pendiente negativa y fuera del origen, que nos separe estas clases de puntos. Por tanto para el bias le asignamos un peso de  $w_0 = -1.5$ ,  $w_1 = 1$  y  $w_2 = 1$ , con estos datos podemos:

- Hacer la función de propagación donde  $h((x_1, x_2)) = (x_1 * 1) + (x_2 * 1) + (-1.5) * 1 = x_1 + x_2 - 1.5$
- Hacer la función de activación escalón  $a((x_1, x_2)) = \text{sgn}(h) = \text{sgn}(x_1 + x_2 - 1.5)$
- Dar la salida donde  $s = a(x)$ , es la identidad de la activación.

$x_1$	$x_2$	h	s
0	0	-1.5	0
0	1	-0.5	0
1	0	-0.5	0
1	1	0.5	1

Para la compuerta *or* es algo muy similar pues podemos igualmente representar la función con cuatro puntos en el espacio cada uno representando una instancia, solo que ahora tres de estos puntos serán positivos y solo uno negativo, los puentes positivos serían  $p_2 = (0, 1)$ ,  $p_3 = (1, 0)$  y  $p_4 = (1, 1)$ , mientras que  $p_1 = (0, 0)$  negativo, el plano lo podemos dividir con una linea recta con pendiente negativa, así le asignamos  $w_0 = -0.5$ ,  $w_1 = 1$  y  $w_2 = 1$ , con esto hacemos los paso que ya sabemos:

- Hacer la función de propagación donde  $h((x_1, x_2)) = (x_1 * 1) + (x_2 * 1) + (-0.5) * 1 = x_1 + x_2 + 0.5$
- Hacer la función de activación escalón  $a((x_1, x_2)) = \text{sgn}(h) = \text{sgn}(x_1 + x_2 + -0.5)$

#### 4. Perceptrón simple

---

- Dar la salida donde  $s = a(x)$ , es la identidad de la activación.

$x_1$	$x_2$	$h$	$s$
0	0	-0.5	0
0	1	0.5	1
1	0	0.5	1
1	1	-1.5	1

Tomando el hecho que en la naturaleza las neuronas van pasando información en una estructura que forma niveles de abstracción, esto lo modelamos como capas de neuronas conectadas entre sí, cada capa haciendo su trabajo de abstracción. El perceptrón simple es un modelo neuronal unidireccional, una capa de entrada y otra de salida, que por si solo no puede separar todas las funciones lógicas pues tenemos el XOR, para resolver esto usaron perceptrones multicapa que se explicará más adelante en el curso.

## Funciones de activación

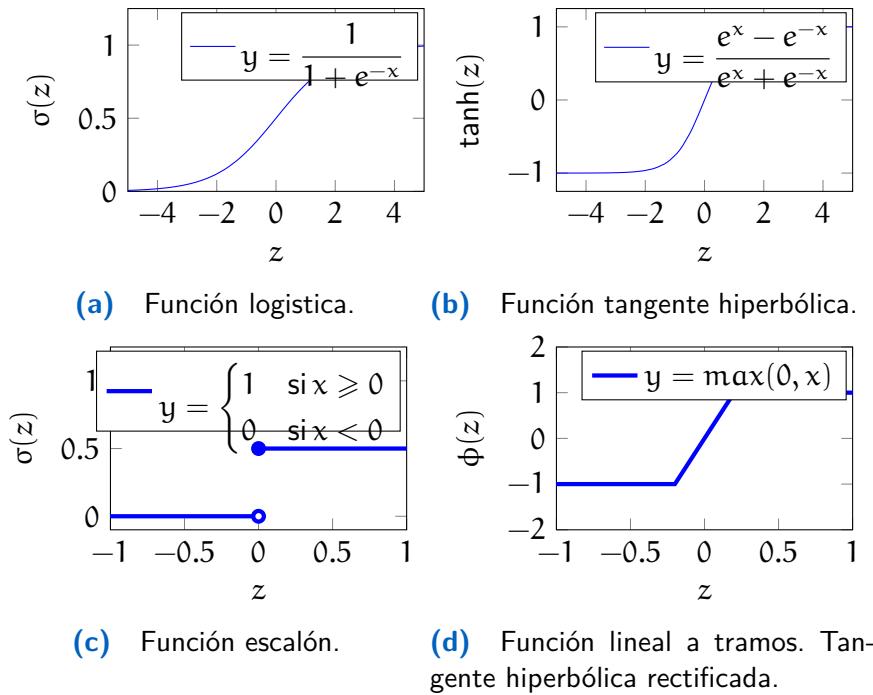
Recordando que la forma de las funciones de activación es  $y = f(x)$ , donde  $x$  representa el potencial postsináptico e  $y$  el estado de activación de la neurona, es decir si va a lanzar un disparo o no. Las funciones de activación más empleadas son:

## Funciones de error

Entonces tomando como base el perceptrón, una vez obtenidas las salidas con una primera iteración nos daremos cuenta de que tan lejos o que tan cerca estuvimos de la respuesta correcta, con esto darnos la oportunidad de que pesos ajustar respecto a sus entradas, en la segunda iteración. Ahora para facilitarnos esto y tomando en cuenta que el entrenamiento consiste en varias iteraciones hasta llegar a aprender la tarea  $T$ , hacemos uso de una función de error que nos ayude a minimizar la diferencia de error en las salidas.

Primero veamos el entrenamiento para una sola neurona, para esto haremos uso de la regla de aprendizaje del perceptrón (learning rule perceptron), donde para cada entrada, en la capa de salida se le calcula la desviación a la función objetivo. El cual utilizamos para ajustar los pesos del perceptrón (ver fig 4.5).

Usualmente al principio del entrenamiento se asignan pesos aleatorios, a medida que avance el entrenamiento, se van modificando con cada iteración, así  $w_i \leftarrow w_i + \Delta w_i$ . Esto con base a [la regla de aprendizaje](#) donde:



**Figura 4.4** Las funciones de activación más usadas son la función sigmoide  $\sigma(z)$  y la tangente hiperbólica  $\tanh(z)$ .

$$\Delta w_i = \alpha(y - y_{out})x_i \quad (4.1)$$

Con  $y$  es la salida deseada,  $y_{out}$  la salida generada,  $\alpha$  la taza de aprendizaje (learning rate) y  $x_i$  la entrada  $i$ . Lo que hace la taza de aprendizaje es moderar el grado en que los pesos son cambiados con cada iteración, se le asigna un valor muy pequeño (0.1 o 0.2) y conforme se logran ajustar los pesos se minimiza aún más.

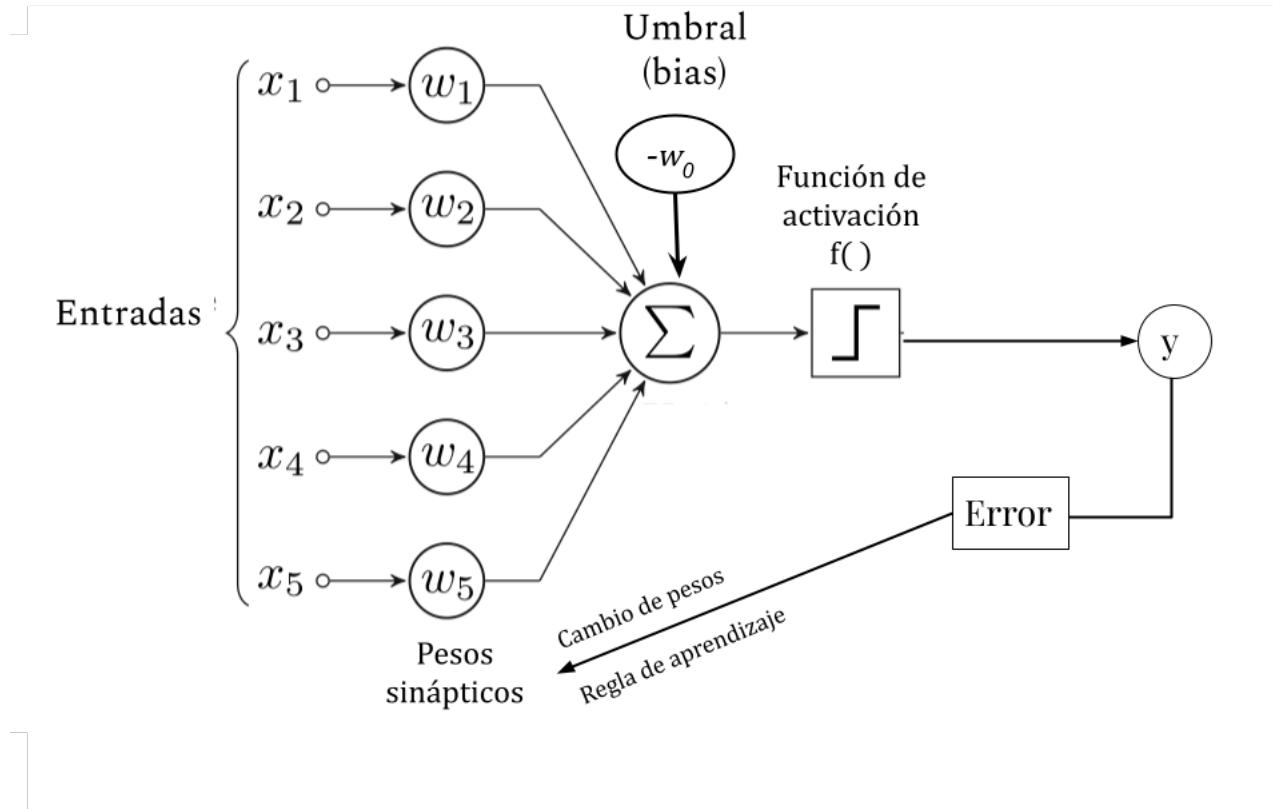
Para entrenar un perceptrón, utilizamos cualquier método de optimización de funciones para encontrar los parámetros  $w$  que minimizan el error con alguna de las siguientes funciones de error:

**Diferencias al cuadrado**, también conocida como regla aprendizaje delta, es la suma de cuadrados de errores que se tuvieron con cada ejemplar el el conjunto de entrenamiento. Los podemos describir como:

$$\frac{1}{2m} \sum_{m=0}^{M-1} (y_m - a_m)^2 \quad (4.2)$$

con  $y_m$  y  $a_m$ , la salida obtenida dado un ejemplar  $m$  y la salida correcta del ejemplar  $m$  respectivamente.

## 4. Perceptrón simple

**Figura 4.5** Modelo estandar de un perceptrón.

**Entropía cruzada.** Se usa para problemas de clasificación ya que su comportamiento es más suave (soft) y permite hacer clasificaciones más certeras. Se define como:

$$L(\Theta) = -\frac{1}{m} \sum_{m=1}^{m=0} (y^{(m)} \log(a_m) + (1 - y_m) \log(1 - a_m)) \quad (4.3)$$

Entonces juntando los conceptos que ya sabemos podemos describir el algoritmo de entrenamiento para el perceptrón de la siguiente forma:

1. Iteramos sobre todos los ejemplares.
2. Para cada ejemplar se calcula la función de propagación, es decir la suma ponderada.
3. Se calcula la función de activación con esta suma.

4. Se calcula la función de salida.
5. Actualización de pesos de acuerdo a la regla de aprendizaje.
6. Repetir de 1-6 hasta que los pesos nos satisfagan, un número de iteraciones establecidas.

## Medidas de rendimiento

Las medidas de rendimiento de una red nos sirven para ver de manera concreta como se comportó nuestra red, durante el entrenamiento, que fue lo que pudo aprender, que sobre aprendio, es decir, memorizo, y que aún le cuesta trabajo aprender, por tanto, será incapaz de predecir. Se usan las siguientes:

**Matriz de confusión** : (Confusion matrix) Es una matriz, donde las celdas representan las predicciones que hizo nuestro modelo de clasificación de clases, respecto a las salidas esperadas. Así siendo las columnas las salidas  $y$ 's del modelo entrenado y las filas las salidas esperadas  $y_{true}$ . Nos facilita a ver cuando un clasificador está confundiendo clases, contabilizando a que clase etiqueto a los diferentes ejemplares. Ahora veamos, que puede estar representando cada celda en la matriz de confusión, estos pueden ser:

1. **VP** Verdaderos positivos (*TP, True Positive*): La clasificación de los ejemplares predichos, condicen con las etiquetas esperadas de los ejemplares.
2. **VN** Verdaderos negativos (*TN, True Negative*): El ejemplar que no es parte de clase, no son asignados a esa clase, son predichos correctamente.
3. **FP** Falso positivo (*FP, False Positivo*): El ejemplar que **no** es parte de una clase i fue clasificado como tal.
4. **FN** False negativo (*FN, False Negative*): El ejemplar que es parte de una clase i no fue clasificado como tal.

**Ejemplo 4.1.** Notemos un ejemplo sencillo, supongamos que tenemos una tarea binaria, donde queremos indicar que una persona está embarazada (de acuerdo a unos estudios). Ahora tenemos que:

- **VP**, sería predecir que una mujer está embarazada y que en efecto esté embarazada. (**Correcto**)
- **VN**, sería con un hombre que no está embarazado y pues en efecto no está embarazado. (**Correcto**)
- **FP**, sería predecir que un hombre está embarazado y **no** este embarazado. (**Error, tipo 1**)

## 4. Perceptrón simple

- ***FN***, sería predecir que una mujer **no** está embarazada y esta embarazada. (**Error, tipo 2**)

Entonces dados los resultados binarios que nos entregó el modelo, los médicos nos dan las respuestas correctas a 10 estudios, representadas en la siguiente tabla.

Ejemplar	1	2	3	4	5	6	7	8	9	10
Sujeto	M	F	M	F	M	M	F	F	F	M
Etiqueta	No	Si	No	Si	No	No	No	Si	No	No
Clase	0	1	0	1	0	0	0	1	0	0
Predicho	0	0	0	1	0	1	0	1	0	0
Valores	VN	FN	VN	VP	VN	FP	VN	VP	VN	VN

Con estos datos, podemos construir nuestra matriz de confusión MC contabilizando las salidas obtenidas respecto a las deseadas. Quedando así de la siguiente manera:

		Salidas y	
		Si	No
Etiquetas	Si	VP = 2	FN = 1
	No	FP = 1	VN = 6

**Tabla 4.1** Matriz de confusión binaria.

Ahora para un modelo que sea multiclase, en el que tengamos que clasificar varias clases de ejemplares. Tendremos una matriz M de  $n * n$  donde n es el número de clases, así los valores **VP**, **VN**, **FP**, **FN**, son calculados para cada clase e identificados en la matriz M de la siguiente forma, también puedes ver la figura 4.6:

- **VP**, para la clase i será la celda  $M[i][j]$  con la  $i = j$ .
- **VN**, para la clase i será la suma de los valores de toda la matriz menos los valores de la fila i ni los valores de la columna i.
- **FN**, para la clase i será la suma de los valores en la fila i, excepto la celda  $M[i][j]$ , con  $i = j$  que es el **VP**.
- **FP**, para la clase i será la suma de los valores en la columna i, excepto la celda  $M[i][i]$  que es el **VP**.

Valores para la clase i		Salidas y del clasificador			
		clase 1	clase 2	clase i	clase n
Etiquetas ytrue	clase 1	VN	VN	FP	VN
	clase 2	VN	VN	FP	VN
	clase i	FN	FN	VP	FN
	clase n	VN	VN	FP	VN

**Figura 4.6** Matriz de confusión para clasificador multiclasses.

**Ejemplo 4.2.** Tenemos un clasificador encargado de identificar las cinco vocales del español, escritas a mano. Entonces tenemos un total de 5 clases representadas por a, e, i, o, u, cada letra representada por un número del 0 al 4. Así la clase 0 = a, la clase 1 = e y así respectivamente. Este fue entrenado con 15 ejemplares etiquetados y se obtuvieron los siguientes resultados:

Ejemplar e	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
Etiqueta	a	e	i	o	u	a	a	e	i	i	o	u	o	u	a
Clase y <sub>true</sub>	0	1	2	3	4	0	0	1	2	2	3	4	3	4	0
Predicho y	0	1	3	3	2	1	1	1	2	2	4	4	3	2	0

Ahora para hacer la matriz de confusión MC notamos que tenemos una matriz de 5x5, donde  $MC[i] = \text{Predicho}$  y  $MC[j] = \text{EtiquetasReales}$ . Entonces necesitamos contabilizar las predicciones y asignarlas a sus respectivas celdas, donde con la tabla anterior notamos que en el ejemplar e3 con  $y_{true}(e3) = 2$ , se predijo que  $y(e3) = 3$ , entonces  $MC[2][3] += 1$ , pues la etiqueta nos posiciona en la fila 2 y lo predicho en la columna 3. Ahora en el ejemplar e4 con  $y_{true}(e4) = 3$  se predijo que  $y(e4) = 3$ , entonces  $MC[2][2] += 1$ . Así con cada ejemplar vamos a ir sumando su clasificación. Para construirla tendríamos un pseudocódigo siguiente:

```
def getMC (Ejemplares, y):
    """
    :param y: salidas obtenidas (int)
    """
    MC = [len(y)] [len(Ejemplares)] # y == Ejemplares
    full_zeros(MC)
    for e in range (0, len(Ejemplares)):
        predicho = y[e]
        y_true = Ejemplares[e].etiqueta
```

## 4. Perceptrón simple

```
MC[predicho][y_true] += 1
return MC
```

Dados los datos anteriores tenemos, la siguiente matriz de confusión:

		Salidas y				
		1	2	3	4	5
Etiquetas	1	2	2	0	0	0
	2	0	2	0	0	0
	3	0	0	2	1	0
	4	0	0	0	2	1
	5	0	0	2	0	1

**Tabla 4.2** MC = Matriz de confusión multiclas.

Para calcular los valores *VP*, *VN*, *FP*, *FN*, por clase se propone el siguiente pseudo-código:

```
MC = getMC(y_salidas, y_true)

VP = VN = FP = FN = [0,0,0,0,0]

for i in range(len(MC)):
    for j in range(len(MC[0])):
        FN[i] = FN[i] + MC[i][j] if (i != j) else FN[i]
        FP[i] = FP[i] + MC[j][i] if (i != j) else FP[i]
        VP[i] = MC[i][j] if (i == j) else VP[i]
        VN[i] = sum(MC) - FN[i] - FP[i] - VP[i]
```

Si quisiéramos saber los valores *VP*, *VN*, *FP*, *FN*, para todo el desempeño total, simplemente sumamos lo obtenido en cada clase así:

```
VP_total = sum(VP)
VN_total = sum(VN)
FP_total = sum(FP)
FN_total = sum(FN)
```

Ahora los valores de cada celda nos representan lo siguiente:

- $VP_{total}$ , toda la diagonal de la matriz. La salida del modelo coincide con lo etiquetado con el ejemplar.
- $VN_{total}$ , todas las veces que no era la clase i y dijó que no era de la clase i.
- $FN_{total}$ , todas las veces que era  $clase_i$  y dijó que era  $clase_x$ . (deseado  $FN = 0$ )
- $FP_{total}$ , todas las veces que predijo que era  $clase_i$  y era  $clase_x$ . (deseado  $FP = 0$ )

Así que, si los valores que no están en la diagonal de la matriz son cero o todas nuestras clasificaciones están en la diagonal, podemos decir que nuestro modelo aprendió a clasificar correctamente todas las clases.

**Precisión y recuperación** : La precisión (*precision*) nos dice la proporción de ejemplares que se logró clasificar correctamente. Mientras que recuperación (*recall*) nos dice cuantas asignaciones de lo que nos interesa pudo clasificar correctamente en otras palabras donde del total de las respuestas correctas que se pueden tener, cuantas respuestas positivas acertadas se tuvo.

$$P = \frac{VP}{VP + FP} = \frac{\text{VerdaderosPositivos}}{\text{ValoresEsperados}} \quad (4.4)$$

$$R = \frac{VP}{VP + FN} = \frac{\text{VerdaderosPositivos}}{\text{ValoresPredichos}} \quad (4.5)$$

Valores para la clase i		Salidas y del clasificador			
		clase 1	clase 2	clase i	clase n
Etiquetas ytrue	clase 1	VN	VN	FP	VN
	clase 2	VN	VN	FP	VN
	clase i	FN	FN	VP	FN
	clase n	VN	VN	FP	VN

**Recall =  $VP / (VP+FN)$**   
Cuantos valores esperados fueron asignados realmente.

**Precisión =  $VP / (VP+FP)$**   
Cuento de lo que clasifico fue correcto

- Cuando el modelo detecta los ejemplares, pero los incluye en otras clases también:  $P$  es bajo y  $R$  es alto.
- Cuando el modelo no detectó bien los ejemplares, pero tampoco los incluyó en otras clases:  $P$  es alto y  $R$  es bajo.

#### 4. Perceptrón simple

- Cuando el modelo detecta los ejemplares y no los incluye en otras clases:  $P$  y  $R$  es alto.
- Cuando el modelo no detecta los ejemplares:  $P$  y  $R$  es bajo.

En ocasiones le daremos más importancia al recall y en otras a la precisión. Retomando el ejemplo previo 4.1, no nos importa tanto los casos de error tipo 1, donde el modelo se equivoca con los ejemplares negativos, nos importa los errores de tipo 2, los Falsos Negativos, en este caso le daremos especial atención al Recall y que este sea alto. Pero si bien nuestra tarea es detectar cuando un correo es spam, no impacta tanto que aunque en ocasiones un correo sea spam, no lo clasifique como tal (error tipo 2 FN), nos es crucial que un correo que no sea spam nos lo clasifique como tal (error tipo 1 FP). Así deseando que los falsos positivos se acerquen a cero, dandonos como resultado una precisión alta aunque el recall sea bajo.

**Exactitud y medida f** : La exactitud (*accuracy*) es una medida de cuántas predicciones correctas en total hizo el modelo para el conjunto de datos completo (no se recomienda usar si tienes clases desbalanceadas, es decir, muchos elementos de una clase y poco de otra pues nos puede fallar totalmente con las clases pequeñas y aun así su valor sería alto). La medida f (*f score*) se utiliza para combinar las medidas de precisión y recall en un solo valor. Es práctico porque hace más fácil el poder comparar el rendimiento combinado de la precisión y la recall. Se dan las ecuaciones a continuación:

$$\text{Accuracy} = A = \frac{VP + VN}{VP + VN + FP + FN} = \frac{VP + VN}{\text{TodosLosValoresClasificados}} \quad (4.6)$$

$$F = \frac{2}{\frac{1}{P} + \frac{1}{R}} = 2 \frac{P * R}{P + R} \quad (4.7)$$

La medida f, también la podemos escribir (con un poco de aritmética) como:

$$F = \frac{2VP}{2VP + FP + FN} \quad (4.8)$$

# 5 | Perceptrón multicapa

## Intro

Perdí el archivo :CCCC

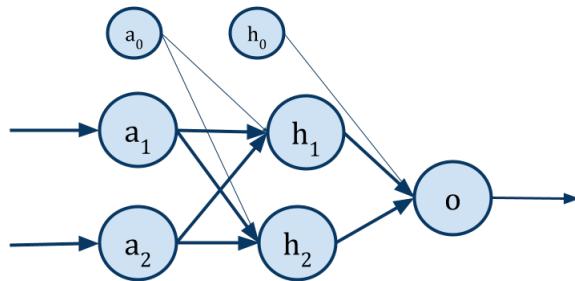
## XOR

Anteriormente, habíamos logrado separar clases de ejemplares siempre y cuando estos se pudieran modelar en un plano linealmente separable. Hasta ahora había sido un reto importante (no logrado) hacer clasificaciones de ejemplares distribuidos en un plano no separable linealmente, como lo es *la función XOR* donde, tenemos que las respuestas positivas se encuentran en una diagonal opuesta a las respuestas negativas y no hay manera de separar a los blancos de los negros con una sola frontera lineal, lo que hace imposible separar con un solo perceptrón. Entonces notamos que necesitamos de dos líneas que nos permitan separar el plano. Así llegamos a la idea que necesitamos más de una capa, que nos permita hacer la siguiente separación del plano.

Entrada $x_1$	Entrada $x_2$	Salida $y$
0	0	0
0	1	1
1	0	1
1	1	0

**Figura 5.1** Tabla de verdad para la función XOR.

## 5. Perceptrón multicapa



**Figura 5.2** Perceptrón para la función XOR, con una capa oculta.

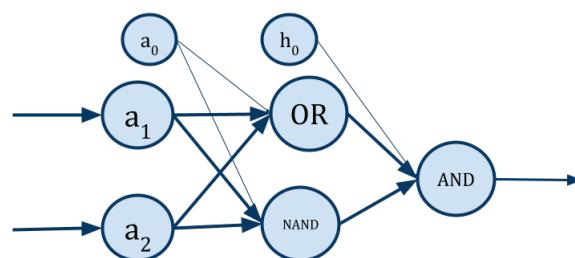
En la tabla de verdad notamos que para la función XOR, cuando la suma de las entradas es:

- par (tomando en cuenta al cero como par) la salida es 0.
- impar la salida es 1.

Para poder aprender la función únicamente tenemos que añadir una capa intermedia, la llamada capa oculta. Esta capa será la responsable de tomar las características de los vectores de entrada.

Ahora para la función tenemos dos entradas, una capa oculta (con dos neuronas) y una salida. Las neuronas en la capa oculta nos permitirán dividir el plano. Así, una solución sería que las neuronas ocultas denotadas por  $h$  (*hidden* de oculto en inglés) comparten pesos, la primera neurona  $h_1$  haría la función OR, haciendo la distinción entre las entradas 00, la segunda neurona  $h_2$  haría la función NAND haciendo posible distinguir el vector de entrada 11, una vez obtenido estos resultados de la capa oculta la neurona de salida o de output se encargara de distinguir la intersección entre estas, ejecutando la función AND. En resumen, la función **XOR** la podemos escribir como:

$$\begin{aligned} \text{XOR} &= (x_1 \vee x_2) \wedge \neg(x_1 \wedge x_2) \\ \text{XOR} &= \text{AND}(\text{OR}(x_1, x_2), \text{NAND}(x_1, x_2)) \end{aligned} \quad (5.1)$$



**Figura 5.3** Solución para perceptrón de la función XOR.

## Propagación hacia adelante manual

En esta parte vamos a ver como podemos evaluar la red para la solución que se dio en la sección anterior. (Recordemos que también se pueden dar otras soluciones para el XOR).

Entonces el procedimiento matemático general para poder evaluar una red cuando tenemos más de un perceptrón. Veamos primero que pasa con el **xor**, y el **nand**, por el momento no tomaremos en cuenta las neuronas de entrada a puesto que solo son usadas para almacenar las entradas. La capa oculta está formada realmente por dos percetrones, que dan su salida a una neurona más que es la capa de salida de nuestra red. Para calcular sus salidas debemos aplicar la suma ponderada de las entradas a una función activación, así podemos ver que la salida de una neurona oculta es la siguiente:

$$z_j = g(\sum_i w_{ij} a_i) \quad (5.2)$$

A estos perceptrones a su vez se les están aplicando la función de activación sigmoide:

$$h_j = \frac{1}{1 + e^{-z}} \quad (5.3)$$

Así la neurona de salida recibe a los perceptrones ya evaluados y listos para aplicar pesos a estos y una función de activación, que en este caso son dados de la siguiente forma:

$$z_o = g(\sum_j w_{jo} h_j) \quad (5.4)$$

$$o = \frac{1}{1 + e^{z_o}} \quad (5.5)$$

Así tomando de ejemplo la función XOR, los valores de la capa oculta son evaluados de la siguiente forma, donde  $x_1$  y  $x_2$  son evaluados por 0 o 1 según sea necesario:

$$\begin{aligned} h_0 &= 1 \\ h_1 &= g(1w_{01} + x_1w_{11} + x_2w_{21}) \\ h_2 &= g(1w_{02} + x_1w_{12} + x_2w_{22}) \end{aligned} \quad (5.6)$$

Entonces hasta aquí ya tenemos los valores de la capa oculta, una vez que ya tenemos este conjunto de valores podemos empezar a trabajar con el tercer perceptor, sus valores de entrada van a estar dados por los valores de activación de  $h_1$  y  $h_2$  y por un sesgo,

## 5. Perceptrón multicapa

pues recordemos que es la función Nand, por tanto necesitamos movernos ligeramente del origen. Lo que vamos a tener aquí la fórmula se ve similar, lo único es que ahora los valores de entrada fueron los valores que obtuvimos en el cálculo de la capa anterior:

$$o = g(h_0w_{0o} + h_1w_{1o} + h_2w_{2o}) \quad (5.7)$$

Como estamos trabajando capa por capa, primero entran los valores con los que van a trabajar todos los perceptrones, después calculamos todos los de la capa oculta que son independientes entre sí, aunque tengan en común las mismas entradas y finalmente hacemos el cálculo de la siguiente capa, que en ese caso es la salida, pero bien podría ser otra oculta.

Por esto estamos hablando de un algoritmo llamado de *alimentación hacia adelante*, este hay que recorrerlo de izquierda hacia derecha, para ir obteniendo los diferentes valores de activación en cada una de las capas. Una característica de este tipo de arquitectura es que, al estar las neuronas conectadas, siempre se está conectando de las neuronas en la capa anterior hacia las neuronas en la capa siguiente. Cuando tenemos la evaluación de las capas ocultas, podemos obtener finalmente la salida final, con una evaluación final.

Si bien esta forma de evaluación nos lleva al resultado correcto, este método se puede simplificar sobretodo en el caso que estemos manejando más entradas y más perceptrones en las capas ocultas. Así nos daremos cuenta que es posible usar matrices, esta forma de evaluación la veremos en la siguiente sección.

## Propagación hacia adelante vectorizada (con matrices)

La sección anterior si bien nos da la idea somera de cómo van a ser las operaciones para el aprendizaje ahora veamos lo de forma matricial. Esto nos ayudará a escribirlo en un momento dado en el lenguaje de nuestra convención.

Retomando la red neuronal de la sección anterior ahora con pesos.

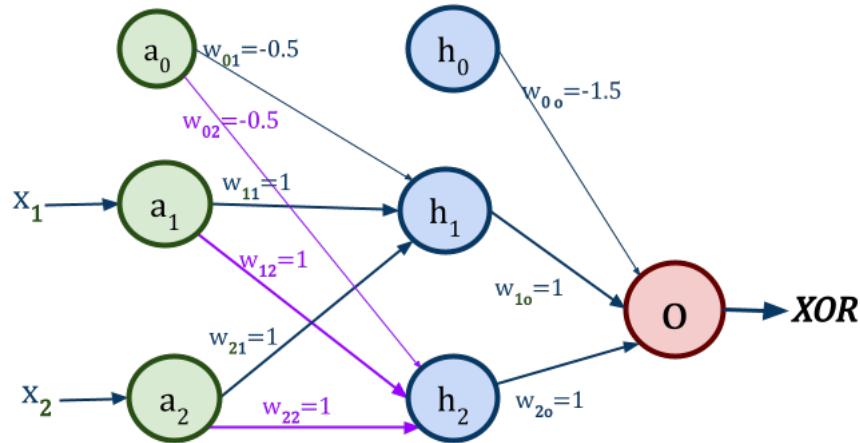


Figura 5.4 Función XOR, con capa oculta

Ahora veamos a las entradas como un vector de la siguiente forma:

$$A = \begin{bmatrix} a_0 \\ a_1 \\ a_2 \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \\ 1 \end{bmatrix}$$

Ahora para poder hacer los siguientes cálculos acomodamos los pesos correspondientes a cada perceptrón. Para el primer perceptrón tomamos los pesos que están conectados desde la neurona de origen  $a_0, a_1$  y  $a_2$  hacia la neurona oculta  $h_1$  (los índices de los pesos indican, el origen y el destino, en ese orden), los colocamos en un renglón de la matriz de pesos, así representamos nuestro primer perceptrón y hacemos lo mismo para  $h_2$ , así obtenemos la siguiente matriz de pesos:

$$W = \begin{bmatrix} w_{01} & w_{11} & w_{21} \\ w_{02} & w_{12} & w_{22} \end{bmatrix} = \begin{bmatrix} -0.5 & 1 & 1 \\ 1.5 & -1 & -1 \end{bmatrix}$$

Ya que tenemos la representación de las entradas y pesos en vectores podemos hacer la representación de la activación de estas neuronas de la capa oculta, que es de la forma  $H = g(W * A)$ , que de forma matricial lo podemos escribir de la siguiente forma:

$$H = \begin{bmatrix} h_1 \\ h_2 \end{bmatrix} = g \left( \begin{bmatrix} -0.5 & 1 & 1 \\ 1.5 & -1 & -1 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \\ 1 \end{bmatrix} \right) = g \left( \begin{bmatrix} 0.5 \\ 0.5 \end{bmatrix} \right) = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

La matriz resultante representa los valores de activación resultantes en cada perceptrón, el primer valor de la matriz representa el resultado de la activación de  $h_1$  y el segundo a  $h_2$ . Así ahora tenemos al vector  $H$  que sera el vector de entrada para la neurona de

## 5. Perceptrón multicapa

salida o, procedemos de la misma forma tomando en cuenta al sesgo  $h_0$  y a los pesos en forma matricial, nos queda de la siguiente forma  $o = g(W_{ho} * H)$ :

$$o = g \left( \begin{bmatrix} -1.5 & 1 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} \right) = g([0.5]) = [1]$$

Hasta este momento ya logramos resolver la función XOR para una sola entrada, esta forma de resolver una red la llamaremos convención 1. Si queremos que nos de la respuesta a varias entradas al mismo tiempo tendríamos que modificar un poco la representación de nuestros valores, trasnspriendo las matrices anteriores que teniamos. Que se verían de la siguiente forma:

$$A^T = [a_0 \ a_1 \ a_2] = [1 \ 0 \ 1]$$

$$W^T = \begin{bmatrix} w_{01} & w_{02} \\ w_{11} & w_{12} \\ w_{21} & w_{22} \end{bmatrix} = \begin{bmatrix} -0.5 & 1.5 \\ 1 & -1 \\ 1 & -1 \end{bmatrix}$$

$$H = g(A^T * W^T)$$

$$H = [h1 \ h2] = g \left( [1 \ 0 \ 1] \begin{bmatrix} -0.5 & 1.5 \\ 1 & -1 \\ 1 & -1 \end{bmatrix} \right) = g([0.5 \ 0.5]) = [1 \ 1]$$

$$o = g \left( [1 \ 1 \ 1] \begin{bmatrix} -1.5 \\ 1 \\ 1 \end{bmatrix} \right) = g([0.5]) = [1]$$

Esta forma se llama la convención 2, lo que estamos haciendo en la convención 2 es más bien poner nuestro ejemplar de manera horizontal. Así cada renglón de A va a representar una entrada el primer valor de cada renglón siempre será el sesgo. Así al tener varias entradas de la siguiente forma:

$$X = \begin{bmatrix} 0 & 0 & 1 & 1 \\ 0 & 1 & 0 & 1 \end{bmatrix}$$

A la matriz A es X transpuesta y le agregamos los sesgos, quedando así:

$$A = \begin{bmatrix} 1 & 0 & 0 \\ 1 & 0 & 1 \\ 1 & 1 & 0 \\ 1 & 1 & 1 \end{bmatrix}$$

La capa oculta quedaría de la siguiente forma  $H = g(W^T A)$

$$H = g \left( \begin{bmatrix} 1 & 0 & 0 \\ 1 & 0 & 1 \\ 1 & 1 & 0 \\ 1 & 1 & 1 \end{bmatrix} \begin{bmatrix} -0.5 & 1.5 \\ 1 & -1 \\ 1 & -1 \end{bmatrix} \right) = g \left( \begin{bmatrix} -0.5 & 1.5 \\ 0.5 & 0.5 \\ 0.5 & 0.5 \\ 1.5 & -0.5 \end{bmatrix} \right) = \begin{bmatrix} 0 & 1 \\ 1 & 1 \\ 1 & 1 \\ 1 & 0 \end{bmatrix}$$

El resultado de esta matriz  $H$  es que, la primera columna está representando la neurona que evalúa la función OR y la segunda columna representa el NAND. Para obtener la salida de la neurona o que es un AND, haríamos las mismas operaciones anteriores solo que con sus respectivos pesos y valores de  $H^T$ , así obteniendo lo siguiente  $o = g(H'W^T) = \text{XOR}$ :

$$o = g \left( \begin{bmatrix} 1 & 0 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 0 \end{bmatrix} \begin{bmatrix} -1.5 \\ 1 \\ 1 \end{bmatrix} \right) = g \left( \begin{bmatrix} -0.5 \\ 0.5 \\ 0.5 \\ -0.5 \end{bmatrix} \right) = \begin{bmatrix} 0 \\ 1 \\ 1 \\ 0 \end{bmatrix}$$

Más adelante esta representación, nos ayudará cuando estemos trabajando con redes neuronales que necesiten *diferentes tipos de características*, como; edades, número de hijos, tamaño de una casa, etc. Esta anotación es mucho más visual para la forma en que vamos a recibir la información (tabla de datos), por eso es más usada esta versión.

## Interpretación matemática del mapeo no lineal

En esta sección vamos a ver por qué agregar una capa más nos permite calcular funciones que no era posible calcular antes es decir no basta con tener más percepciones en realidad necesitamos tener percepciones que reciban como entrada la salida de percepciones anteriores y esto lo vamos a notar aquí por eso el selector es una función tan útil para poder explicar estos conceptos de entrada que teníamos al inicio de la función sort eran los valores 0 0 0 1 1 0 y 1 1 estos solamente los podría bueno puede utilizar perceptor únicamente para evaluar funciones que sean linealmente separables entonces cuando hemos aplicado la primera capa donde pudimos poner tantas operaciones como que hicimos lo que hicimos fue tomar estas entradas y mapear las a un conjunto distinto y esté ocurriendo en este nuevo conjunto

Observen que las entradas ahora son 01 11 11 10 es decir estas dos están duplicadas ya no tengo el 0 0 por eso es que se dice que trabajamos con un mapeo no lineal nuestra función sigmoide o nuestra función escalón transformaron nuestro espacio de manera que no tenemos ahora una relación 1 a 1 qué efecto va a tener. Aquí tenemos una vez más cuáles son las entradas originales las entradas que se obtuvieron es bueno salidas de las capas ocultas que van a hacer ahora las entradas para el último perceptor entonces veamos acá qué pasó con este cuadro aquí teníamos cuatro datos que no eran linealmente

separables después de la primera capa fue como haberlo plegado a otro espacio de dos dimensiones, pero observen ahora que como que esto es bueno vamos a ver quién se mató a quién el 0 0 se convirtió en 0,1 es decir este para acá el 0,1 se convirtió en 1,1 este fue para acá el 1010 también se convirtió en 1,1 es otro de acá y el 11 se convirtió en el 10 está aquí entonces nuestros puntos negros como que giraron un poquito y quedaron aquí y nuestros dos puntos blancos quedaron mapeados uno encima del otro en el mismo punto de este lado.

Entonces ahora sí lo puedo separar con un plano qué pasa en medio de los dos cualesquier que lo logre es bueno entonces estos tres valores quedaron mapeados ahora a un nuevo espacio que por cierto tiene una dimensión solo hay un valor. Entonces el paso fundamental radicó en que pudimos mapear este espacio hacia un espacio nuevo donde si es posible separar a nuestros datos linealmente y sobre esto pues ya simplemente aplicamos lo que podía hacer cualquier perceptual matemáticamente ese es el poder que nos está añadiendo una capa de en medio y básicamente lo podemos generalizar ahora si a cualquier función. Es posible quitar los sesgos si utilizamos la función escalón como función de activación y definimos precisamente el menor o igual para la parte del escalón en esta ocasión vamos a obtener valores que están exactamente en este punto de salto así es que vamos a tener que decir quién queda a la izquierda quien quiere a la derecha y entonces podemos hacer esto igual se necesitan dos capas, pero podemos quitarnos la parte de los sesgos porque nuestras fronteras si están pasando por ser un coma se hace solo por curiosidad bien aquí tenemos entonces ya esa interpretación vemos que los ceros negativos se van a tomar como ceros y solamente los positivos van a quedar como unos y automáticamente se puede hacer esta versión resumida.

## Propagación hacia adelante para el perceptrón multicapa

Ahora para el modelo de una red en general tenemos nuestra arquitectura base que va a ser una red en capas también se le conoce como el perceptron multicapa (tipo feed-forward) en esta primera versión tenemos la capa de entrada que realmente solo recibe las entradas y el sesgo. Las salidas de cada capa sirven de entradas a la capa inmediatamente posterior en la red multicapa. Por lo general todas las salidas de una capa se distribuyen a todas las neuronas de la siguiente capa, formando capas completamente conectadas (fully-connected layers). Cuando la red multicapa incluye más de una capa oculta, se dice que la red neuronal es profunda [deep neural network]. Los niveles de actividad de las neuronas de cada capa vienen dados por una función de activación no lineal de los niveles de actividad de las neuronas de la capa inferior.

(Insertar imagen de red).

Anteriormente habíamos estado asignando pesos a ojo/intuición, esto porque eran pocos los valores de entrada, esto normalmente no es así y vamos a desconocer en la totalidad los pesos necesarios que se aproximen a nuestra función objetivo.

La red neuronal es en sí, es una función que nos va permitir pasar un vector de n dimensiones a uno de m dimensiones, con n la cantidad de características en nuestros datos y m la cantidad de características que necesitamos obtener.

Entonces vamos a utilizar un tipo de codificación que se utiliza para la clasificación, llamado One-hot encoding, donde nos vamos a enfocar en el valor más elevado en nuestros resultados. Ahora nuestros problemas involucran más de dos clases, donde cada renglón de nuestra matriz de salida nos va indicar si pertenece o no a una clase, algún ejemplar dado. Así tenemos cada dimensión en la matriz va a representar una clasificación, por ejemplo si queremos distinguir en una imagen entre un coche, casa, animal lo podríamos codificar de la siguiente forma:

$$\text{Salida1} = \text{coche} = [1 \ 0 \ 0] \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}$$

$$\text{Salida2} = \text{casa} = [0 \ 1 \ 0] \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}$$

$$\text{Salida3} = \text{vaca} = [0 \ 0 \ 1] \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}$$

Toda la parte del aprendizaje, de reconocer distintos patrones, obtención de características varias, va a quedar entre *la capa de entrada y las capas ocultas*, para que al final nos quedemos únicamente con el problema de separar las clases unas de otras tantas sean necesarias, *en la capa de salida* y tengamos nuestra información clasificada. Así cada neurona de salida es una clase, específica, entonces si se activa esa neurona nos indica que nuestro ejemplar pertenece a dicha clase. Es importante no escatimar en el número de perceptrones necesarios para la clasificación a la salida, pues esto podría resultar en asignación conjunta de características, que se podrían interpretar como similitudes entre clases y en caso de no existir, dificultaan enormemente la distinción entre características y generar fallos a la red. Siempre asignar tantas dimensiones de salida (perceptrones) como clases.

Los cálculos para cada capa intermedia realmente son las mismas que hemos visto anteriormente representándose de la siguiente forma:

$$a_j^{(l+1)} = g \left( \sum_i w_{ij}^{(l)} a_i^{(l)} \right) \quad (5.8)$$

Con  $l + 1$  representando el número de capa a la que vamos,  $g$  la función activación (normalmente podemos usar la función sigmoide) y  $a_i$  los valores de las entradas.

# 6 | Entrenamiento por retropropagación

## Esquema general entrenamiento

En este capítulo vamos a trabajar con el entrenamiento, en esta sección abordaremos primero el esquema general del entrenamiento así pues veamos la derivación del primer algoritmo de entrenamiento para perceptores multicapa le llamó primero porque también es ahorita el más utilizado el más común y el primero que se estableció históricamente vamos a ver cómo se entrena las redes neuronales primero que nada vamos a definir en qué consiste aprender algo cuando estamos hablando concretamente de redes neuronales recordemos que ya habíamos hablado antes de los espacios de y poder que realmente aprender se trata de encontrar una función dentro de un espacio de funciones posibles y la característica que tiene esta función es que aproxima una función que modela algún fenómeno que nos está interesando a aprender entonces vamos a definir el problema de aprendizaje para un percepción multicapa de la siguiente manera dada una arquitectura de red neuronal encontrar los pesos tales que la función definida por la red neuronal approxime hasta cierta tolerancia la función de interés  $f(x)$  voy a entrar en detalles sobre esta anotación  $x$  son nuestros datos de entrada pueden consistir en diferentes características dependiendo del problema del que estamos trabajando puede ser por ejemplo segmentos de alguna canción pueden ser imágenes de audio pueden ser valores referentes a características de una casa tratando de predecir sus preciosos en su región entonces la  $x$  son precisamente todas esas características en general  $x$  puede ser un suelo ejemplar que sería el que va a comenzar nuestra función de tal manera que tengamos  $f(x)$  que va a ser mapeado a alguna clase estas redes que estamos viendo aquí lo que están haciendo es resolver problemas de clasificación queremos saber si dada una imagen es por ejemplo una casa una vaca un perro dada una canción si es una canción triste una canción alegre una canción tranquila entonces estamos mapeando podríamos pensarlo así al espacio de clases entonces la función  $f(x)$  sería aquella función ideal que recibe una imagen y nos dice inmediatamente que es recibe un caco de una canción y nos dice inmediatamente a qué género pertenece o qué estado de ánimo representa lo que nosotros vamos a hacer con la red neuronal es tratar de aproximar esta función que asumimos que existe en el universo y ahora estamos tratando de ver cómo calcularla este  $h$  viene precisamente del nombre hipótesis porque es la hipótesis que está planteando nuestra red neuronal puede ser que sea una buena aproximación de esta función o puede ser que sea mala lo que

nosotros queremos hacer es encontrar pesos en esta red neuronal tal que nuestra aproximación sea lo mejor posible entonces aquí entra también el elemento de cierta tolerancia es probable que a veces no podamos aprender exactamente esta función y que siempre encontraremos errores o alguno que otro ejemplar que no va a ser clasificado correctamente bueno dependiendo de nuestro contexto tendremos que decir hasta qué momento se pueden considerar perdonables esos errores o si tenemos que continuar busca recordemos que en general tener demasiada precisión en los datos con los que entrenamos puede significar que estamos aprendiendo ruido y que después cuando recibamos datos nuevos no vamos a tener buenas predicciones entonces aquí vamos a tener que ir aprendiendo a lo largo del curso cómo podemos evaluar bueno cómo fijar hasta dónde entraría está tolerancia entonces uno se acostumbra entrenar a las redes neuronales así insistir una vez más en este punto una red neuronal con un conjunto de pesos dados ya define una función es a la que le estamos llamando  $h$  de zeta en los videos anteriores aprendimos a asignar estos pesos inclusive manualmente eligiendo nosotros una recta en una representación gráfica entonces ponerle los pesos a la red o tener sin presencia mente una red con pesos ya es una función a veces podemos encontrar directamente esos pesos de tal manera que nos modelan correctamente la función que nosotros queremos a veces el problema va a ser precisamente que no conocemos pesos y los queremos encontrar ahí es donde va a entrar el elemento de aprendizaje entonces a qué le vamos a llamar entrenamiento o cómo vamos a entrenar a una red para que aproxime cada vez más a la función que nosotros queremos el mejor nuestro punto inicial es una función que no se parece como cambiamos eso como lo arreglamos la estrategia de entrenamiento para encontrar un conjunto de pesos que satisfaga este requerimiento o sea aproximar correctamente la función consiste en lo siguiente uno definir una función de error o pérdida le vamos a llamar  $J$  en este momento otra anotación común que encontrarán en internet es el de las perdidas que mida la distancia entre los valores deseados y los valores obtenidos con un conjunto de pesos dados  $z$  entonces recordemos aquí que estamos en un tipo de entrenamiento supervisado este es para casos de aprendizaje supervisado donde tenemos un conjunto de valores deseados es decir ya sabemos que tenemos una función que queremos aproxima el siguiente punto es utilizar alguna técnica de optimización de funciones que minimice este error y en este momento es querido enfatizar la parte alguna técnica puede ser cualquiera cualquier técnica de optimización que nos permite minimizar el error la podemos intentar utilizar para entrenar redes neuronales algunas funcionarán mejor que otras dependiendo del espacio en el que estemos buscando pero podemos elegir entre todas ellas en particular en algún momento vamos a ver otra familia que no entra en las creamos ahorita como son los algoritmos genéticos también es una técnica de optimización entonces también se puede llegar a utilizar entonces tanto la arquitectura considerar como la función de error y la técnica de optimización dependerá del tipo de función que se desee aproximar ahorita vamos a verlo de manera abstracta de manera general un poco a poco a lo largo del curso iremos viendo precisamente diferentes tipos de funciones y qué arquitecturas se recomiendan qué estrategias funcionan para entre vamos a ver ahorita un esquema general este esquema es el que se conoce como entrenamiento por retropropagación y propagación hacia atrás también se le suele llamar en inglés realmente el único nombre es barack corporation también frecuentemente se encontrarán que en la literatura en español algunos

## 6. Entrenamiento por retropropagación

no lo traducen y lo dejan tal cual en inglés vamos a utilizar aquí usualmente el término retro propagación que significa bueno ustedes habrán dado cuenta de que tenemos un primer alimento que era alimentación hacia adelante en inglés feedforward lo que vamos a hacer para entrenar a la red es recorrer los pesos en dirección inversa desde el resultado hacia los datos de entrada y por eso lleva precisamente el nombre contrario al algoritmo de evaluación de la red normal bien entonces este algoritmo consiste en la siguiente combinación este es históricamente en primer lugar utilizar la regla de la cadena para obtener el gradiente de la función de error con respecto al conjunto de pesos dado un conjunto de datos de entrenamiento  $x$  con sus etiquetas el segundo paso utilizar descenso por el gradiente para encontrar un mínimo lo suficientemente bueno de la función de error aquí es donde entra precisamente la perspectiva histórica utilizar descenso por el gradiente es la técnica más sencilla que existe por eso es interesante empezar por aquí precisamente a plantear la metodología pero en la actualidad se utilizan ya métodos de optimización que dependen del gradiente bastante más avanzados entonces en qué consiste conceptualmente esta idea en que nosotros al inicio no sabemos cuánto deben de valer los pesos para modelar correctamente la función entonces vamos a comenzar con una asignación aleatoria de pesos suponiendo que éste sea una versión caricaturesca de una semana muy simple vamos a suponer que esta parte de aquí abajo representa el conjunto de pesos en nuestra red neuronal elegimos algún punto al azar todo ese punto calculamos el error que está cometiendo la red al comparar los datos que obtuvo la red neuronal al aplicar feedforward sobre nuestros datos de entrada con respecto a la respuesta que queríamos que nos dieran si nosotros asignamos valores al azar evidentemente no tenemos por qué esperar que la aproximación sea buena lo que va a ser ahora por nosotros el descenso por el gradiente es ir modificando poco a poco estos pesos para que nosotros tengamos ahora una nueva aproximación y cada vez nos acerque más a una región donde el error sea más pequeño y aquí importa también muy bien en la parte que dice para encontrar un mínimo lo que es el descenso por el gradiente es lo siguiente comenzamos con un conjunto de propuestas para nuestros parámetros de tal calculamos el gradiente es decir las derivadas de la función de error con respecto a los parámetros de aquí vamos a obtener el gradiente este gradiente lo que nos va a dar es la dirección de máximo ascenso de la función de error al agregar este signo menos lo que estamos haciendo es invertir ese vector de la dirección de máximo acceso y obtener por consiguiente la dirección de máximo descenso entonces con este vector que básicamente aproxima en la tangente de la curva pero además nos da la dirección en la que va a descender lo más rápido posible lo que hacemos es modificar los parámetros que habíamos propuesto exactamente en esa dirección para que nos lleve lo más rápido posible por donde está descendiendo la pena y con eso calculamos la nueva propuesta de parámetros para la red aquí les dije que es un dibujo un poco caricaturesco porque es un segmento de paraboloide en el caso de las redes neuronales cuando ya tenemos varias neuronas y muchas conexiones la curva no es un paraboloide con un único mínimo en realidad vamos a tener un montón de mínimos y máximos locales extendiéndose en varias direcciones junto con algunas regiones un poco peligro grandes picos y es un paisaje bastante interesante lo que va a hacer por nosotros entonces descenso por el gradiente es acercarnos al mínimo más cercano que nosotros tengamos como se entrena entonces realmente una red neuronal bueno como el mínimo

el que lleguemos depende de el conjunto original que hayamos dado de parámetros es para el entrenamiento lo que vamos a hacer es inicializar aleatoriamente estos datos para obtener puntos en diferentes regiones del espacio y repetir el entrenamiento varias veces aquí el entrenamiento que nos logre llevar al mejor mínimo que satisfagan nuestros requerimientos de haber aproximado lo suficientemente bien la función ese es el conjunto con el que nos vamos a quedar una vez que ya hemos evaluado nuestra red y que hemos visto que es buena para hacer predicciones inclusive en otros conjuntos de datos que no sean con los que entrenamos y así que para eso nos sirve el conjunto de validación entonces podemos decir que ya tenemos una red que podemos utilizar para hacer adiciones en otras versiones algunos sistemas a la mejor hasta podrían guardar una colección de las redes que quedaron mejores entrenadas y utilizar algún sistema de votación para tomar decisiones al momento de utilizarlas ya en producción lo que estamos viendo aquí a la derecha es una imagen de los cortes que se podrían hacer horizontalmente de esta misma imagen entonces también podemos visualizar a nuestro conjunto de parámetros  $z$  y coordenadas sobre este plano y en estas que son las curvas de nivel podemos ir tratar todo visualizar buenas ya que hasta el mínimo hacia que estaría como está subiendo el paraboloide y podemos ver también como la dirección negativa del gradiente pues nos iría dirigiendo poco a poco hacia este centro donde tenemos uno entonces actualmente se puede cambiar la técnica de optimización por algún otro método más avanzado también dependiente del gradiente el hecho de que depende del gradiente pues va a ser que en la siguiente sección sea sumamente importante que es como calculamos el gradiente de la función de error la derivación original utilizaba diferencias al cuadrado como función de error recordemos que es aquella en la que deseamos cuánto es lo que quiero obtener menos mi hipótesis de lo que evalúa un mínimo función de la red neuronal sobre los datos elevados al cuadrado y después sumamos sobre todas las y todos los ejemplares entonces esta fue la primera que se utilizó sin embargo se ha comprobado que esta función es adecuada para problemas de regresión pero para problemas de clasificación es mejor utilizar otra que se le conoce como entropía cruzada y ven artículos un poco viejitos o inclusive de otros campos van a ver que es muy frecuente encontrar diferencias al cuadrado por ello si quieren ver esa derivación se puede consultar en el libro de rosalino link les recomiendo la tercera edición creo que la anotación es más clara en la tercera edición que en la segunda lo que vamos a hacer nosotros aquí es ver la derivación con entropía cruzada que es básicamente ya el estándar cuando queremos empezar a trabajar con problemas de clasificación vamos a ver también que tiene algunas propiedades bastante bonitas y nos va a facilitar las cosas.

## Función de error: Entropía cruzada

Vamos a introducir entonces a nuestra función de error para problemas de clasificar y la entropía cruzada tiene su origen en teoría de la información mide la cantidad de bits necesarios para identificar una clase dada la hipótesis de la red siendo que éstas provienen de la función de  $x$  este  $g(x)$  vendría a ser precisamente los valores reales a

## 6. Entrenamiento por retropropagación

nuestras etiquetas lo que ocurre en el mundo experimental esta sería la red neuronal que está tratando de codificar con bits la información de a qué clase pertenecen estos datos entonces lo que estamos tratando de predecir es si esta función se parece a esta entonces va a ser muy fácil codificar en cualquier este clase de estos datos utilizando esta función si no se parecen entonces necesitaríamos un montón de bits y va a ser muy eficiente básicamente no nos va a funcionar si lo queremos visualizar más gráficamente me gusta utilizar la siguiente explicación estos cheques que tenemos aquí son nuestras etiquetas es decir los valores correctos que nosotros creímos vemos que en esta operación aparecen dos términos aparece y aparece uno menos y bueno si estamos tratando de un problema de clasificación recordemos que entonces nuestras neuronas de salida van a tener dos posibles valores 0 y 1 la respuesta correcta es 0 o es 1 como tenemos una función de activación logística en este caso nunca vamos a tener exactamente 0 o exactamente 1 pero sí nos vamos a acercar lo más que podamos dependiendo de los pesos que peta recordemos que las magnitudes afectaban s entonces estos días van a representar precisamente esos valores 0 y 1 y lo que tenemos acá es el logaritmo de lo que obtuvimos vamos a considerar entonces los dos casos supongamos que teníamos un cero pero lo que predijo nuestra red fue uno lo que iba a suceder entonces bueno de estos dos términos observemos que si lleva el es cero este término simple y sencillamente no va a contar nos vamos a quedar únicamente con este en este caso pasa lo siguiente si lleva el es cero aquí es un 1 observen que realmente el error va a estar dado por esta función de acá si lo que obtuvimos aquí otra vez insistes cercano a uno porque tenemos una función logística entonces qué va a suceder con este logaritmo vamos a tener un número cercano a cero que ocurre con el logaritmo cuando nos vamos a cero bueno este le funciona carisma lo que está pasando es que nos estamos disparando hacia menos infinito observemos este signo menos que se coloca acá en la definición de entropía cruzada el efecto que va a ser entonces es que si queríamos un valor cero y obtuvimos algo cercano a uno bueno entre más nos acerquemos a uno esto más se nos va a disparar hacia infinito como quería uno una razón y nos estamos yendo para acá muy bien y pasa entonces si si queríamos un 0 y si obtuvimos un 0 bueno lo que vamos a tener aquí entonces es el logaritmo de 1 es cero entonces si la respuesta era correcta nuestra función de error va a valer cero perfecto otra ves como es una función sigmoide bueno nunca vamos a tener exactamente cero es aquí el error nunca va a ser exactamente cero pero va a estar cerca vamos a ver qué ocurre ahora en este otro caso bueno supongamos ahora que queríamos que valiera pero lo que obtuvimos fue algo cercano a cero en ese caso una vez más esto se convierte en uno entonces realmente la función de error está dada por el término de este lado tenemos entonces el logaritmo de algo cercano a 0 que se nos va a disparar hacia menos infinito multiplicado por el menos que tenemos acá afuera bueno pues vamos a completar de hecho el otro lado de esta gráfica lo que tenemos es que estamos hacia acá arriba y queríamos un valor 1 en el aire y obtuvimos un 1 bueno aquí se queda el 1 aquí nos queda el logaritmo de 1 que es 0 entonces vamos a estar hacia este lado de esta manera lo que podemos pensar entonces es que esta función de error tiende a ser lo siguiente si nos estamos equivocando se dispara hacia infinito y estamos en lo correcto vamos a estar muy cerca de cero cuando estamos hablando ya de un sistema complejo como una red neuronal con muchísimas neuronas varias capas varias neuronas en la capa de salida pues la gráfica en realidad no va a ser

este paraboloide bonito va a ser una gráfica con bastantes rugosidades mínimos locales etcétera pero digámoslo así es menos rugosa que si hubiéramos utilizado diferencias al cuadrado entonces por eso vamos a utilizar esta función bueno ya leyendo los detalles de la anotación vamos a fijarnos entonces que estamos obteniendo un promedio sobre el número de ejemplares de entrenamiento o el número de ejemplares en general sobre los cuales estamos evaluando en este momento el error que comete nuestra red y eso va a ser  $m$  es por eso tenemos la suma desde uno hasta  $m$  de todas las contribuciones de error y  $1/m$  para que estemos hablando de un promedio  $s$   $I$  es el número de neuronas y capa de salida recordemos que en general nuestro perceptor multicapa va a tener una neurona por cada clase que nosotros estemos tratando de evaluar entonces tiene sentido que sumemos el error para cada una de las clases con las cuales estamos trabajando y observemos que cada una de esas clases efectivamente está representada por algún valor 0.1 en su neurona correspondiente entonces aquí estamos sumando sobre todas las neuronas de salida que corresponde a todas las clases que estamos tratando de identificar con esta red y que es el valor deseado para la carísima neurona de salida y toma valores en 0.1 bueno entonces cada una de estas neuronas tiene que tener etiqueta corresponde y finalmente  $htc$  está es la red neuronal evaluada en el ejemplar y estima.

## Derivada de la función logística

Podemos entrar ahora a la parte más interesante de las matemáticas como calculamos el gradiente de la función que acabamos de definir este sería la notación general que vamos a utilizar para nuestra red neuronal vamos a empezar ahorita fijándonos en un solo perceptor y vamos a hacer un cálculo que nos va ayudar bastante a simplificar al rato los datos vamos a fijarnos en este perceptor está recibiendo estos otros como entrada observemos que tiene y una serie de pesos que lo están alimentando y recordemos que nuestra función de activación ahorita es la logística  $1/(1 + e^{-z})$  a la menos  $x$  entonces cuando nosotros hablamos de sacar la derivada recordemos qué vamos a querer ser nuestras variables al momento de entrenar en realidad son los pesos no son los datos de entrada vamos a considerar para un ejemplar dado cuál es el error que está cometiendo mi percepción y voy a tratar de modificar los pesos para reducir ese error entonces el ejemplar de entrenamiento está fijo inclusive cuando entrenamos utilizando ya varios ejemplares de entrenamiento el conjunto de ejemplares de entrenamiento está fijo lo que nosotros estamos modificando son los pesos por ello cuando nosotros calculemos la derivada lo que nos va a interesar es bueno voy a cambiar esta  $x$  por  $z$  donde recordemos que  $z$  era precisamente la combinación lineal de los valores de entrada multiplicados por los pesos y sumados todos hechos entonces eso es lo que vamos a hacer en la siguiente diapositiva recordemos simplemente por ilustración cómo se ve la función sin muy bien aquí precisamente ya lo que estoy haciendo es sustituir por la multiplicación del ejemplar por los pesos observen que al tener la notación matricial esto me va a permitir más adelante utilizar aquella anotación en la que  $x$  tenía hacia abajo todos los ejemplares de entrenamiento que nosotros queramos hacia la derecha todas las dimensiones que

## 6. Entrenamiento por retropropagación

necesite para las características de entrada y recordemos también que el vector  $z$  lo que tenía era sobre las columnas los pesos correspondientes para calcular las entradas a cada una de las neuronas de salida vamos a decir que tenemos  $sl$  neuronas entonces para pensarlo ahorita de manera un poco cómoda me limito al caso de un solo ejemplar este  $x$   $z$  sería una versión abreviada de escribir  $x \beta_0 \beta_1 \dots \beta_n$  pero con ceros unos con unos hasta  $x_n$  entonces esto es lo que tengo realmente escrito todo es afectado por el signo menos entonces vamos a calcular la derivada de esta función con respecto a cualquiera de estos parámetros y lo vamos a llamar  $z'$  bueno por la regla de la cadena comenzamos calculando la derivada de  $1$  entre una función sería menos  $1$  entre lo que tengo aquí abajo elevado al cuadrado por la derivada de lo de adentro el  $1$  sería constante entonces no pasa nada me queda solamente ese término derivada del exponencial pues es la misma exponencial derivada de  $x$  la derivada de lo que esté en el exponente ahora en el caso del exponente recordemos que es precisamente el menos multiplicando a todo esto que está acá entonces el menos lo vamos a sacar aquí y como la parcial le estoy sacando con respecto a uno de los pesos lo único que va a sobrevivir de la derivada de este exponente va a ser precisamente la  $\sigma$  que esté acompañando a la  $\theta$  con respecto a la cual estoy derivando entonces por eso tenemos aquí la presencia de este  $xy$  solito ya que hicimos esto el resto va a ser utilizar un poco de trucos algebraicos para escribirlo de una manera que nos resulte mucho más cómoda realmente el cálculo de la derivada ya lo terminamos hasta aquí vamos a ver ahora qué propiedad descubrimos si tratamos de reescribir esto bueno en primer lugar vamos a cancelar este menos con este menos vayan dando los cargando los dio más este que este que vamos a separar ahorita el dividendo de manera que quede uno entre este elemento la vamos a subir encima del otro observemos que como se están multiplicando pues realmente no hemos hecho nada sólo fue una reescritura y la  $\sigma$  y la vamos a tener ahorita acompañándonos de adorno todo el tiempo así que así recuerden que está aquí en la derecha y ya no nos tenemos que volver a preocupar de ella de ahí en fuera entonces qué vamos a hacer ahora con estos dos elementos lo siguiente es darnos cuenta de que lo que escribimos intencionalmente aquí pues es exactamente  $\sigma$  ah qué bonito entonces la derivada de esto es lo mismo por otro término que está acá pero resulta que de este término también podemos decir algo interesante el truco favorito de los profesores de cálculo vamos a sumar y restar  $1$  lo cual es un  $0$  entonces realmente no modificado nada y otra vez vamos a separar convenientemente cada uno de los términos en particular vamos a observar otra vez este patrón no se hace conocido a la pista vamos a poner este término del lado derecho y en verdad yo quiero agarrar con uno con este les va a quedar menos  $\sigma$  otra vez y los otros dos pues lo ponemos aquí pero a que estamos viendo son idénticos perfecto entonces ya tenemos signo esto es realmente dan uno y estoy triste que siguen otra vez pero con signo negativo entonces lo que tenemos es una propiedad bastante interesante de la función logística de hecho acá lo volví a escribir para que se vean más claro simplemente estoy calculando la derivada de  $\sigma$  con respecto al exponente osea que sería esto completo en lugar de meterme los detalles y entonces lo que vemos es que la derivada de la sigmoida es la sigmoida por  $1 - \sigma$  la sigmoida en este caso bueno como nos interesa en los pesos por eso tenemos aquí un nivel más en la regla de la cadena y nos va a parecer este  $xy$  si tienen un poco de curiosidad de la forma que toma esta

derivada pues aquí tenemos precisamente la gráfica gracias a esta propiedad estamos en la forma más general vamos a ver que obtener la derivada del gradiente completo se nos va a facilitar mucho entonces vamos a pasar después a la derivada de la función de error.

Las funciones logísticas se utilizan a menudo en redes neuronales para introducir no linealidad en el modelo o para sujetar señales dentro de un intervalo específico . Un elemento de red neuronal popular calcula una combinación lineal de sus señales de entrada y aplica una función logística limitada como función de activación al resultado; este modelo puede verse como una variante "suavizada"de la neurona umbral clásica .

Las funciones logísticas se utilizan en varios roles en estadística. Por ejemplo, son la función de distribución acumulativa de la familia logística de distribuciones y, un poco simplificadas, se utilizan para modelar la posibilidad que tiene un jugador de ajedrez de vencer a su oponente en el sistema de clasificación

Estas relaciones dan como resultado implementaciones simplificadas de redes neuronales artificiales con neuronas artificiales . Los médicos advierten que las funciones sigmoidales que son antisimétricas con respecto al origen (por ejemplo, la tangente hiperbólica ) conducen a una convergencia más rápida cuando se entrena redes con retropropagación.

La función logística es en sí misma la derivada de otra función de activación propuesta, el softplus Una opción común para la activación o ."plastamiento"funciones, usadas para el clip para grandes magnitudes para mantener la respuesta de la red neuronal limitada

## Parcial con respecto a los pesos en la penúltima capa

una vez más para ser más legible los cálculos vamos a comenzar a calcular el gradiente considerando que ahorita solamente estamos evaluando a la red en un ejemplo y vamos a tener que resolverlo por etapas lo que estamos pensando es lo siguiente nuestra función de la red neuronal como aquí las entradas produjo aquí las etiquetas correspondientes y ahora lo que nuestra función de errores están viviendo es la distancia entre lo que obtuvimos en esta capa y una colección de etiquetas que era lo que nosotros queríamos que saliera entonces nuestra función de error está actuando directamente sobre esta capa nada más pero cómo se llegó a este resultado bueno pues se llegó después de haber realizado cálculos sobre lo que teníamos evaluar acá utilizando los pesos en esta capa cómo se llegó este resultado pues habiendo evaluado lo que había ocurrido acá involucrando también a los pesos de esta otra capa entonces podemos pensar que los valores que obtuvimos aquí a la salida dependen de los pesos en cada una de las capas sin embargo la dependencia es ligeramente distinta en el sentido de que la distancia entre este peso y este es inmediata aquí para poder saber cómo depende este resultado final de un peso que tengo acá atrás pues tengo que irme dos niveles más hacia acá recordando que este influyó en el cálculo de este número entonces el gradiente lo vamos a ir calculando cappa porque de ahí viene precisamente el nombre de bach preparation entonces lo primero y lo más sencillo va a ser calcular como dependió la función de error

## 6. Entrenamiento por retropropagación

de los pesos que se encuentran inmediatamente detrás de las neuronas de salida y así verdad necesito explicar poco a poco la anotación lo que vamos a hacer en este momento es considerar que en el algoritmo feedforward la primera capa la vamos a tomar con los índices y la siguiente capa índices  $j$  la tercera capa índices  $k$  por eso estamos pensando que los pesos de la última capa conectan a la neurona  $j$  con la neurona  $k$  pero para considerar que vamos a ir haciendo esto en reversa nos importa que esta es la última capa no le vamos a llamar la capa  $l$  y luego vamos a trabajar con la capa anterior que sería la capa  $l$  menos 1 luego vamos a ir a la capa todavía anterior  $l$  menos 2 y esta anotación también nos va a permitir darnos cuenta que si hubiera más capas hacia atrás pues siempre y sencillamente iríamos restando más números otro punto importante en la anotación vamos a asignarle a los pesos el índice correspondiente a la capa anterior a la cual tuvieron que multiplicar entonces cada capa va asociada con el bloque de pesos que tiene adelante el cálculo final entonces para la capa  $l$  depende de los valores de activación en la capa  $l$  - sólo  $\times$  los pesos en la matriz  $l$  menos 1 y así obtenemos lo que está en la campaña memorizar a la anotación porque la vamos a necesitar o bien entonces lo que tenemos en esta línea inmediatamente es la función de error observen que quite ahorita la suma sobre los diferentes ejemplares de entrenamiento por eso solamente lo vamos a calcular para un ejemplar y al resto al ratito ya será sencillo generalizar a varios ejemplares bien entonces si tenemos un suelo ejemplar de entrenamiento lo que si nos sigue importando es el hecho de que tenemos varias clases de entre las cuales podría haber estado la respuesta y por eso tenemos la suma desde que hay igual a 1 hasta  $s_l$  entonces a través un poco de notación  $s_l$  es el número de neuronas que tenemos en la capa  $l$   $s_l$  menos 1 sería el número de neuronas en la capa  $l$  menos suelo y así sucesivamente y aquí viene la misma definición que ya teníamos antes el valor que queremos llegar y el valor que realmente obtuvimos de hecho podemos quitar ahorita este 6 porque solamente hay un ejemplar de entrenamiento entonces ya tenemos aquí un solo ejemplo y ahora lo que queremos hacer es calcular la derivada con respecto a los pesos justamente antes de haber evaluado el valor de la última capa la notación que estamos utilizando aquí en vez de la  $h$  de hipótesis es el hecho de que tal cual el valor de la hipótesis es el valor de activación de la neurona entonces estamos hablando del valor de activación de cada una de estas y esas son las áreas si se fijan cuando evaluamos el algoritmo feedforward ya obtuvimos los valores de activación de hecho de todas las neuronas entonces los podemos utilizar aquí y precisamente lo que vamos a hacer para poder calcular la derivada es recordar hecho eso de nuestro valor de activación fue calculado con las reglas de activación de un perceptor entonces comenzamos a calcular la derivada de la entropía cruzada en primer lugar vamos a observar que estoy calculando la parcial del error con respecto a uno de los pesos bueno aquí tenemos la suma sobre todo de las neuronas de salida no pero observemos qué ocurre con este peso solamente está contribuyendo a la salida de esta neurona a esta neurona no le afecta en lo más mínimo lo que haya ocurrido aquí y esta neurona tampoco le afecta luego entonces la parcial de estos dos elementos con respecto a este peso vale cero porque no hay una dependencia son constantes en este caso si para la neurona que está conectada con este peso es decir la neurona acá sí importa lo que ocurrió con este peso por ello de los tres elementos que teníamos aquí aparece en el dibujo es el en general solamente va a sobrevivir uno de los términos y

es aquel donde estamos trabajando exactamente con la neurona carísima entonces si de momento ya no tenemos esta suma podemos trabajar solamente con lo que hay aquí adentro tenemos nuestro signo menos acá afuera después llegué acá 1 - que realmente son constantes entonces aquí sale la aie de acá la derivada del logaritmo es 1 entre lo que está dentro del logaritmo por eso tenemos aquí acá en el denominador regla de la cadena derivada de lo de adentro entonces tenemos otra vez la derivada con respecto a  $z_{jk}$  del valor de activación acá y aquí es importante que recordemos cómo era calculado este valor y era precisamente con la función de activación aplicada sobre la combinación lineal de todos los pesos que venían desde la capa anterior y que nos conectaban con este valor y estaban multiplicados pues por los valores de activación precisamente de la capa anterior es decir estamos considerando que esta neurona tiene las contribuciones de todas estas neuronas de acá y entrando para acá y eso es lo que estamos escribiendo aquí de manera explícita igualmente este es un término que ya habíamos calculado para poder evaluar el fit forward las setas que son las combinaciones lineales de los valores de activación de atrás por los pesos entonces si recordamos eso podemos simplificar nos un poco la nota y haciendo lo mismo con este otro término otra vez tenemos este que es una constante después la derivada del logaritmo 1 entre lo de acá por la derivada parcial de lo que esté aquí adentro la deriva de parcial de 1 pues va a ser un 0 de este menos lo sacamos y lo estamos poniendo aquí dejamos indicado de momento la derivada parcial de aquí tenemos estoy escribiendo en versión resumida lo mismo que tenemos acá porque es exactamente bien ahora si observamos este elemento otra vez podemos aprovechar el hecho de que solamente nos interesa este peso entonces realmente lo que ocurre con estas otras dos neuronas de momento no nos está afectando solamente nos interesa el valor de  $a$  que conecta con este esa jota exactamente los 102 las fotos primas que tenemos aquí solamente nos vamos a quedar con esto si hacemos eso entonces ya no necesitamos considerar aquí toda la zona para el siguiente término es el que me falta por aquí es decir calculé la derivada parcial de todo lo que estaba acá adentro con respecto a  $z$  acá pero pues ahora que calculemos la derivada pues vamos a tener que seguir pues entonces como se va a ver esto en primer lugar este término se copia tal y como está el menos aquí lo seguimos cargando que habíamos visto anteriormente acerca de la derivada de la función logística pues que es la función logística por 1 - ella misma es esa primera parte la vamos a dejar aquí indicada por otro lado tenemos que continuar con la regla de la cadena calculando la derivada de esta suma que tenemos aquí adentro pero lo que acabamos de ver es que solamente va a sobrevivir pues el elemento donde jota prima es jota y los otros se van a morir son constantes con respecto a este peso entonces la derivada de este elemento con respecto a éste es uno y el término que queda sobreviviendo es esta y es precisamente entonces le estoy poniendo aquí afuera de una vez porque nos va a salir otro idéntico del lado derecho recuerdan de la  $x$  que teníamos cuando calculamos la deriva de la signo y bueno aquí está precisamente el término que está sobreviviendo muy bien entonces repitiendo la misma idea de este lado tenemos el mismo factor tenemos que prima y también aquí al aplicar la regla de la cadena va a salir una idénticamente entonces por eso le estaba factorizando de una vez lo siguiente que vamos a hacer es bueno factorizar este también lo estamos poniendo aquí afuera y otra vez lo que vamos a hacer es trabajar con estos dos para tratar de escribirlos de una manera más amigable

## 6. Entrenamiento por retropropagación

ponemos entonces como un denominador estos dos se multiplican este que está aquí lo multiplicó por el denominador de acá estoy acá lo multiplicó por el de acá tenemos el signo menos los voy a llegar por 1 menos acá menos acá por un número chica una vez que hacemos esto las cosas se vuelven hermosas esto de aquí es rica esto de acá es un término cruzado allí - una al menos por menos más un término cruzado allí éste se va con éste ahora tenemos que recordar también que propiedad bonita tenía la derivada prima pues era  $g$  por uno menos 100 pero la función  $g$  evaluada en la combinación lineal que es pues es exactamente la que calculamos cuando estábamos haciendo el fit for work así salió entonces este es que es que tengo aquí le sustituyó simple y sencillamente por la sas que fue lo que calculamos en el instante en el que recuerdo esto o miren esteban y fue los que nos quedó entonces pues únicamente la diferencia entre lo que quería y lo que salió multiplicado por el valor de activación de la neurona en la capa anterior bien porque es bonito ser entropía cruzada consigo les basta con legis ticas muy bien entonces a esto que esté aquí se le acostumbra poner el nombre del está acá porque podemos decir que fue el error que se cometió en la última capa y aquí quedó bastante literal simplemente tomamos lo que salió bueno lo que queríamos - lo que salió este elemento que está acá es lo que está contribuyendo precisamente el hecho de que este peso no se estaba conectando con alguien en la parte de atrás y ya tenemos entonces la manera de calcular todas las parciales del error con respecto a los pesos y en esta capa vamos a ver entonces en el próximo capítulo cómo calcular las derivadas con respecto a los pesos que están en la capacidad.

### Parcial con respecto a los pesos en la última capa

si ya nos quedó claro cómo funciona la derivación para el peso anterior entonces trabajar con esta capa ya no va a ser tan complicado vamos a tener un poco más de elementos pero va a ser sencillo agregarlos en primer lugar vamos a tener que observar quienes se ven afectados por el nuevo peso con respecto al cual queremos derivar bien aquí tenemos un poco más de transmisión de efecto este peso solamente afecta a esta neurona pero a través de esta neurona ahora si estamos conectados con absolutamente todas las otras neuronas entonces aquí vamos a tener que hacer los toman en cuenta más términos clientes para poder hacer esta derivación voy a tener que volver a empezar desde un principio es decir vamos a volver a empezar con la función de error original y vamos a tener que empezar a derivar desde aquí vamos a recordar ahora que cada vez que calculamos los valores de activación en la última capa bueno estos provienen en realidad también dependen de el valor de activación en esta neurona de acá atrás entonces eventualmente vamos a tener que regresar hasta acá y el valor de esta neurona pues dependió ahora sí directamente del peso que nos está interesando entonces vamos a ir desarrollando poco a poco esa composición de funciones entonces comencemos con el mismo paso de antes calculamos la derivada parcial de jota con respecto al peso correspondiente pero ahora si no puedo eliminar la suma que tengo al inicio porque todos los valores de salida dependen de el peso con el que estamos trabajando entonces lo que

vamos a tener aquí ahorita es que nos quedamos con la suma lo único que estoy haciendo ahora es pasarla junto con el signo menos recordemos que la derivada de una suma es la suma de las derivadas entonces podemos meter inmediatamente nuestro problema de derivar a la parte de adentro y simplemente dejar la suma que fuera entonces la primera parte se ve idéntica que antes esto es una constante derivada del logaritmo eso no entre acá para poder calcular esta parcial entonces tenemos que recordar como estaba escrita y lo mismo va a ocurrir de este lado va a salir del signo menos por lo que teníamos acá adentro y tenemos entonces  $1 - \frac{1}{z_j}$  por la parcial de estar acá es lo que estamos sustituyendo acá ya hicimos todas las cuentas todo lo que se cancela entonces no es necesario volverlo a hacer simple y sencillamente llegamos a lo que ya teníamos antes era  $y - \frac{1}{z_j}$  llega a menos acá todo esto a lo que se le había llamado delta acá y lo que vamos a empezar ahora es a sacar un poquito los detalles estábamos antes derivando con respecto a  $z_j$  con respecto a  $j$  y solamente nos había quedado un término aquí en este caso estamos derivando con respecto a  $y$  y jota y quien depende de  $y$  y jota pues son las  $a_s$  es básicamente estos pesos son constantes cuando los consideramos con respecto a este que nos interesa ahora y como ahora pues recordamos que nuestra red está siendo alimentada por todos estos caminos lo que vamos a tener es que por cada uno de estos términos si no se está sobreviviendo el peso que acompaña a esta  $j$  entonces por eso aquí vamos a tener nuestro término delta que a pues está acompañado por el respectivo peso que conecta a la neurona o está con la neurona que sale acá y ahora lo que tenemos que hacer es calcular la parcial de este valor de activación la buena noticia es que efectivamente solamente nos interesa un valor de activación que es el de esta neurona  $j$  que es con la que nos estamos conectando todos los demás en este momento no nos van a afectar pero entonces de aquí están saliendo precisamente estos términos y ahora para poder calcular esta derivada parcial pues vamos a repetir lo mismo de antes que es recordar cómo fue calculada y ahí es donde va a aparecer ahora si explícitamente el peso y jota que es el que nos interesa bien para ir simplificando el bueno de este paso para este paso simple y sencillamente se hace en la sustitución ahora vamos a estar trabajando sobre la  $s_j$  que serían todas estas neuronas de acá que fueron las que contribuyeron con sus distintos pesos a poder calcular este valor este que está aquí entonces sería la combinación lineal esta neurona por eso le pusimos aquí el índice  $j$  y ahora lo que vamos a tratar de hacer es ir simplificando poco a poco esta notación entonces ahora sí vamos a aprovechar que esto se llama delta acá y aquí lo que tenemos es la misma suma pero de esto escribimos delta acá y este es exactamente el mismo peso que estamos colocando acá de aquí bueno regla de la cadena recordemos las propiedades que tienen gente simplemente dejamos indicado que prima con respecto a  $z_j$  sería que prima evaluado en  $z_j$  y la derivada de lo que está acá dentro de lo que esté que dentro ahora sí a diferencia de este caso donde nos quedamos con varios pesos en este solamente hay uno que es el que basa un elemento que va a sobrevivir estamos derivando con respecto a  $y$  y  $j$  entonces ahora las variables son éstas las que están en que son las constantes y solamente hay una es cuando  $y$  vale exactamente acá y cuando calcula la derivada con respecto a sí misma pues va a ser 1 y sobrevive una desde toda esta suma solamente sobrevive precisamente acá es la que tenemos aquí entonces aquí está ahí que sobrevivió a todo este elemento que tenemos aquí de hecho observemos que cómodamente está ahí

## 6. Entrenamiento por retropropagación

no depende del acá realmente si la podemos actualizar a todo esto y le íbamos a llamar del pj y entonces podemos reescribir ahora si nuestras ecuaciones de manera un poco más cómoda ya teníamos entonces de la primera capa que el error cometido por la última capa lo podíamos ver directamente como la diferencia entre lo que queríamos y lo que logramos calcular esto tiene una característica bastante interesante es exactamente por cada neurona tenemos uno de estos después venía la parcial de j con respecto a cada uno de los pesos observemos que por cada uno de estos errores vamos a tener de hecho tantos componentes de éstas como pesos estaban conectados con la carísima neurona entonces de estos tenemos un montón y lo único que teníamos que hacer era multiplicar este valor único por el valor de activación de la neurona con la cual estaba conectado y bueno aquí el signo menos que venimos cargando por la definición de la función de error y ya está ahora equipo se ve con estas tres neuronas bien en la definición de este delta j donde tenemos un producto de las delta casa que venían capa que está más hacia adelante multiplicados por los pesos correspondientes tenemos una suma sobre todos los elementos en última capa y una vez que hicimos esto viene multiplicar por función prima evaluada en la receta j observemos una vez más que está nada más depende de jota no depende de las casas entonces por eso lo podríamos poner entonces realmente le suma nada más afecta este ya que tenemos entonces este producto este que está aquí otra vez tenemos uno por cada neurona en la capa de en medio entonces aquí tenemos las casas aunque vamos a tener las jotas y tenemos uno por cada una de estas este se suele interpretar también como la contribución el error o el error que cometieron todas las neuronas en la capa de en medio y observen que el error de cada neurona pues realmente tiene una es la suma de como contribuyó al error de todas las neuronas que estaban en la capa siguiente pues tiene mucho sentido no se están participando en todos lados pues su error es la suma de todos los errores a los cuales contribuyó y entonces para calcular el gradiente otra vez nos queda una fórmula bastante sencilla es multiplicar este único error por el valor de activación de la neurona con la cual estaba conectada en la capa anterior y de esta manera podemos obtener todas las parciales con respecto a los pesos que estaban conectados con estas neuronas y aquí tenemos todos esos y este que quedaría la parte de hacer el cálculo directamente el siguiente problema que vamos a tener es bueno si podríamos implementar perfectamente ya con un algoritmo de descenso por el gradiente con estas fórmulas que tenemos aquí simplemente tendremos un montón de ciclos fort para estar calculando todas estas sumas y multiplicaciones sin embargo la forma en la que se acostumbra a trabajar ahora con las redes neuronales no es directamente calculando esto componente por componente sino que lo vamos a utilizar connotación matricial esto va a tener varias ventajas por un lado la notación va a ser muchísimo más compacta y por otro lado va a permitir en la implementación de los algoritmos con procesadores como duda lineal con gpu a piece como q da para procesadores con gpu de manera que estas operaciones se están realizando en paralelo y entonces trabajamos muchísimo más rápido con las redes neuronales y realmente el estar utilizando notación matricial va a tener un impacto directo sobre el tiempo que tardan en ejecutarse nuestros algoritmos y aquellos que están preocupados un poco por la ecología nos dirán que si tardaremos poco tiempo pero vamos a estar gastando muchísima electricidad entonces igualmente hay que seguir tratando de optimizar esto lo más que pueda hasta aquí bien en la primera

parte de la derivación del día.

## Vectorización

vamos a ver ahora una parte que es sumamente interesante es como vamos a paralizar nuestro trabajo mediante el uso de notación con matrices y vectores y vamos a aprovechar ahora sí para generalizar nuestras fórmulas en la derivación matemática de ingrediente se hizo considerando que tuviéramos un solo ejemplar de entrenamiento pero realmente vamos a entrenar a nuestra red con varios ejemplares entonces vamos a agregar eso también tenemos otra vez escrita nuestra forma fórmula para el error de entropía cruzada pero ahora sí estamos agregando este término donde estamos calculando un promedio sumando los errores sobre todos los ejemplares de entrenamiento aparecen entonces estos índices extra y en la parte de arriba para indicar en qué ejemplar de entrenamiento nos encontramos y para poder empezar a trabajar con matrices vamos a empezar a ver cómo están escritos nuestros datos todo recordemos como habíamos dicho que íbamos a escribir nuestro nuestras entradas de entrenamiento la idea es que cada ejemplar es un renglón tenemos  $n$  características este  $x_0$  está pensando en que estamos agregando aquí los sesgos y después observemos lo que ocurre con nuestras etiquetas de clasificación recordemos que en redes neuronales utilizamos el one hot en coding entonces aunque si tenemos cinco clases tenemos cinco neuronas de salida y eso quiere decir que nuestra etiqueta tiene algo así este sería una etiqueta en la que la clase correcta es la que está en la tercera neurona otro ejemplar podría ser de esta manera y pues aunque solamente haya uno que sea distinto de cero la forma en la que van a venir empaquetadas nuestras respuestas correctas va a ser precisamente en forma de matriz donde tenemos sobre los renglones tantos bits como neuronas haya en la última capa y tenemos hacia abajo los  $m$  ejemplares de entrenamiento después recordemos cómo teníamos escritas las matrices de pesos teníamos por cada columna los pesos que contribuyen a la a una neurona en la siguiente capa entonces sobre cada renglón tendríamos a tener tantos pesos como neuronas haya en la capa cl y hacia abajo vamos a tener tantos pesos como neuronas había en la capa  $s_l - 1$  y que contribuyeron al siguiente elemento en ese  $l$  así que tendríamos a las diferentes neuronas que van a estar en nuestra nueva capa eso es lo que va a hacer es que cuando multipliquemos  $x$  y  $z$  nos queden otra vez los valores de cada en la siguiente capa con los renglones correspondientes a los ejemplares de entrenamiento y horizontalmente los valores de activación de cada neurona en la siguiente bien entonces partes importantes ejemplares de entrenamiento hacia abajo para estos dos aquí el décima capa hacia acá está para  $l$  menos uno ahora están copiadas acá del lado izquierdo las fórmulas tal y como las obtuvimos ejemplar por ejemplar y lo que vamos a hacer ahora es escribirlas en forma matricial bueno que ya me adelanté ya se las escribí entonces va a ser más fácil si ahora explicó exactamente qué es lo que estamos viendo aquí para eso voy a utilizar ahora otro programa y aquí está vamos a utilizar aquí para que pueda dibujar qué es lo que está primero que teníamos en delta está en la capa no en la capa  $l$  bueno voy a dibujar aquí ahorita una pequeña red neuronal en este conector y entonces nos estamos preguntando

## 6. Entrenamiento por retropropagación

por qué pasa con él aquí está él bueno lo que decíamos era que vamos a tener una de estas del estás por cada uno de nuestras neuronas de salida y lo que vamos a ver es qué tenemos los ejemplares hacia abajo y las neuronas horizontales los valores que tengo aquí verticalmente los acosté y es lo que tengo acá ya vimos también entonces que la llega realmente tiene varios bits horizontalmente y la sas también entonces si yo restó enrique menos acá observen que esto en realidad es un vector tengo una componente por cada neurona de salida entonces aunque aquí les saquemos componente por componente pues tengo uno de cada uno entonces de todo esto me va a salir todo un vector y si además considero que hay ejemplares de entrenamiento pues voy a tener  $m$  de esos entonces la forma más cómoda para acomodar eso va a ser la matriz delta  $\Delta$  donde voy a restar dos matrices voy a restar la matriz que tiene esta forma menos la matriz  $\Delta$  que de hecho tiene exactamente la misma forma la diferencia es que como ésta está evaluada con una logística pues los valores no son ceros y unos en realidad son valores entre 0 y 1 entonces cuando haga la resta pues voy a obtener otra vez una matriz con la misma forma pero pues con números entre 0 y 1 entonces delta  $\Delta$  va a tener los componentes hacia la derecha y  $m$  ejemplares dentro y entonces ya podemos ver porque simplemente basta con que yo reste menos en forma matricial y ya automáticamente tengo acomodados todas las del test que aquí tendría que haber calculado una por una y como estamos repitiendo eso para todos los ejemplares de entrenamiento pues ya tengo de una buena vez todos acomodados en la misma matriz y vamos a hacer ahora para el gradiente este es un efecto mucho mucho muy interesante porque en primer lugar aquí este nos conviene el hecho de que tenemos varios ejemplares de entrenamiento y nosotros lo que queremos hacer es sacar un promedio sobre todos ellos entonces aparte de que afectan los signos menos entonces el signo menos se conserva el 1 entre  $m$  se conserva porque es para poder sacar promedio pero ahora vamos a tratar de matar varios pájaros de una pedrada observemos lo siguiente tenemos que multiplicar y sumar a todas las valores de activación por el delta observen que este es en la capa anterior este es el error que cometió en la capacidad ente el hacer las combinaciones de todos contra todos es cuando obtengo el efecto que ocurrió con cada peso pero quiero sumar estos productos para todos los ejemplares de entrenamiento y sumarlos bueno eso es prácticamente lo que hace una multiplicación de matrices ahora para poder hacer todo esto en un solo paso lo que tenemos que hacer es acomodar las acorde mente nos vamos a hacer lo siguiente ya habíamos dicho que esta delta que realmente tiene forma de matriz como ésta de tal manera que mira aquí hacia acá tenemos todos los es el es y así aquí abajo tenemos todos nuestros ejemplares de entrenamiento bueno y que sabemos de la  $a_j$  y también son los valores de activación y queremos combinar cada  $j$  de cada ejemplar de entrenamiento con su respectiva que del mismo ejemplar de entrenamiento y después multiplicarlos y sumarlos bien entonces vamos a acomodar a la sas de la siguiente manera y vamos a poner hacia acá a los ejemplares de entrenamiento y así acá a los elementos son ordenadas iba a suceder si yo hago esto observen que se multiplican ejemplares de entrenamiento contra ejemplares de entrenamiento misma  $a_j$  contra mismos del  $t_{jk}$  pero de diferentes ejemplares de entrenamiento cada uno con su respectivo y además se suman para definir la coordenada que va a quedar aquí cuando yo tomé este contra la siguiente columna entonces otra vez vamos a estar multiplicando en ejemplares contra mi ejemplares los vamos a sumar y nos

va a dar el dato que tengamos aquí los conservamos las columnas que teníamos acá si cuando empieza a hacerlo con los renglones voy a tener entonces lo que ocurre con las s menos el c lm no son tan clones que teníamos acá y esto es sumamente interesante porque esto que está aquí debería de resultarnos un poco conocido si si - solo vamos a ver que teníamos acá si si - 1 tiene exactamente la misma forma en la matriz de pesos y no es coincidencia recordemos que es lo que estamos tratando de calcular estamos tratando de calcular el gradiente el ingrediente es la parcial con respecto a cada uno de los pesos entonces lo que acabamos de obtener es ese gradiente que tiene acomodadas cada una de las parciales en la posición que corresponde al mismo peso pero en la matriz de pesos bastante bonito muy bien entonces por eso cuando ponemos la matriz a él pero transpuesta para que tengamos ahora si los ejemplares de entrenamiento de forma horizontal entonces podemos obtener en una multiplicación de matriz todos los componentes del gradiente para los pesos en esa capa lo que ocurre con las siguientes bueno de hecho con los siguientes gradientes es que se van a hacer exactamente igual lo único que nos queda ligeramente entretenido es cómo vamos a calcular los errores de las capas intermedias hay entonces que observar lo siguiente en primer lugar las ventas van a tener la misma forma que éste es que primas entonces simple y sencillamente van a ser dos matrices donde se tienen que multiplicar componentes a componentes las vamos a poner con este símbolo que significa circo digo se llama cirque y significa multiplicación componente a componente entre estas dos matrices ahora como vamos a escribir ésta se parece a lo anterior pero ahora sobre lo que queremos es sumar es sobre las neuronas de salida recordemos que la contribución al error de una neurona de la capa intermedia pues es la suma sobre los errores a los que contribuyó en todas las neuronas en la capa siguiente entonces del ejercicio anterior ya debemos de ver como pista en que si queremos hacer una suma sobre estos elementos pues como que esto es lo que vamos a tener que tener en la parte horizontal de la primera matriz y vertical de la segunda entonces de aquí empiezan a salir precisamente los tips de cómo escribir estas 2 estás deltas pues ya tienen automáticamente de manera horizontal lo que está ocurriendo con cada una de las neuronas entonces se quedan exactamente iguales y lo que estaba ocurriendo con los pesos es que entonces vamos a crear los componentes de ese l sobre los renglones los necesitamos tomar la transpuesta para que cumpla con esa condición y en principio y básicamente ya quedaría que ese entonces está en notación extraña que tenemos en la parte de abajo recordemos que si estamos utilizando sesgos si estamos utilizando sesgos tenemos una neurona en esta capa que no está conectada en nadie con nadie en la parte de atrás entonces no existen pesos que hayan hecho que las neuronas de acá contribuyan al error de esta ésta no tenía error es un sesgo entonces aquí no hay nada que calcular por eso lo que tenemos que hacer en la parte de aquí es quitar a los elementos que corresponderían a las así pues el hecho de que los riesgos no están conectados con la capa anterior es básicamente lo que nos está diciendo es que quitemos de la matriz de pesos a todas las conexiones que parten de esta neurona porque estas conexiones pues no bueno ésta no va a estar contribuyendo cada día más bien las que estén en esta capa no están contribuyendo al error de esta neurona entonces lo que ocurría con esta neurona en este momento a estas no les importa realmente porque no están pasando para acá entonces por eso se quita estos elementos que estén aquí nos

## 6. Entrenamiento por retropropagación

dice entonces que nos brincamos algo que está en el índice cero y nos vayamos nada más a los que están en el índice 1 entonces teniendo este cuidado de ir completando los datos sin los elementos de sesgo cuando no hay esos elementos presentes pues ya podemos implementar todo esto con operaciones ahora para poder ponerle un poco en limpio observamos que simplemente el que es diferente es el error cometido en la última capa nada más es yemen osa 1 después las fórmulas no sólo para la capa 1 menos uno sino que si tuviéramos más capas hacia atrás se verían idénticas entonces les podemos escribir así tenemos que el error va a ser esta delta 1 multiplicada por la matriz de pesos quitando las los que afectan en los sesgos x eje prima pero recordemos que ese prima pues tenía truco porque son los valores de activación que ya tenemos multiplicación componente a componente por 1 - estos valores de activación y finalmente para los componentes del gradiente ahora sí lo único que vamos a hacer era multiplicar los valores de activación de la capa anterior por los errores de la capa siguiente y eso nos iba a dar ahora sí la participación del gradiente para cada uno de los pesos y esto lo podemos repetir obteniendo uno de una matriz de estos por cada matriz de pesos que hayamos utilizado los con esta técnica básicamente obtenemos el gradiente de pequeños bloques hitos de material de matrices una matriz por cada matriz de pesos si quieren escribir el gradiente entonces como se escribe usualmente en cálculo es como un solo vector pues lo que hay que hacer es aplanar todas esas pequeñas matrices y poner cada una de las componentes en una componente del vector y nos va a quedar un vector gigantesco que tiene tantas tantas componentes como la suma de componentes en todas estas matrices bien y con eso quedaría todo el último detalle que les quiero mencionar es bueno este signo siempre se nos hace travesuras en esta versión recuerden que lo que acabo de calcular aquí es el gradiente y cuando nosotros hallamos descenso por el gradiente lo que necesitamos entonces es la dirección inversa entonces nada más finalmente entonces podemos decir en qué consiste el algoritmo de propagación hacia atrás lo único que estamos haciendo es obtener los ejemplares de entrenamiento las respuestas correctas necesitamos saber cuántas que pasan vamos a inicializar nuestras matrices de errores con ceros bueno podremos hacer lo mismo también con las del gradiente empezamos en la [Música] haciendo feedforward que es lo que nos indican estas primeras líneas asignando nuestros datos de entrenamiento a la primera capa la capa de chocolate a partir de ahí utilizamos feedforward para ir calculando todos los valores de activación de las siguientes capas hasta llegar a la última y a partir de ese momento podemos empezar a trabajar ahora sí con el entrenamiento de la red tomar en cuenta que era lo que queríamos para sustituir en los cálculos de los errores y los gradientes con las fórmulas que tenemos en la parte de atrás y nos importó precisamente las etiquetas y ya que tenemos ese entonces podemos actualizar ahora si nuestros pesos en paralelo tienen que actualizarse todos al mismo tiempo no para actualizar primero unos y después otros porque eso ya no es el gradiente en paralelo tenemos que actualizar tomando los pesos que ya teníamos menos alfa por el gradiente d función de error repitiendo esto el suficiente número de veces la idea es que eventualmente estos pesos nos permitan encontrar un mínimo de la función de error y eso sería todo en esta fórmula de hecho fue obtenida de andrew james en la fórmula para el error roselyn orvin presentan la derivación con diferencias del cuadrado en lo que se hizo aquí fue combinar ambas para obtener entonces la derivación pero ya

con entropía cruzada y si quieren ver un poco más acerca de esta técnica pues pueden checar y sobre todo la parte de la interpreta.

# 7 | Optimización del entrenamiento

## Problemas en redes profundas

Vacio.

## Gradiente desvaneciente (o que explota)

Vacio.

## Entrenamiento en línea vs en lotes

Vacio.

## Normalización y normalización por lotes

Vacio.

## Regularización

Vacio.

# 8 | Caso de análisis e interpretación

## Red Hinton árbol familiar con numpy (entrenamiento)

Vacio.

## Red Hinton árbol familiar con pytorch

Vacio.

## 9 | Entrenamiento con genéticos

### Algoritmos genéticos

Vacio.

### Neuroevolución

Vacio.

### Antecedentes: Aprendizaje por refuerzo en videojuegos

Vacio.

### Arquitectura para estimar la función de recompensa

Vacio.

### Entrenamiento

Vacio.

# 10 | Mapeos autoorganizados

## Introducción

Vacio.

## Aprendizaje no supervisado

Vacio.

## Mapeos auto-organizados

Vacio.

## Kohonen

Vacio.

# 11 | Redes Neuronales Convolucionales

**Convolución**

**Redes Convolucionales**

**Softmax**

**MNIST**

## **Parte III**

## **Redes con ciclos**

# 12 | Redes Neuronales Recurrentes

**Derivadas ordenadas**

**Retropropagación en el tiempo**

**Sistemas dinámicos y despliegue del grafo**

**Arquitectura recurrente universal**

**Función de error**

**Forzamiento del profesor**

## 13 | Atención

# 14 | LSTM

# 15 | GRU

## **16 | Casos de análisis: etiquetado de palabras y conjugación de verbos**

16. CASOS DE  
ANÁLISIS: ETIQUETADO  
DE PALABRAS Y  
CONJUGACIÓN DE  
VERBOS

## **Parte IV**

# **Redes no dirigidas**

# 17 | Redes de hopfield

## Entrenamiento

# 18 | Máquinas de Boltzman

## Entrenamiento

Partículas y partículas de fantasía

Máquinas de Boltzman Restringidas

# 19 | Redes adversarias

## GANs

# A | Ecuaciones diferenciales